

# TEMPORAL ANALYSIS AND FUTURE PREDICTION OF BILLION TREE TSUNAMI FORESTS: A CASE STUDY OF GARHI-CHANDAN PAKISTAN

MATEEN, S. – NUTHAMMACHOT, N.\* – TECHATO, K.

*Division of Environment, Faculty of Environmental Management, Prince of Songkla University, Songkhla 90110, Thailand*

*\*Corresponding author  
e-mail: narissara.n@psu.ac.th*

(Received 23<sup>rd</sup> Nov 2023; accepted 9<sup>th</sup> Feb 2024)

**Abstract.** This article investigates the temporal analysis of billion tree tsunami forests in Garhi Chandan area of Pakistan based on three supervised methods, namely random forest algorithm (RFA), principal component analysis (PCA) combined RFA and support vector machine (SVM). As a first step, the Sentinel-2 and Landsat-8 data fusion is performed to enhance the spatial resolution of the data to 10 m. The overlapping features in the data may compromise the classification accuracy, thus, to overcome this limitation, PCA is utilized. As a second step, classification is performed using RFA, PCA-RFA and SVM methods by using the data of the years 2016 and 2023. The change map analysis is done by using the aforementioned methods. As a next step, ground data matching is performed for the classified samples using each method. Finally, by utilizing logistic regression, future prediction for the years 2028, 2030 and 2033 is performed. The PCA-RFA technique achieved the best overall accuracy of 95% with a Kappa hat score of 0.93. The second-best result is achieved by RFA classifier, with overall accuracy of 92% with a Kappa hat score of 0.92. SVM showed moderate matching with an overall accuracy of 72% with a Kappa hat score of 0.55.

**Keywords:** *temporal analysis, random forest classifier; billion tree tsunami project, principal component analysis, support vector machine, image classification*

## Introduction

Machine learning algorithms are being widely utilized in remote sensing applications such as land use and land cover mapping (LULC). Amongst several machine learning classifiers, the random forest algorithm (RFA), deep learning and support vector machine (SVM) have attracted the attention of the research community and its utilization in remote sensing applications is growing day by day (Sheykhmousa et al., 2020). With the advancement of satellite missions, the classification is becoming more difficult due to large volume of satellite data, training data set imbalance and landscape non-uniformity (Pouteau, 2012).

To date the three most commonly utilized machine learning algorithms in remote sensing applications are SVM, RFA and deep learning methods (Habib et al., 2009; Boulesteix et al., 2012; Mainali et al., 2023). Although deep learning methods are more effective as compared to SVM and RFA (Heydari and Mountrakis, 2018, 2019; Mohammadimanes et al., 2019), however it is difficult to optimally select and train the hidden layers. Due to ease of implementation and low computational cost, SVM and RFA methods are still attracting the attention of the research community. Support vector machine (SVM) was first introduced in 1970s and it found wide applications in remote sensing (Mountrakis et al., 2011). SVM is robust against the distribution of data (Mather and Tso, 2011). Since the classes are binary so a line acts as a separation boundary between discrete classes. The portion of training

inputs that closely matches the feature space acts as support vector for SVM (Bazia and Melgani, 2006). Practically, in order to separate the classes from each other, using linear SVM, a hurdle is data overlapping, and thus the basic linear SVM do not perform optimally in such situations (Scholkopf and Smola, 2018). In order to address the aforementioned problem, kernel trick method is proposed in the basic SVM (Kavzoglu and Colkesen, 2009). The kernel function is implemented using radial basis and sigmoid functions (Mountrakis et al., 2011; Khemchandani et al., 2008).

Random forest algorithm (RFA) is a supervised ensemble method that was first introduced by Breiman in 2001 (Breiman, 2021). Amongst multiple models/trees, a majority voting method is utilized to choose the best output. RFA uses bagging methods (Fawagreh et al., 2014; Breiman, 1996). The bagging is more robust to the over fitting problem as compared to boosting. RFA utilizes bagging method and it has lower computational complexity as compared to other known machine learning methods. Kandpal and Kumar (2022) utilized the RFA for detecting medical plants in the Himalayan region. In order to solve the imbalanced data set classification problem, a filter based RFA is proposed by Khosravi reported in Khosravi (2019). RFA is also applied in hazard analysis such as landslide mapping using multisource data and the details are given by Liu et al. (2019). Land use and land cover area estimation is vital for management of resources, urban and town planning; in this regard RFA also finds application in such classification problem reported in Sales (2022). Miao et al. (2010) presented a detailed comparison between RFA and adaboost tree method for an ecosystem classification in Mojave Desert. Liyanage et al. (2019) utilized RFA to predict the landslide hazards in Sri Lanka.

While classifying the multi spectral data (MSD), usually the collinear dependencies of the classes will not allow the classifier to operate in cost effective way (Uddin et al., 2020). Principal Component Analysis (PCA) is a widely utilized method applied in extracting the most important features and in the dimension reduction of MSD data. Singh and Harrison (1985), Mackiewicz and Ratajczak (1993), and Eklundh and Singh (1993) utilized PCA for reducing the overlapping features in the data. Based on the above cited work, this research article has the following objectives:

1. Three supervised classification methods, namely SVM, RFA and PCA-RFA are utilized for the classification of the Sentinel-2 and Landsat-8 OLI data for the years 2016-2023.
2. PCA is utilized for extracting the most important features from the Sentinel-2 and Landsat-8 OLI fused data set and minimization of the dimensionality and correlation of the data sets.
3. After applying PCA to the original data set, RFA is applied to the resultant data set and the obtained results are compared with RFA and SVM methods.
4. Change map analysis is performed.
5. The classified data is compared with the ground samples collected from Google earth.
6. Future prediction for the years 2028, 2030 and 2033 is performed.

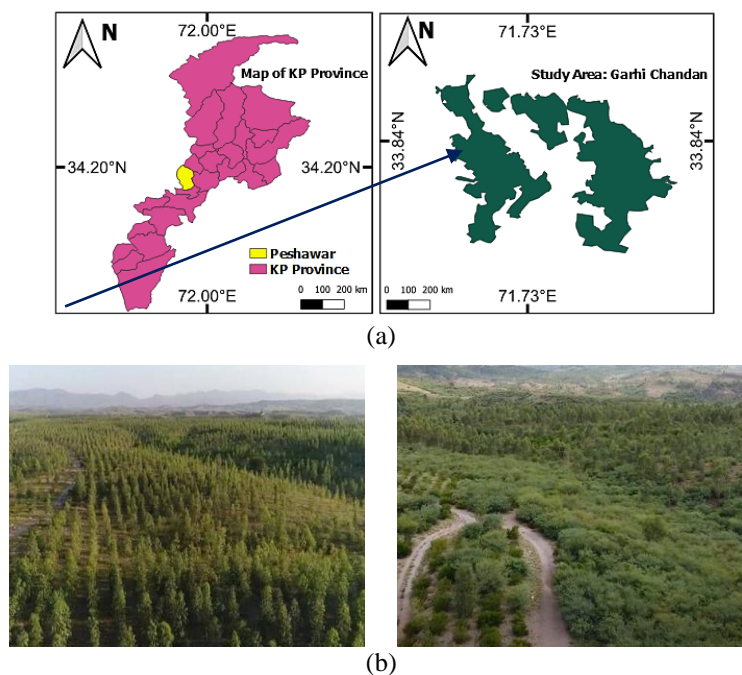
The rest of the article is presented as follows. The next section presents materials and methods, followed by the results; discussion on the results sections, while a conclusion is drawn based on the presented results in the last section.

## Materials and methods

This section presents details about our proposed study area, classification methods, estimation of accuracies and details about reference data for the classification.

### Study area

In 2014, Government of Pakistan initiated billion tree tsunami forests plantation (afforestation) in several regions of Khyber Pakhtunkhwa province. For this research work, we selected Ghari Chandan forests region as our main study area. Our study area is approximately 3141.6 ha/31.42 million m<sup>2</sup>, located 33°50'0" North and 71°42'0" East and in the vicinity of the capital of Khyber Pakhtunkhwa province. *Figure 1a* shows the map of Khyber Pakhtunkhwa province (left figure), while our study area is shown in the right column of *Figure 1a*. *Figure 1b* shows the experimental images of the study area.



**Figure 1.** Study area (a) map (b) experimental images

### Classification methods

In this research work, RFA and SVM methods are utilized for classification of land cover in our study area. In the subsections given below, the two classification methods are explained with more details.

#### Random forest algorithm

RFA classifier utilizes ensemble learning, and multiple models/trees are integrated together to enhance the classification accuracy. Two widely utilized ensemble learning methods are boosting and bagging (Fawagreh et al., 2014; Breiman, 1996). Boosting ensemble method is associated with a problem of over fitting. Bagging is another ensemble learning method, that is more robust and the over fitting problem is fixed.

Figure 2 shows the bagging ensemble learning which serves as a foundation of the RFA algorithm. The steps describing the implementation procedure of the RFA are explained in Figure 2.

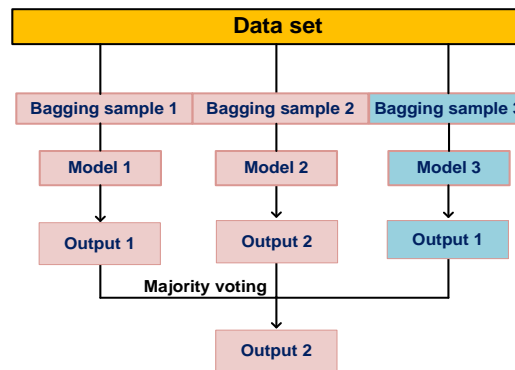


Figure 2. Block representation of random forest algorithm

### Support vector machine

SVM is a machine learning method that is primarily based on the kernel based algorithms. SVM uses convex quadratic optimization to achieve global optimality and it is a non-parametric technique which is robust to the distribution underlying data sets. The linear SVM classifier is the basic version, which introduces a linear boundary between classes optimally. As shown in Figure 3, the expression  $\vec{A} \cdot \vec{x} + B = 0$ , denotes the separation linear boundary between the two classes, and one green and two black dots show the support vectors. As shown in the figure,  $\vec{A}$  is a normal vector and  $B$  is the scalar offset. From Figure 3, the objective function  $sign(\vec{A} \cdot \vec{x} + B)$  is optimally solved using the objective function (1) detailed in (Liu et al., 2019).

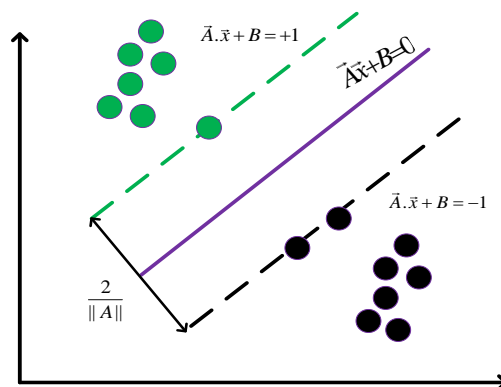
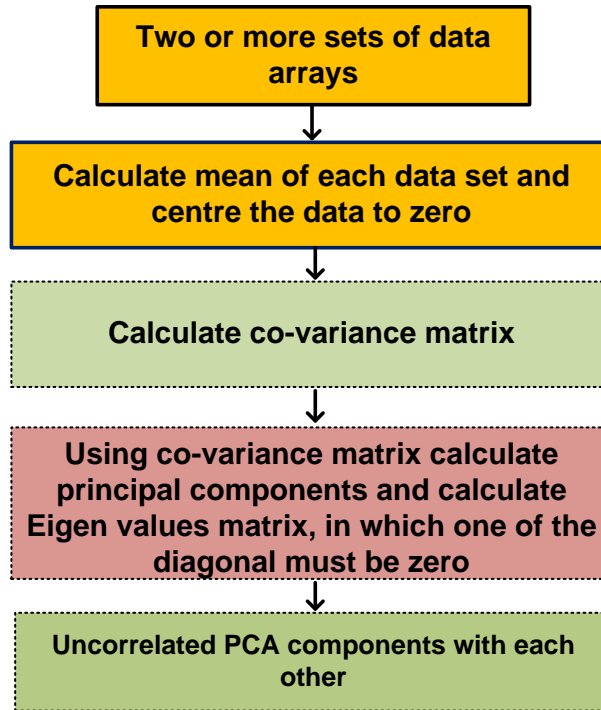


Figure 3. Block representation of SVM

### Principal component analysis (PCA)

PCA is a mathematical data dimensionality reduction tool used to map large number of correlated variables into small number of uncorrelated variables. The uncorrelated variables are called the principal components (PCs) (Uddin et al., 2020). The block diagram representation of PCA is shown in Figure 4.



**Figure 4.** Block representation of PCA

Let us assume two data sets (Uddin et al., 2020) as follows:  $X = [x_1 \ x_2 \ x_3 \ \dots \ x_m]^T$  and  $Y = [y_1 \ y_2 \ y_3 \ \dots \ y_k]^T$ , where  $X$  is  $n * m$  order and  $Y$  is  $p * k$  order data sets. In order to understand the PCA, consider the following steps.

**Step 1:** In the first step, the two data sets given in  $X$  and  $Y$  are centered to the origin and expressed as follows:

$$\begin{aligned} X_c &= [x_i - x_{mean}]^T ; i = 1:m \\ Y_c &= [y_j - y_{mean}]^T ; j = 1:k \end{aligned} \tag{Eq.1}$$

In *Equation 1*,  $X_c$  and  $Y_c$  represent the data sets shifted to the origin of two-dimensional plan,  $x_i$  and  $y_j$  show the elements of the original data sets, while  $x_{mean}$  and  $y_{mean}$  depicts the mean values of the original data sets.

**Step 2:** In the second step, the covariance matrix  $\text{cov} = \begin{bmatrix} \text{COV}_{11} & \text{COV}_{12} \\ \text{COV}_{21} & \text{COV}_{22} \end{bmatrix}$  is calculated based on  $X_c$  and  $Y_c$  using the following expressions:

$$\begin{aligned} \text{cov}_{11} &= \frac{1}{m-1} \sum_{i=1}^m (x_{ci} - x_{c(mean)})^2 \\ \text{cov}_{22} &= \frac{1}{k-1} \sum_{j=1}^k (y_{cj} - y_{c(mean)})^2 \\ \text{cov}_{21} = \text{cov}_{12} &= \frac{1}{m-1} \sum_{i=1}^m (x_{ci} - x_{c(mean)}) \cdot (y_{ci} - y_{c(mean)}) ; i = j \end{aligned} \tag{Eq.2}$$

In Equation 2,  $\text{COV}_{11}$ ,  $\text{COV}_{12}$ ,  $\text{COV}_{21}$  and  $\text{COV}_{22}$  show the elements of covariance matrix,  $x_{ci}$  and  $y_{cj}$  show the elements of the origin shifted data sets, while  $x_{c(\text{mean})}$  and  $y_{c(\text{mean})}$  depicts the mean values of the origin shifted data sets.

**Step 3:** Eigenvalues and eigenvectors are calculated from the covariance matrix as follows:

$$\det |(\text{cov} - \lambda I)| = 0 \rightarrow \begin{bmatrix} \text{COV}_{11} - \lambda & \text{COV}_{12} \\ \text{COV}_{21} & \text{COV}_{22} - \lambda \end{bmatrix} = 0 \quad (\text{Eq.3})$$

$$\text{COV} \cdot v = \lambda \cdot v \quad (\text{Eq.4})$$

In Equations 3 and 4,  $\lambda$  represents the eigenvalues and,  $v$  shows an eigenvector.

**Step 4:** A new matrix  $D$  is formulated as followed:  $D = [X_c \ Y_c]^T$ , then the principal component matrix is calculated as follows:

$$D \cdot v = \begin{bmatrix} x_{1(pc)} & y_{1(pc)} \\ x_{2(pc)} & y_{2(pc)} \\ x_{3(pc)} & y_{3(pc)} \\ \cdot & \cdot \\ \cdot & \cdot \\ x_{m(pc)} & y_{k(pc)} \end{bmatrix} \quad (\text{Eq.5})$$

In Equation 5,  $X_{pc} = [x_{1(pc)} \ x_{2(pc)} \ x_{3(pc)} \ \dots \ x_{m(pc)}]^T$  and  $Y_{pc} = [y_{1(pc)} \ y_{2(pc)} \ y_{3(pc)} \ \dots \ y_{k(pc)}]^T$  represents the elements of principal component matrix. Note that the covariance matrix of the principal component matrix is represented as follows (Uddin et al., 2020; Mackiewicz and Ratajczak, 1993):

$$\text{COV}_{(pc)} = \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix}$$

### Sample design for accuracy assessment

Let  $a_i$  represents the percent area proportion of each class,  $\delta_o$  is the target standard deviation, and  $\delta_i$  shows standard deviation of each class, where  $i$  represents the classes such as forest, bareland and vegetation. To design the total number of samples for accuracy assessment, the following expression is utilized (Olofsson et al., 2014; Congedo, 2021).

$$N = \left( \sum_{i=1}^n a_i \frac{\delta_i}{\delta_o} \right)^2 \quad (\text{Eq.6})$$

In *Equation 6*,  $N$  shows the total number of samples calculated and the target standard deviation is chosen as:  $\delta_o = 0.01$  (Mateen et al., 2022).

The number of samples for each class is calculated based on the average of the following: 1. Equal number of samples ( $Nn^{-1}$ ) 2. Area proportion-based samples ( $Na_i$ ). Thus, the sample calculation of each class is expressed as follows (Olofsson et al., 2014; Congedo, 2021).

$$N_i = \frac{(Na_i + Nn^{-1})}{2} \quad (\text{Eq.7})$$

In *Equation 7*,  $N_i$  shows the number of samples of each class and  $n$  shows the number of classes.

### **Satellite data fusion and band sharpening**

In this work, data set for our study area is downloaded from two satellite missions namely Sentinel-2 and Landsat-8 OLI. The satellite data is downloaded using Semi-automatic plugin (SCP) of the QGIS software. The satellite data is collected on February 15, 2023. The fusion of the two data sets is done using the technique proposed in (Sigurdsson et al., 2022). The proposed method is applied to the data fusion with additional MATLAB-QGIS interface (Mateen et al., 2022). In order to enhance the spatial resolution of all sentinel-2 and Landsat-8 OLI bands to 10 m, the following cost function is utilized (Mateen et al., 2022).

$$F = \sum_{i=1}^n 0.5 \|\lambda_i - D_i B_i k_i\|^2 + \sum_{j=1}^m \xi_j \mu_w(y_j) \quad (\text{Eq.8})$$

In *Equation 8*,  $F$  shows the cost function to be optimized,  $n$  depicts the number of bands,  $\lambda_i$  contains all original bands (Sentinel-2 and Landsat-8 OLI),  $k_i$  shows the estimated bands,  $D_i$  represents the down sampling factor and  $B_i$  is the blurring factor. Moreover  $K = [k_i^T]_{i=1}^n$  is an  $n*k$  orthonormal matrix. The tuned parameter  $\xi_j$  adjusts the spatial regularization term  $\mu_w()$ . Since the data set is downloaded using the semi-automatic classification (SCP) plugin of QGIS and the data fusion and sharpening is performed in MATLAB software so a linking methodology is developed (Mateen et al., 2022). Further details are available in Mateen et al. (2022). The detailed methodology of the data fusion is shown in *Figure 5*.

### **Proposed methodology**

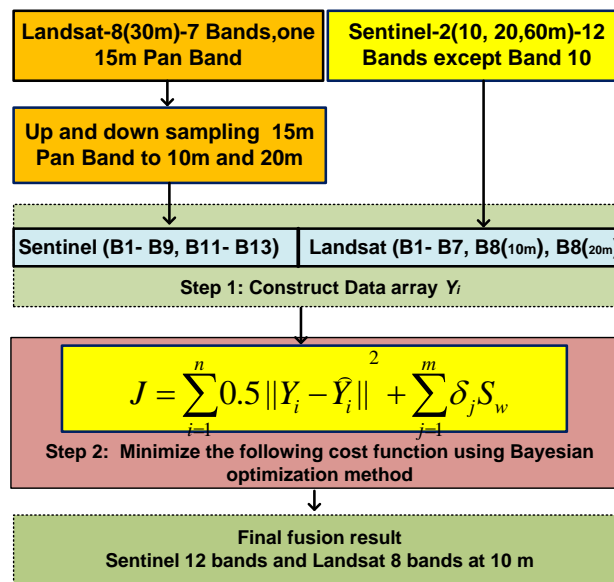
The proposed combined methodology utilized for this research work is shown in *Figure 6*. The proposed methodology consists of three steps. The first step involves the image fusion and sharpening of Sentinel-2 and Landsat-8 bands to a spatial resolution of 10 m. The second step involves the image classification, accuracy assessment and round data validation. In the 3<sup>rd</sup> step, change map is generated and the cellular automata model is utilized to predict the future maps of our study area for the years 2028, 2030 and 2033.

## Results

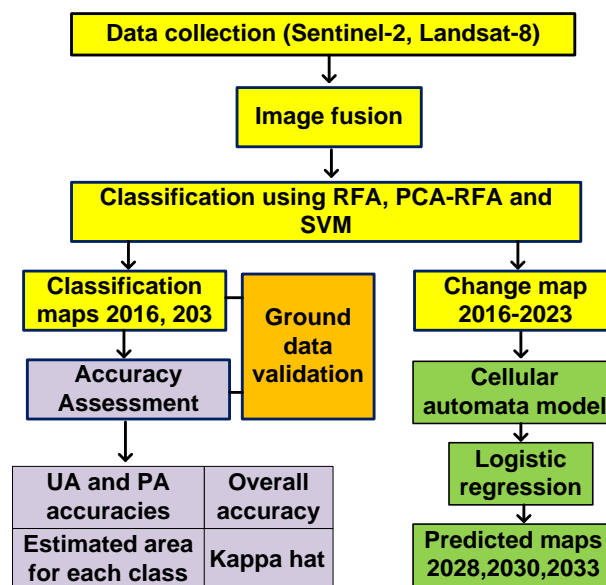
In the results section, the following topics are included.

### Image sharpening

From *Figure 7a* and *b*, the images in left column represent the original bands 1 and 9 of Sentinel-2 satellite with a spatial resolution of 60 m, while in the right column, the sharpened bands 1 and 9 are shown with a spatial resolution of 10 m. Similarly, *Figure 7c* shows the original and sharpened bands of Landsat-8 OLI in pseudocolor green colors combination.

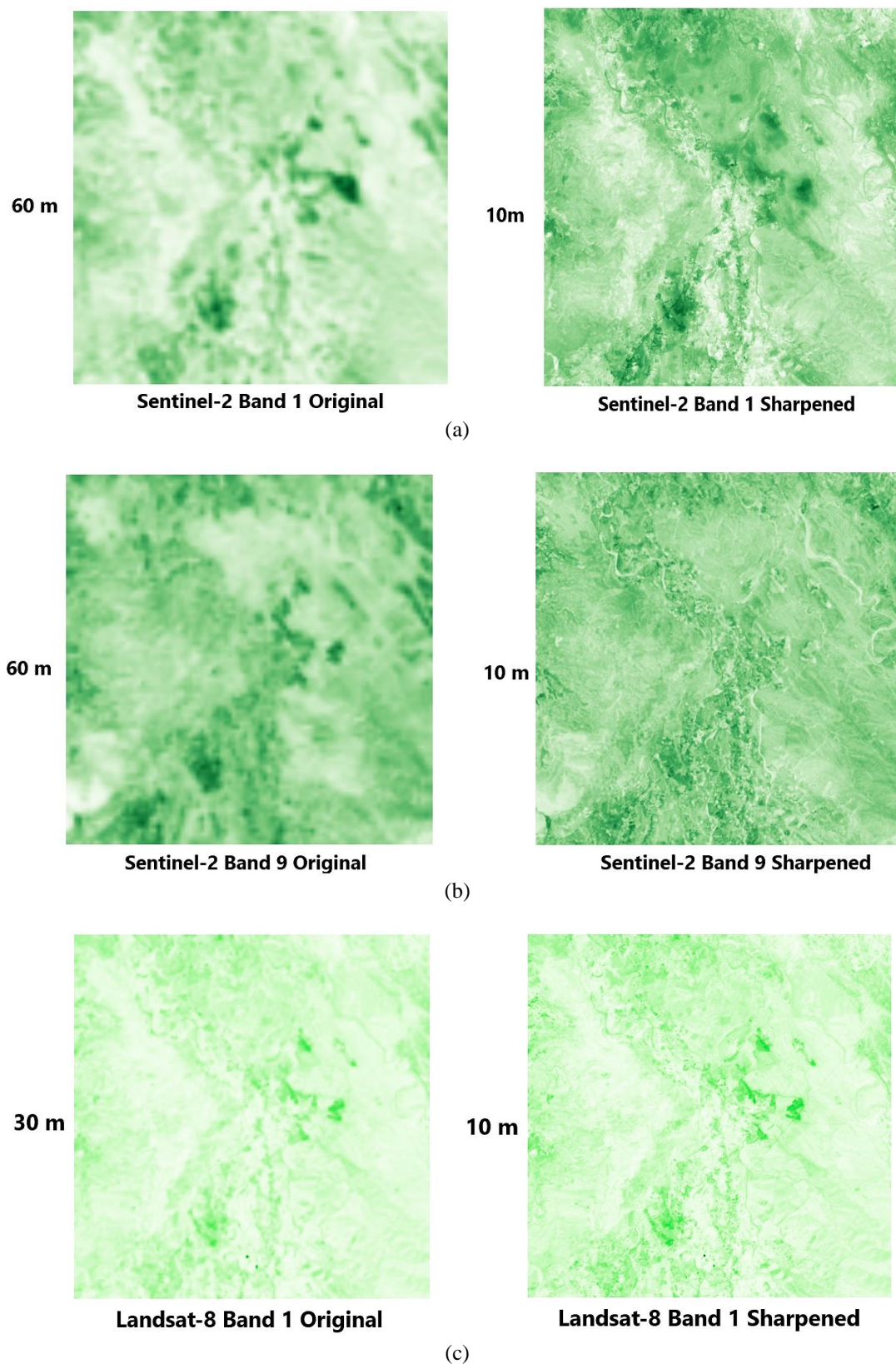


*Figure 5. Block representation fusion process*



*Figure 6. Proposed methodology*





**Figure 7.** Image sharpening (a) Left: Sentinel-2 band 1 original, Right: Sentinel-2 band-1 sharpened image. (b) Left: Sentinel-2 band 9 original, Right: Sentinel-2 band 9 sharpened image. (c) Left: Landsat-8 band 1 original, Right: Landsat-8 band 1 sharpened image (Mateen et al., 2022)

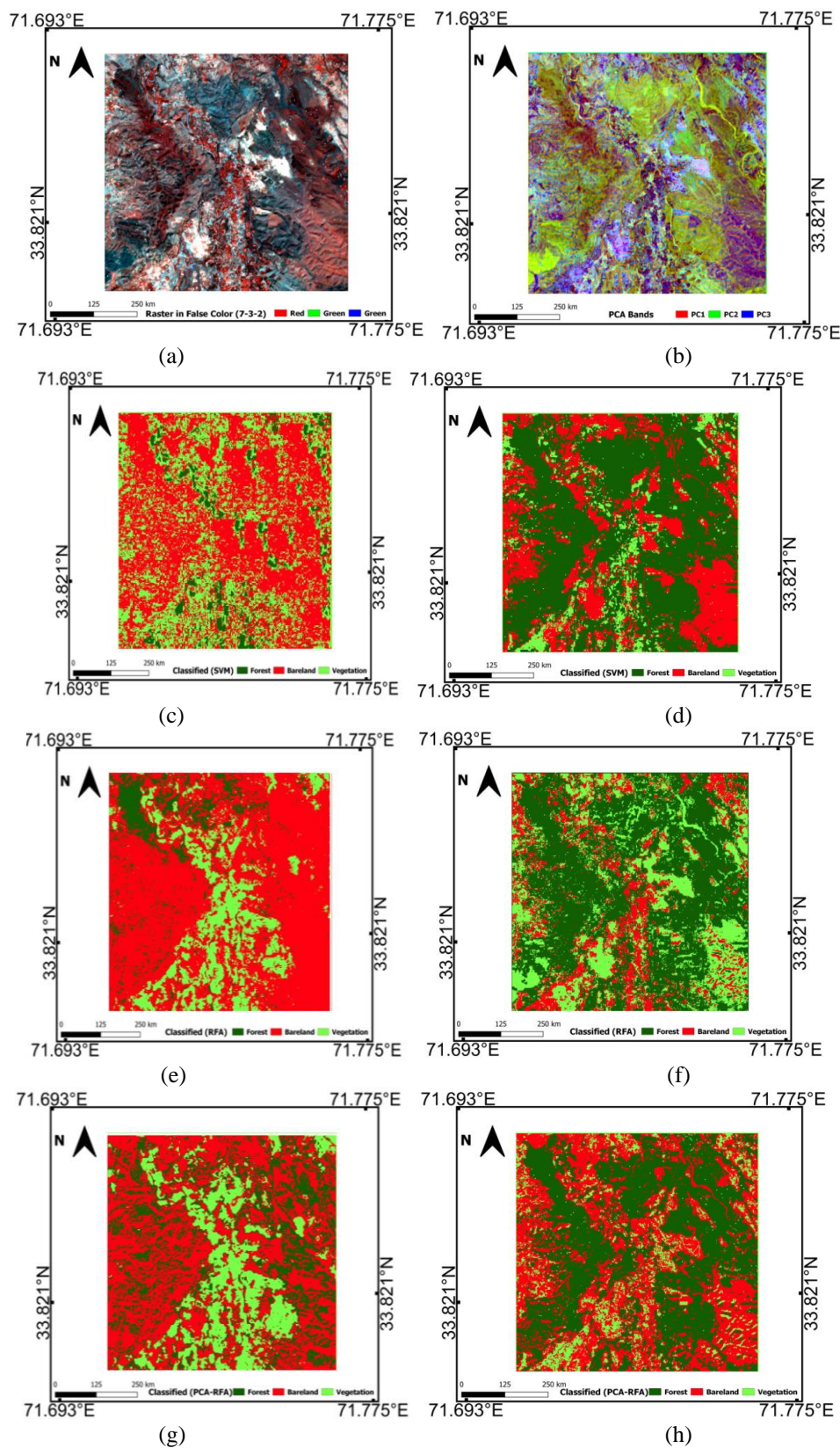
For 20 m bands, the ERGAS, SAM and RMSE scores are estimated as 7.1, 7.05 and 0.001 respectively, while for 30 m bands the respective observed scores are 1.23, 1.23 and 0 (Mateen et al., 2022). Similarly, for 60 m bands, the ERGAS, SAM and RMSE scores are estimated as 0.47, 0 and 0. Note that a smaller ERGAS value shows minimum distortion in the fused image. The lowest ERGAS score of 0.47 is estimated in case of 60 m bands, which depicts low distortion. The average observed UIQI scores for 20 m, 30 m and 60 m bands are recorded as 0.52, 0.65 and 0.70. A UIQI score of 1 interprets a good image quality. In case of 60 m band, a UIQI score of 0.70 means the best quality of the sharpened image as compared the original data. The average observed SAM scores for 20 m, 30 m and 60 m bands are recorded as 7.05, 1.23 and 0. A SAM score 0 interprets a good spectral angle quality. In case of 60 m band, a SAM score of 0 values means the best spectral angle quality of the sharpened image as compared the original data.

### ***Image classification***

For image classification, a total number of 70 training samples are utilized for each SVM, RFA and PCA-RFA methods. For our study area, the training samples represent labeled dataset from the forests, bareland and vegetation classes. The number of trees for the RFA method with the above training samples are set to 100 and the maximum number to split parameter is set to 10. While for SVM, the regularization parameter is set to 1.2 and a linear SVM is utilized for the training purpose. For both RFA and SVM methods, the maximum number of iterations are set to 500. *Figure 8a* shows raster image representing the false color composite 7-3-2 and it is created from the sharpened images of the fusion between the Sentinel-2 and Landsat-8 OLI data. *Figure 8b* shows an image raster created from the three PCA bands PC1, PC2 and PC3. The eigen values, eigen vectors and variance of PCA bands are given in *Table 1*. The highest accounted and cumulative variance are recorded for PCA band 1 with a score of 76.56, for PCA band 2 and 3, with scores of 11.28 and 87.85 respectively.

*Figure 8c, e* and *g* show data classified for the year 2016 using SVM, RFA and PCA-RFA methods respectively. While *Figure 8d, f* and *h* show the data classified for the year 2023 using the aforementioned three methods. From the classified data presented in *Figure 8c-h*, the classification reports are generated and shown in *Tables 2* and *3* respectively. Both the tables show the area proportion for each class as a pixel sum, % area and in square meters. Moreover, the classified area plots for the data of the year 2016 and 2023 using the aforementioned methods are shown in *Figure 9*.

*Figure 9a* and *b* show the classified area for each class using SVM method for the data of the years 2016 and 2023 respectively. For forest class, percent areas of 7.09% and 60.14% are recorded for the data in the years 2016 and 2023 respectively. While for the bareland class, percent areas of 50.89 and 30.87% are recorded for the data in the years 2016 and 2023 respectively. Similarly, *Figure 9c* and *d* show the classified area for each class using RFA method for the data of the years 2016 and 2023 respectively. For forest class, percent areas of 9.54% and 56.55% are recorded for the data in the years 2016 and 2023 respectively. While for the bareland class, percent areas of 72.5% and 23.39% are recorded for the data in the years 2016 and 2023 respectively. *Figure 9e* and *f* show the classified area for each class using PCA-RFA method for the data of the years 2016 and 2023 respectively. For forest class, percent areas of 7.4% and 49% are classified for the data in the years 2016 and 2023 respectively. While for the bareland class, percent areas of 65.57% and 40.57% are classified for the data in the years 2016 and 2023 respectively.



**Figure 8.** (a) Original raster (b) raster representing PCA components (c)2016 data classified using SVM (d) 2023 data classified using SVM (e) 2016 data classified using RFA (f) 2023 data classified using RFA (g) 2016 data classified using PCA- RFA (h) 2023 data classified using PCA-RFA

**Table 1.** Eigen vectors, eigen values and variance

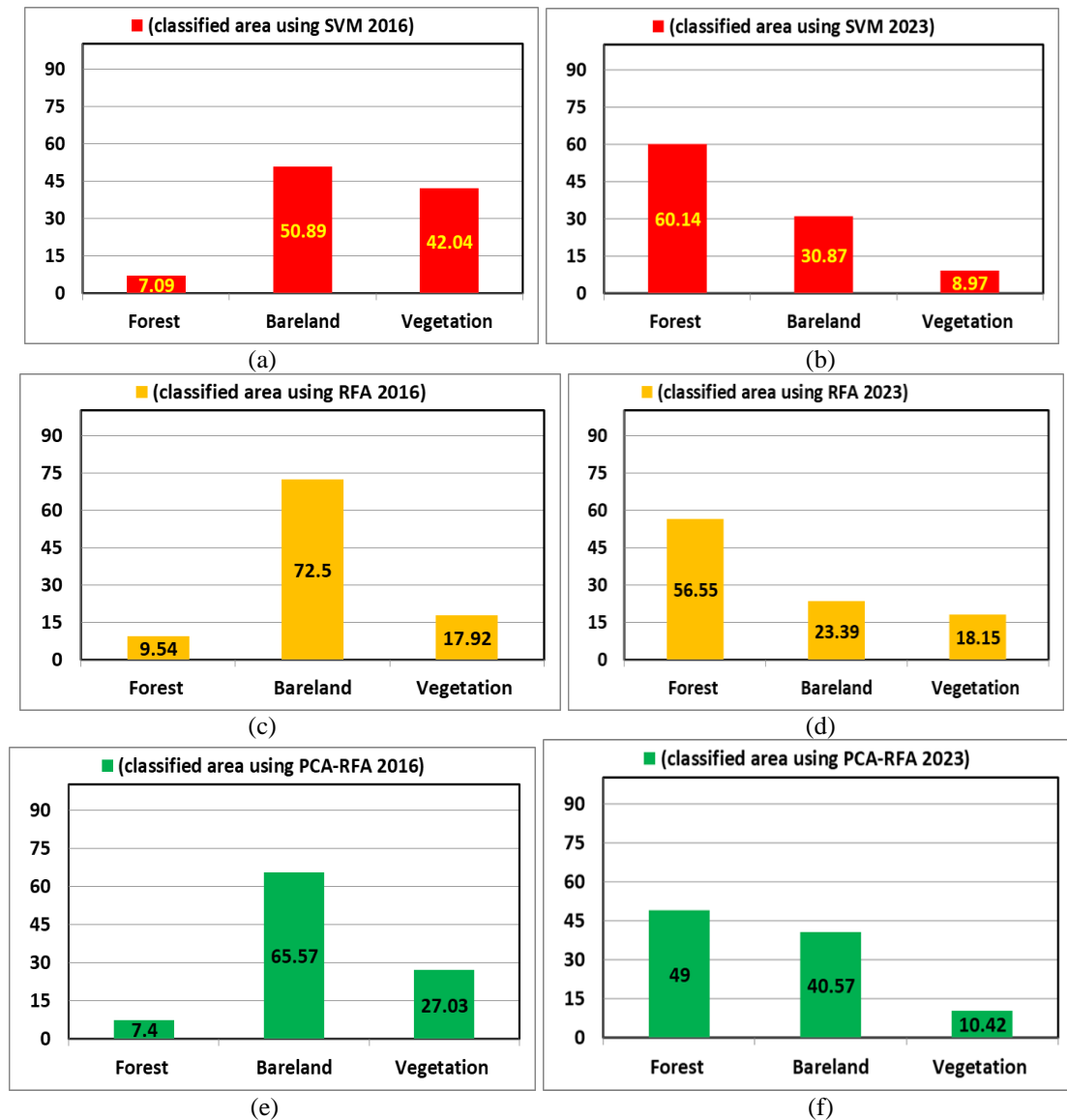
Bands	Vector1	Vector2	Vector3
1	0.1152	-0.0065	0.3401
2	0.1794	0.2535	0.1836
3	0.2423	0.2709	0.1943
Eigen values		Variance	Cumulative variance
0.0150		76.56	76.56
0.0022		11.28	87.85
0.0011		5.74	93.60

**Table 2.** Classified areas for each class using SVM, RFA and PCA-RFA (2023)

SVM			
Class	Pixel Sum	Percentage%	Area ( $\times 10^6$ m <sup>2</sup> )
Forest class	188,966	60.14	18.89
Bareland class	97,007	30.87	9.70
Vegetation class	28,187	8.972	2.81
RFA (Mateen et al., 2022)			
Forest class	177,679	56.55	17.76
Bareland class	79,766	25.39	7.97
Vegetation class	56,725	18.15	5.67
PCA-RFA			
Forest class	153,940	49.00	15.39
Bareland class	127,484	40.57	12.74
Vegetation class	32,736	10.42	3.27

**Table 3.** Classified areas for each class using SVM, RFA and PCA-RFA (2016)

SVM			
Class	Pixel Sum	Percentage%	Area ( $\times 10^6$ m <sup>2</sup> )
Forest class	220,96	7.09	2.200
Bareland class	19,990	50.87	15.99
Vegetation class	13,200	42.04	13.20
RFA (Mateen et al., 2022)			
Forest class	30047	9.56	3.00
Bareland class	227790	72.50	22.77
Vegetation class	56323	17.92	5.63
PCA-RFA			
Forest class	232,00	7.40	2.32
Bareland class	206,00	65.57	20.6
Vegetation class	85,736	27.03	8.5



**Figure 9.** (a) Classified area (a)2016 using SVM (b) 2023 using SVM (c) 2016 using RFA (d) 2023 using RFA (e) 2016 using PCA-RFA (f) 2023 using PCA-RFA

### Accuracy assessment

For the estimation of the accuracy matrices, the target standard deviation is chosen as  $\delta_o = 0.01$  and the standard deviation of each class is selected based on area proportion and calculated as follows:  $\delta_{forest} = 0.1, \delta_{bareland} = 0.2, \delta_{vegetation} = 0.3$ , where  $n = 3$  represents the number of classes naming forest, bareland and vegetation. From the classification report tabulated in Table 2, the number of samples for each class are calculated and given in Table 4.

Table 5 shows the estimated accuracy assessment matrices of the year 2023. Using SVM classifier, an overall percent accuracy score of 72.89% is estimated with an overall Kappa hat score of 0.55, while with RFA classifier; the overall percent accuracy of 92.87% and Kappa hat score of 0.87 is observed. The best overall percent accuracy is estimated with PCA-RFA technique with a score of 95.8% and with an overall Kappa

hat of 0.93. Using PCA-RFA technique, the percent producer accuracy PA [%] scores are measured as 98.28, 97.42 and 82.41, while the percent user accuracy UA [%] are estimated as 96.26, 94.79 and 98.24 for the forest, bareland and vegetation classes respectively. Similarly using RFA method, the percent producer accuracy PA [%] scores are measured as 96.28, 87.42 and 90.12, while the percent user accuracy UA [%] are estimated as 96.22, 92.72 and 82.60 for the forest, bareland and vegetation classes respectively. Using SVM method, the percent producer accuracy PA [%] scores are measured as 91.4, 72.4 and 40.1, while the percent user accuracy UA [%] are estimated as 69.9, 74.6 and 86.90 for the forest, bareland and vegetation classes respectively. Similarly, *Table 6* shows the estimated accuracy assessment matrices of the year 2016.

**Table 4.** Samples stratification using SVM, RFA and PCA-RFA

<b>SVM</b>			
<b>Class</b>	$Na_i$	$Nn^{-1}$	<b>Average</b>
Forest class	133	73	103
Bareland class	68	74	71
Vegetation class	19	73	46
<b>Total</b>	<b>220</b>	<b>220</b>	<b>220</b>
<b>RFA (Mateen et al., 2022)</b>			
Forest class	147	87	117
Bareland class	66	86	76
Vegetation class	47	87	67
<b>Total</b>	<b>260</b>	<b>260</b>	<b>260</b>
<b>PCA-RFA</b>			
Forest class	127	87	107
Bareland class	106	86	96
Vegetation class	27	87	57
<b>Total</b>	<b>260</b>	<b>260</b>	<b>260</b>

In order to compare the estimated and classified areas for each class and with each classifier, the results are plotted in *Figure 10*. *Figure 10a* and *b* show the estimated area for each class using SVM method and for the data of the years 2016 and 2023 respectively. From the presented results, it is concluded that using SVM method, the classified and estimated areas for bareland class closely match each other. For the data of the year 2016, the classified and estimated areas for bareland class are 50.89% and 47.12% respectively, while for 2023, the aforementioned areas are 30.87% and 31.81% respectively. Moreover, for the forest classes, the classified and estimated areas for the years 2016 and 2023 are 7.02%, 2.2% and 60.14%, 45.97% respectively. Thus, the matching of the forest class is not very satisfactory. For vegetation class, the classified and estimated areas for the years 2016 and 2023 are 42.04%, 50.67% and 8.97% and 22.23% respectively. From *Figure 10c-d* and *e-f* it is evident that a better match is provided for all classes with both PCA-RFA and RFA methods. Further details are available in the aforementioned figures. The overall accuracy and Kappa hat scores for the data of the years 2016 and 2023 are shown in *Figure 11a-b* and a comparison is provided between SVM, RFA and PCA-RFA methods. From the presented results for

the year 2016, PCA-RFA method has the highest overall accuracy and Kappa hat scores (91% and 0.92 respectively), while with RFA method, the overall accuracy is 89% with a Kappa hat score of 0.88. In this particular case, SVM showed poor performance. Similarly for the year 2023, PCA-RFA method has the highest overall accuracy and Kappa hat scores (95% and 0.93 respectively), while with RFA method, the overall accuracy is 92% with a Kappa hat score of 0.87. From *Figure 11b*, in case of SVM, the overall accuracy is 72% with a Kappa hat score of 0.55.

**Table 5.** Accuracy assessment parameters (2023)

<b>SVM</b>				
	Reference			
		Forest class	Bareland class	Vegetation class
Classified	Forest class	0.4205	0.0759	0.1051
	Bare land class	0.0391	0.2305	0.0391
	Vegetation class	0.0001	0.0117	0.0780
	Total% classified area	0.4597	0.3181	0.2223
	Area error	0.030	0.025	0.026
	Producer accuracy [%]	91.4	72.4	40.1
	User accuracy [%]	69.9	74.6	86.9
	Kappa hat	0.44	0.62	0.81
	Estimated area ( $\times 10^6$ m <sup>2</sup> )	14.4389	9.9940	6.9830
<b>RFA (Mateen et al., 2022)</b>				
	Reference			
		Forest class	Bareland class	Vegetation class
Classified	Forest class	0.5442	0.0142	0.0071
	Bare land class	0.0092	0.2354	0.0092
	Vegetation class	0.0118	0.0196	0.1492
	Total% classified area	0.5652	0.2692	0.1655
	Area error	0.0131	0.0126	0.0142
	Producer accuracy [%]	96.28	87.42	90.12
	User accuracy [%]	96.22	92.72	82.60
	Kappa hat	0.91	0.90	0.79
	Estimated area ( $\times 10^6$ m <sup>2</sup> )	17.7564	8.4600	5.1995
<b>PCA-RFA</b>				
	Reference			
		Forest class	Bareland class	Vegetation class
Classified	Forest class	0.4717	0.0092	0.0092
	Bare land class	0.0085	0.3847	0.0127
	Vegetation class	0.0001	0.0018	0.1024
	Total% classified area	0.4803	0.3957	0.1143
	Area error	0.0097	0.0101	0.0100
	Producer accuracy [%]	98.28	97.42	82.41
	User accuracy [%]	96.26	94.79	98.24
	Kappa hat	0.92	0.91	0.98
	Estimated area ( $\times 10^6$ m <sup>2</sup> )	15.0841	12.4295	3.9022

**Table 6.** Accuracy assessment parameters (2016)

<b>SVM</b>				
	Reference			
		Forest class	Bareland class	Vegetation class
Classified	Forest class	0.0105	0.0700	0.1051
	Bare land class	0.0105	0.2305	0.3391
	Vegetation class	0.0011	0.1707	0.1625
	Total% classified area	0.0221	0.4712	0.5067
	Area error	0.0488	0.0375	0.0863
	Producer accuracy [%]	78.4	75.4	55.1
	User accuracy [%]	77.9	74.6	76.9
	Kappa hat	0.55	0.67	0.49
	Estimated area ( $\times 10^6$ m <sup>2</sup> )	0.694	14.8	15.91
<b>RFA (Mateen et al., 2022)</b>				
	Reference			
		Forest class	Bareland class	Vegetation class
Classified	Forest class	0.0348	0.6776	0.0127
	Bare land class	0.0000	0.0097	0.1695
	Vegetation class	0.0830	0.0126	0.0001
	Total% classified area	0.1179	0.6999	0.1823
	Area error	0.0108	0.0130	0.0074
	Producer accuracy [%]	70.44	96.81	93.05
	User accuracy [%]	86.81	93.44	94.59
	Kappa hat	0.85	0.78	0.93
	Estimated area ( $\times 10^6$ m <sup>2</sup> )	3.7026	21.9875	5.7257
<b>PCA-RFA</b>				
	Reference			
		Forest class	Bareland class	Vegetation class
Classified	Forest class	0.0148	0.2162	0.0107
	Bare land class	0.0220	0.2111	0.1191
	Vegetation class	0.0831	0.2126	0.1101
	Total% classified area	0.1199	0.6399	0.2399
	Area error	0.045	0.015	0.030
	Producer accuracy [%]	78.44	89.81	91.05
	User accuracy [%]	81.81	91.44	92.59
	Kappa hat	0.89	0.86	0.92
	Estimated area ( $\times 10^6$ m <sup>2</sup> )	3.766	20.10	7.53

### Change map

The change map comparison for the years 2016-2023 is shown in *Figure 12a-f* using SVM, RFA and PCA-RFA methods respectively. *Figure 12a-b* shows the respective % area changes and the change map using SVM method. From the presented results it is evident that there is a gross gain of 17.4% and 18.9% for the vegetation 2016–forest 2023 and bareland 2016–forest 2023 cross classes respectively. While a gross loss of -31.2% is observed in the bareland 2016-bareland 2023 cross class. Similarly, *Figure 12c-d* and *e-f* show the respective % area changes and the change map using RFA and PCA-RFA methods respectively.



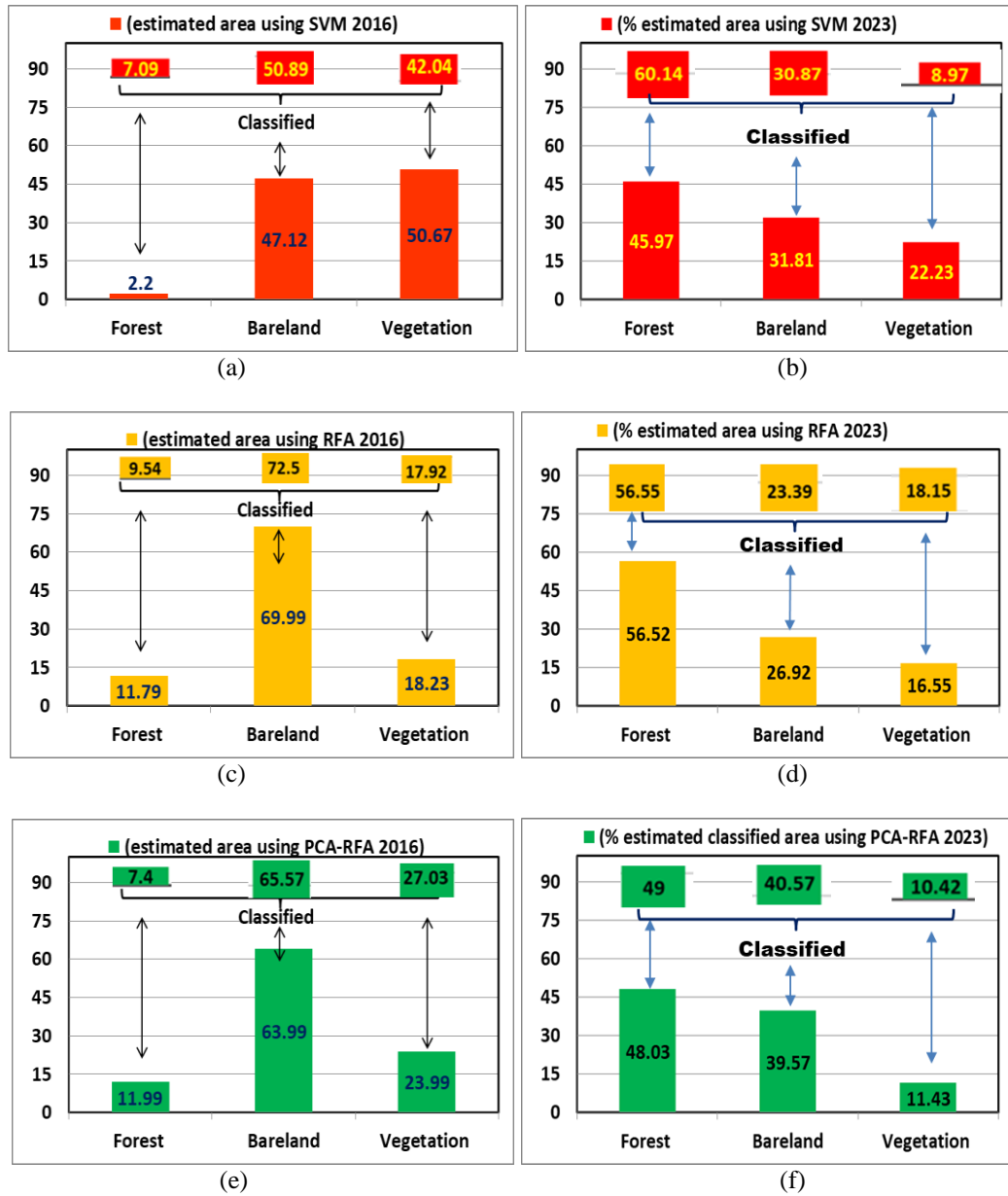


Figure 10. Estimated vs classified area. (a) 2016 using SVM (b) 2023 using SVM (c) 2016 using RFA (d) 2023 using RFA (e) 2016 using PCA-RFA (f) 2023 using PCA-RFA

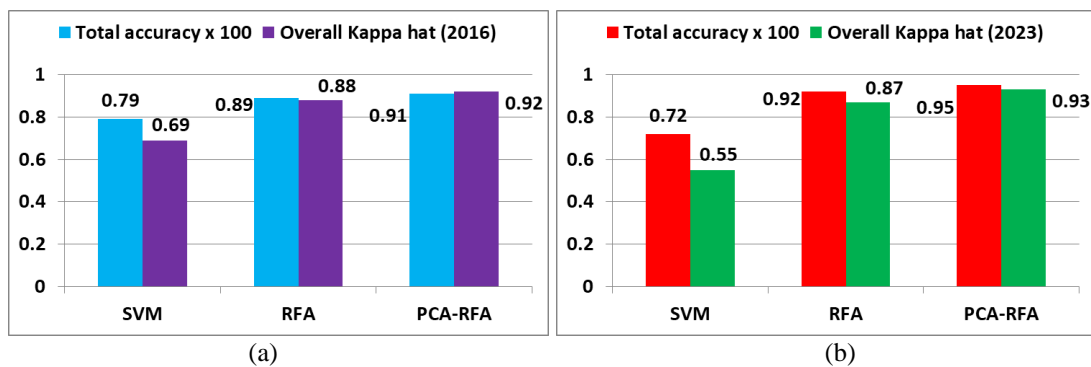
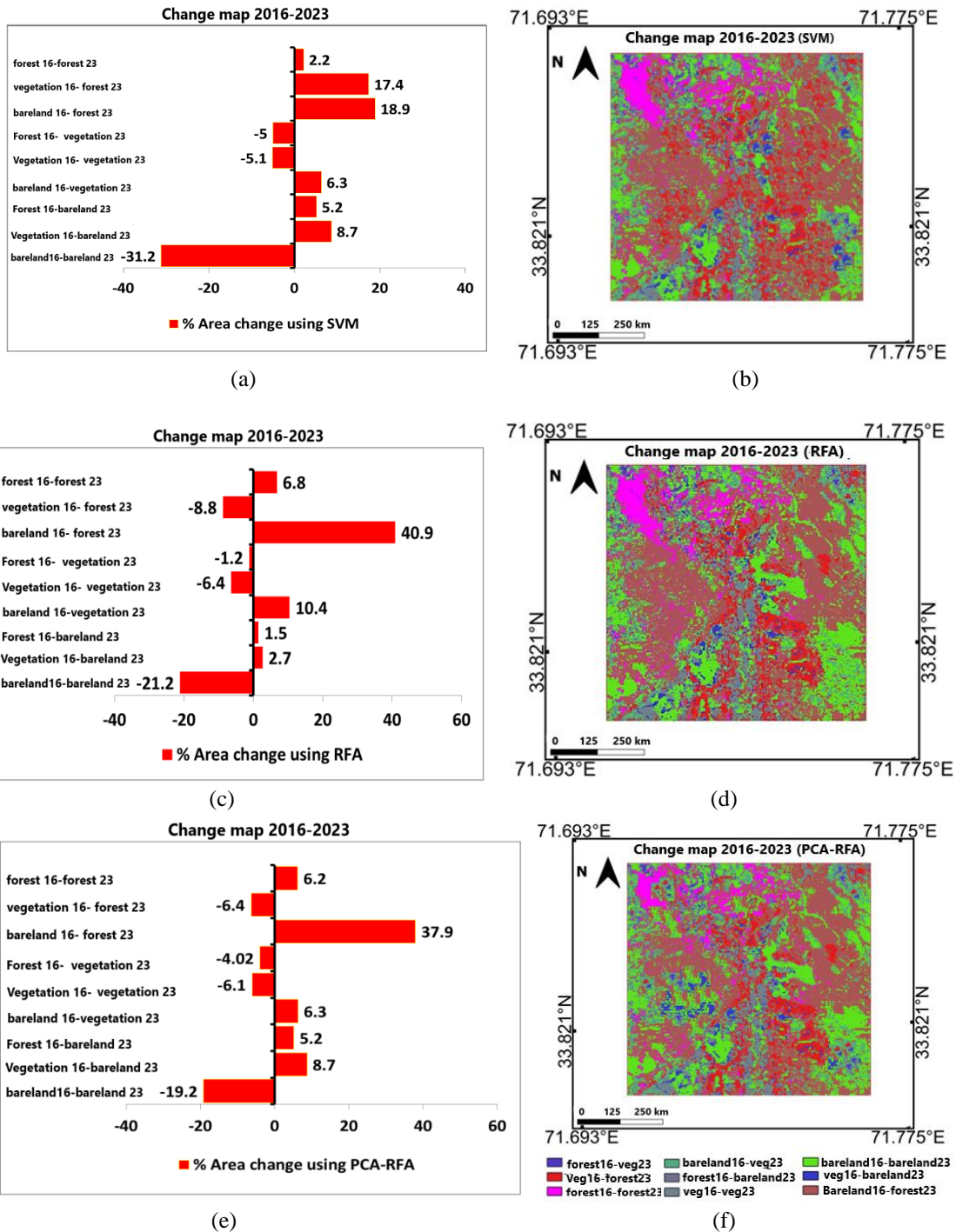


Figure 11. Overall accuracy and Kappa hat score (a) 2016 (b) 2023



**Figure 12.** (a) % area change using SVM (b) change map using SVM (c) % area change using RFA (d) change map using RFA (e) % area change using PCA-RFA (f) change map using PCA-RFA

From the presented results it is evident that there is a gross gain of 40.9% for the bareland 2016–forest 2023 cross class, while a gross loss of -21.2% is observed in the bareland 2016–bareland 2023 cross class. From *Figure 12e-f* a gross gain of 37.9% for the bareland 2016–forest 2023 cross class, while a gross loss of -19.2% is observed in the bareland 2016–bareland 2023 cross class.

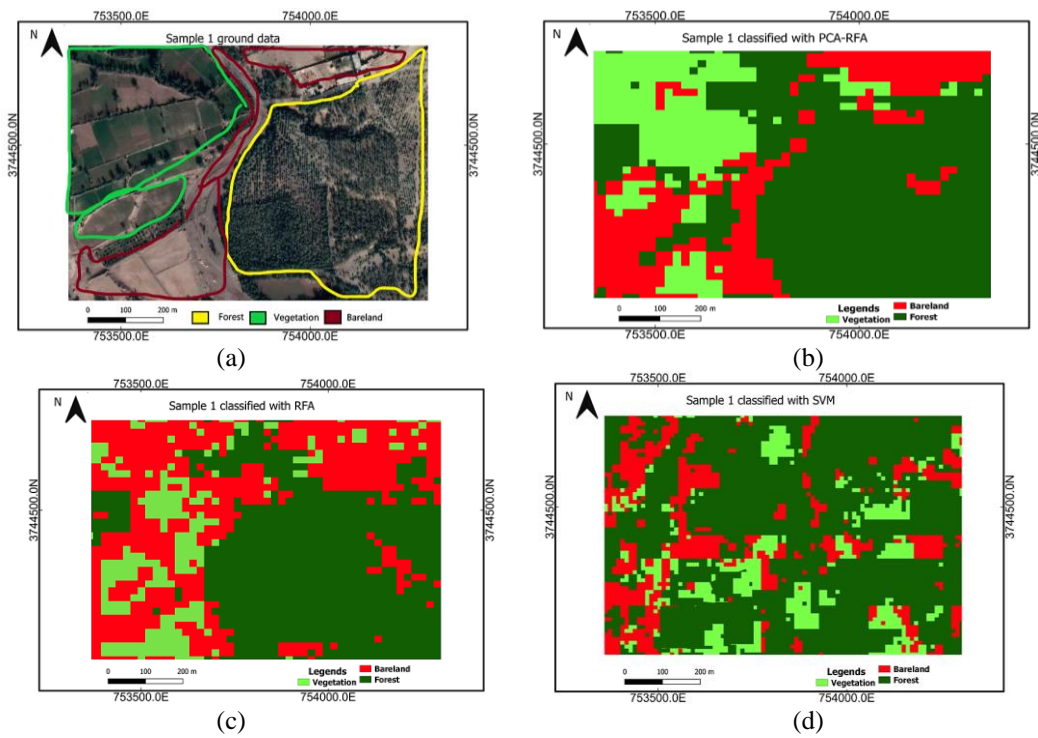
### **Ground data matching**

For ground data validation, the ground samples are collected using Google earth and then each ground sample is compared with the corresponding classified sample. It is hard to provide a qualitative comparison, thus for a quantitative analysis, the area proportion for each classified class in the corresponding is calculated based on the classification reports. While the area proportion of each class in the ground sample is calculated using the area calculator tool of QGIS. From the ground area and the classified area for each class the percentage area error is calculated. *Figure 13a* shows ground sample 1, where the area encircled by yellow polygons shows the forests area, red polygons mean bareland while the vegetation is encircled using green polygons. Comparisons of the ground and classified Sample 1 using PCA-RFA, RFA and SVM are shown in *Figure 13b, c* and *d* respectively. *Figure 17a* shows the % area error in each class using SVM, RFA and PCA-RFA methods. The percentage area error in the forest class using PCA-RFA and RFA methods is observed as 4.6% and 4.2% respectively. While with SVM the percent error in the forest class is noted as 15.7%. Similarly, the percentage area error in the bareland class using PCA-RFA and RFA methods is observed as 2.1% and 8.1% respectively. While with SVM the percentage error in the forest class is noted as 23.8%. *Figure 14a* shows ground sample 2. The classified Sample 2 using PCA-RFA, RFA and SVM are shown in *Figure 14b, c* and *d* respectively. *Figure 17b* shows the % area error in each class using SVM, RFA and PCA-RFA methods. The percentage area error in the forest class using PCA-RFA and RFA methods is observed as 7.2% and 13.2% respectively. While with SVM the percent error in the forest class is noted as 19.3%. Similarly, the percent area error in the bareland class using PCA-RFA and RFA methods is observed as 5.1% and 18.1% respectively. While with SVM the percent error in the forest class is noted as 21.9%. Similarly, *Figure 15a* shows ground sample 3. The classified Sample 3 using PCA-RFA, RFA and SVM are shown in *Figure 15b, c* and *d* respectively. *Figure 17c* shows the % area error in each class using SVM, RFA and PCA-RFA methods. The percentage area error in the forest class using PCA-RFA and RFA methods is observed as 4% and 9.2% respectively. While with SVM the percentage error in the forest class is noted as 8.4%. Similarly, the percentage area error in the bareland class using PCA-RFA and RFA methods is observed as 3.2% and 4.1% respectively. While with SVM the percent error in the forest class is noted as 13.5%. *Figure 16a* shows ground sample 4. The classified Sample 4 using PCA-RFA, RFA and SVM are shown in *Figure 16b, c* and *d* respectively. *Figure 17d* shows the % area error in each class using SVM, RFA and PCA-RFA methods. The presented results in *Figure 17d* show that a similar percent area error trend is observed for each class in sample 4 as it was observed in the aforementioned samples.

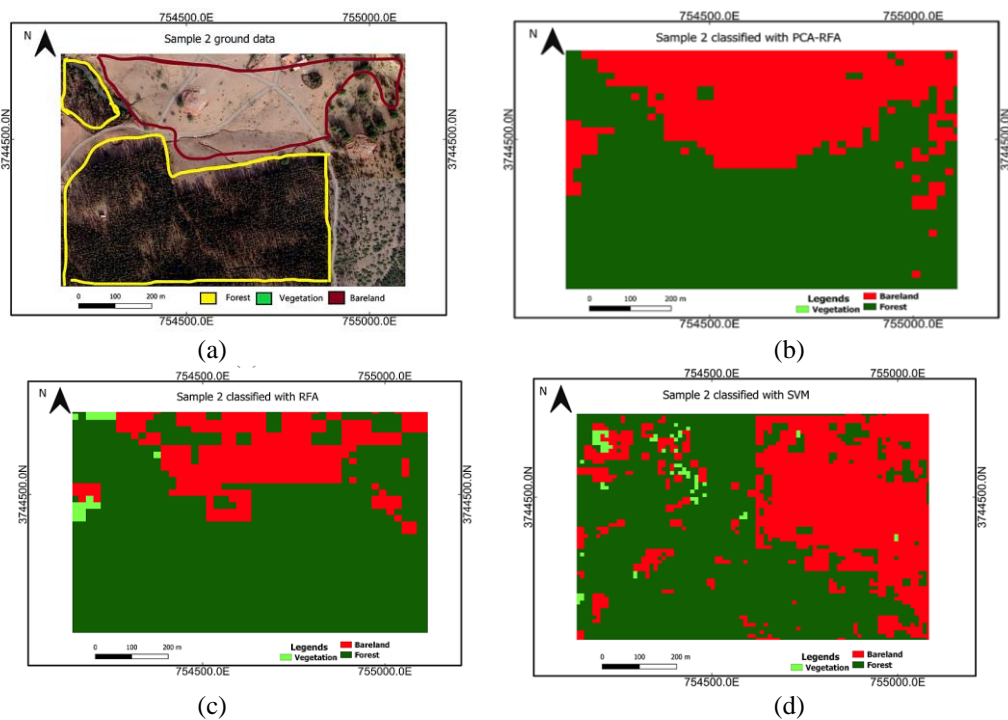
### **Change prediction for the years 2028, 2030 and 2033.**

In order to predict the future classified mapping of our study, Molusce toolbox of QGIS is utilized. The transition modelling was done using logistic regression method with 1 pixel size and 1000 iterations. As a first step the classified data for the years 2016 and 2018 are utilized to predict the data for the year 2020. As a second step, the predicted and the classified data are validated with overall Kappa score of 0.96. Afterwards, the classified data for the years 2016 and 2023 are utilized to predict the data for the year 2028. Later on, the classified data for the years 2023 and the predicted data 2028 are utilized to predict the data for the year 2030. Similar procedure is used to predict the data for the year 2033. *Figure 18a-c* shows the predicted data for the years 2028, 2030 and 2033 respectively.

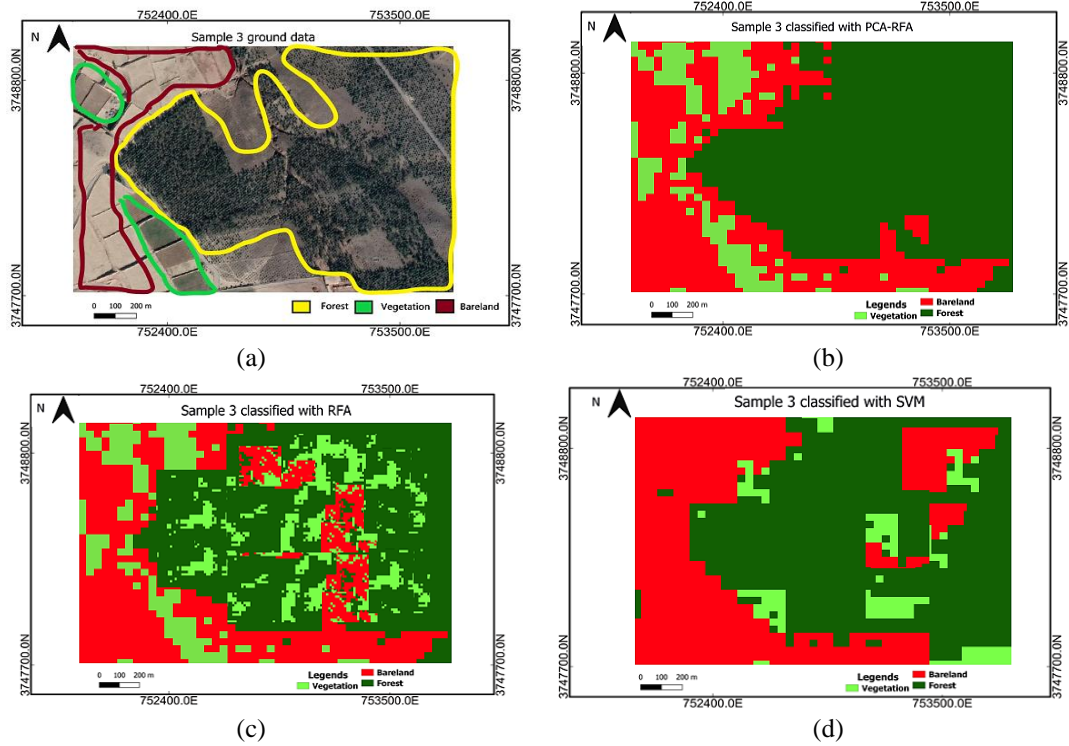
From the qualitative analysis of the predicted data it is concluded the bareland class pixels show decreasing trend while the forest class pixels show increasing trend.



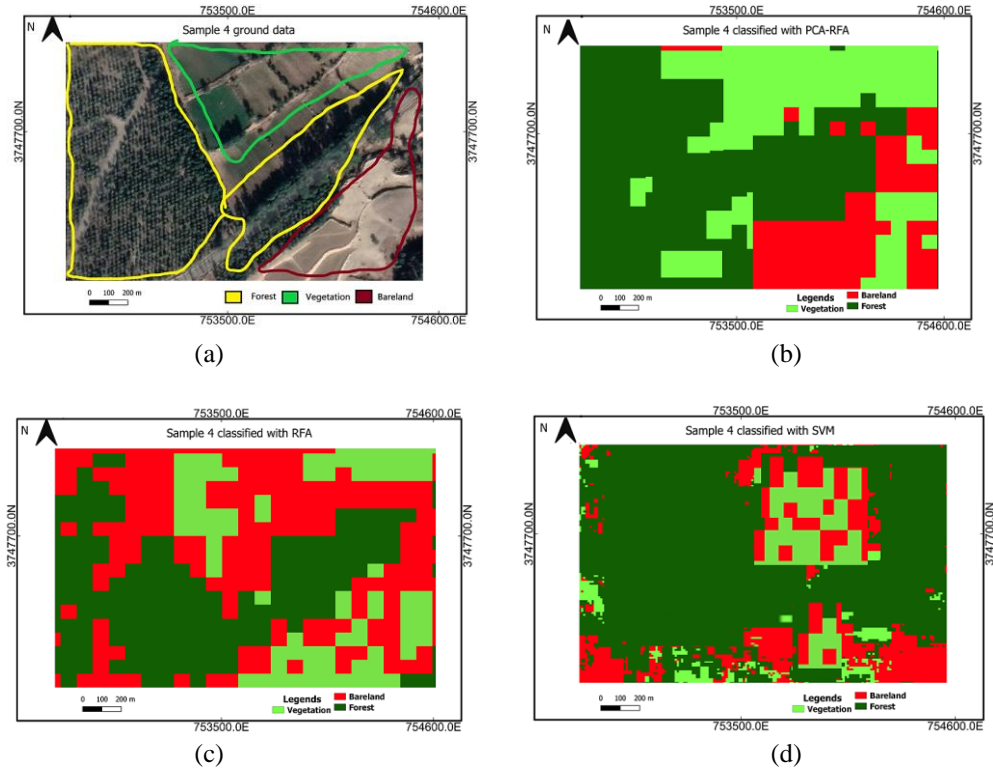
**Figure 13.** (a) Ground sample1 (b) classified sample 1 using PCA-RFA (c) classified sample1 using RFA (d) classified sample 1 using SVM



**Figure 14.** (a) Ground sample2 (b) classified sample 2 using PCA-RFA (c) classified sample2 using RFA (d) classified sample 2 using SVM



**Figure 15.** (a) Ground sample 3 (b) classified sample 3 using PCA-RFA (c) classified sample 3 using RFA (d) classified sample 3 using SVM



**Figure 16.** (a) Ground sample 4 (b) classified sample 4 using PCA-RFA (c) classified sample 4 using RFA (d) classified sample 4 using SVM

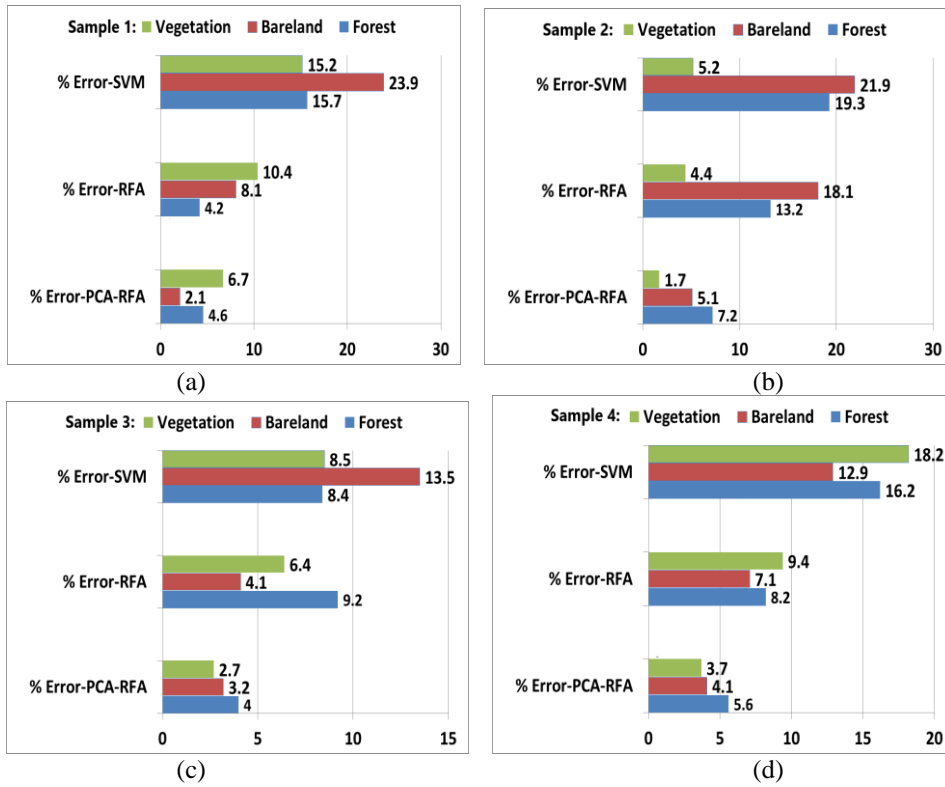


Figure 17. (a) % Area error sample 1 (b) % area error sample 2 (c) % area error sample 3 (d) % area error sample 4

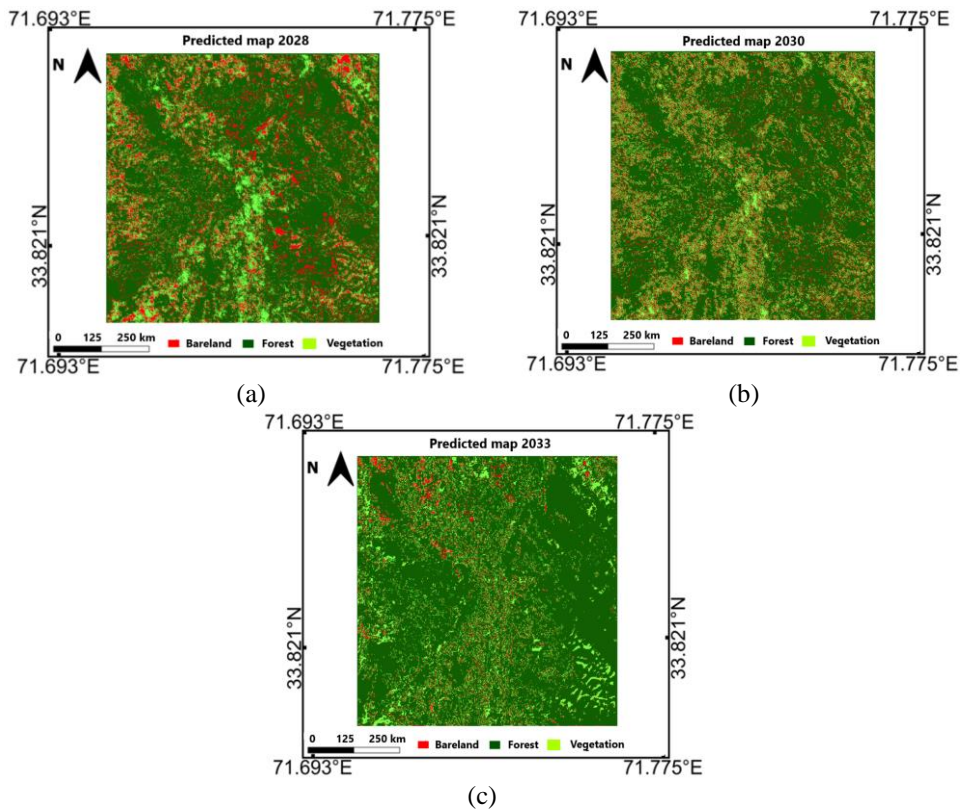


Figure 18. Predicted map (a) 2028 (b)2030 (c) 2033

## Discussion

From the presented results, it is worth mentioning that with PCA-RFA has attained the higher accuracy and Kappa hat scores because the covariance matrix of the principal components has all zero elements except the main diagonal, which means that the principal components are no more related to each other (Liyanage et al., 2019; Uddim et al., 2020). The covariance matrix for the three principal components of the PCA is calculated as follows:

$$\text{cov}_{(pc)} = \begin{bmatrix} 0.0150 & 0 & 0 \\ 0 & 0.0022 & 0 \\ 0 & 0 & 0.0011 \end{bmatrix}$$

For PCA-RFA, the overall accuracy and Kappa hat scores are estimated as 95.8719% and 0.9302 respectively. The lowest overall accuracy and Kappa hat scores of 72.8982% and 0.5523 are recorded for SVM. While in the case of RFA classifier the overall accuracy and Kappa hat score of 92.87% and 0.9777 are observed. RFA performs better as compared to SVM and this finding is according to the research data presented in Adugna et al. (2022) and Avci et al. (2023). PCA-RFA outperforms RFA, and the results obtained in this research work are according to the analysis given in Xia et al. (2017) and Jin and Bie (2006). The collected ground samples are shown in *Figures 13a, 14a, 15a and 16a*. The classified samples are shown in *Figure 13b-d, 13(b-d), 14(b-d) and 15(b-d)*. The % area error of each class by using the aforementioned three methods (*Fig. 17a-d*) confirms the best match provided by the PCA-RFA classifier for all four samples (Xia et al., 2017; Jin and Bie, 2006), while the performance of RFA classifier is also good (Adugna et al., 2022; Avci et al. 2023), however the SVM classifier showed poor mapping performance. Moreover, the predicted maps for our study area for the years 2028, 2030 and 2033 confirm that forest class will dominate the whole study area with fewer percentage of the vegetation class.

## Conclusions

From the presented results, it is concluded that the best accuracy is achieved using PCA-RFA classifier with an overall accuracy of 95.87% and Kappa hat score 0.93. The second-best technique is RFA with overall accuracy and Kappa hat score of 92.87% and 0.97. SVM showed moderate mapping performance with an overall accuracy and Kappa hat scores of 72.89% and 0.55 respectively. PCA-RFA provided the best scores of statistical parameters and less % area error in the validation test. Thus, the temporal analysis with PCA-RFA is the most accurate. A potential future extension to this work is to generate future land cover changes with the Molusce toolbox of QGIS.

## REFERENCES

- [1] Adugna, T., Xu, W., Fan, J. (2022): Comparison of random forest and support vector machine classifiers for regional land cover mapping using coarse resolution FY-3C images. – *Remote Sensing* 14(3): 574.

- [2] Avci, C., Budak, M., Yağmur, N., Balçık, F. (2023): Comparison between random forest and support vector machine algorithms for LULC classification. – *International Journal of Engineering and Geosciences* 8(1): 1-10.
- [3] Bazi, Y., Melgani, F. (2006): Toward an optimal SVM classification system for hyperspectral remote sensing images. – *IEEE Transactions on Geoscience and Remote Sensing* 44(11): 3374-3385.
- [4] Boulesteix, A., Janitza, S., Kruppa, J., König, I. R. (2012): Overview of random forest methodology and practical guidance with emphasis on computational biology and bioinformatics. – *WIREs Data Mining and Knowledge Discovery* 2(6): 493-507.
- [5] Breiman, L. (1996): Bagging predictors. – *Machine Learning* 24(2): 123-140.
- [6] Breiman, L. (2001): ST4\_Method\_Random\_Forest. – *Mach. Learn.* 45(1): 5-32.
- [7] Congedo, L. (2021): Semi-automatic classification plugin: a Python tool for the download and processing of remote sensing images in QGIS. – *Journal of Open Source Software* 6(64): 3172.
- [8] Eklundh, L., Singh, A. (1993): A comparative analysis of standardised and unstandardised principal components analysis in remote sensing. – *International Journal of Remote Sensing* 14(7): 1359-1370.
- [9] Fawagreh, K., Gaber, M. M., Elyan, E. (2014): Random forests: from early developments to recent advancements. – *Systems Science & Control Engineering* 2(1): 602-609.
- [10] Habib, T., Inglada, J., Mercier, G., Chanussot, J. (2009): Support vector reduction in SVM algorithm for abrupt change detection in remote sensing. – *IEEE Geoscience and Remote Sensing Letters* 6(3): 606-610.
- [11] Heydari, S. S., Mountrakis, G. (2018): Effect of classifier selection, reference sample size, reference class distribution and scene heterogeneity in per-pixel classification accuracy using 26 Landsat sites. – *Remote Sensing of Environment* 204: 648-658.
- [12] Heydari, S. S., Mountrakis, G. (2019): Meta-analysis of deep neural networks in remote sensing: a comparative study of mono-temporal classification to support vector machines. – *ISPRS Journal of Photogrammetry and Remote Sensing* 152: 192-210.
- [13] Jin, X., Bie, R. (2006): Random forest and PCA for self-organizing maps based automatic music genre discrimination. – *Proceedings of the 2006 International Conference on Data Mining, DMIN 2006, Las Vegas, Nevada, USA, June 26-29, 2006.*
- [14] Kandpal, K. C., Kumar, A. (2022): Identification and classification of medicinal plants of the Indian Himalayan region using hyperspectral remote sensing and random forest techniques. – *2022 IEEE Mediterranean and Middle-East Geoscience and Remote Sensing Symposium (M2GARSS).*
- [15] Kavzoglu, T., Colkesen, I. (2009): A kernel functions analysis for support vector machines for land cover classification. – *International Journal of Applied Earth Observation and Geoinformation* 11(5): 352-359.
- [16] Khemchandani, R., Jayadeva, Chandra, S. (2008): Optimal kernel selection in twin support vector machines. – *Optimization Letters* 3(1): 77-88.
- [17] Khosravi, I., Jouybari-Moghaddam, Y. (2019): Hyperspectral imbalanced datasets classification using filter-based forest methods. – *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing* 12(12): 4766-4772.
- [18] Liu, G., Mao, S., Kim, J. H. (2019): A mature-tomato detection algorithm using machine learning and color analysis. – *Sensors* 19(9): 2023.
- [19] Liu, J., Li, P. (2019): Landslide mapping and analysis using multi-source data and one-class random forest. – *IGARSS 2019 - 2019 IEEE International Geoscience and Remote Sensing Symposium.*
- [20] Liyanage, L. C., Weerakoon, O. S., Palliyaguru, S. T., Wimalaratne, G. D. S. P. (2019): Towards prediction of landslide susceptibility using random forest for Kalutara District, Sri Lanka. – *2019 IEEE R10 Humanitarian Technology Conference (R10-HTC)(47129).*
- [21] Maćkiewicz, A., Ratajczak, W. (1993): Principal components analysis (PCA). – *Computers & Geosciences* 19(3): 303-342.



- [22] Mainali, K., Evans, M., Saavedra, D., Mills, E., Madsen, B., Minnemeyer, S. (2023): Convolutional neural network for high-resolution wetland mapping with open data: variable selection and the challenges of a generalizable model. – *Science of the Total Environment* 861: 160622.
- [23] Mateen, S., Nuthammachot, N., Techato, K., Ullah, N. (2022): Billion tree tsunami forests classification using image fusion technique and random forest classifier applied to Sentinel-2 and Landsat-8 images: a case study of Garhi Chandan Pakistan. – *ISPRS International Journal of Geo-Information* 12(1): 9.
- [24] Mather, P., Tso, B. (2016): Classification methods for remotely sensed data. – <http://dx.doi.org/10.1201/9781420090741>.
- [25] Miao, X., Heaton, J. S. (2010): A comparison of random forest and Adaboost tree in ecosystem classification in east Mojave Desert. – 2010 18th International Conference on Geoinformatics.
- [26] Mohammadimanesh, F., Salehi, B., Mahdianpari, M., Gill, E., Molinier, M. (2019): A new fully convolutional neural network for semantic segmentation of polarimetric SAR imagery in complex land cover ecosystem. – *ISPRS Journal of Photogrammetry and Remote Sensing* 151: 223-236.
- [27] Mountrakis, G., Im, J., Ogole, C. (2011): Support vector machines in remote sensing: a review. – *ISPRS Journal of Photogrammetry and Remote Sensing* 66(3): 247-259.
- [28] Olofsson, P., Foody, G. M., Herold, M., Stehman, S. V., Woodcock, C. E., Wulder, M. A. (2014): Good practices for estimating area and assessing accuracy of land change. – *Remote Sensing of Environment* 148: 42-57.
- [29] Pouteau, R., Stoll, B. (2012): SVM selective fusion (SELF) for multi-source classification of structurally complex tropical rainforest. – *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing* 5(4): 1203-1212. <http://dx.doi.org/10.1109/JSTARS.2012.2183857>.
- [30] Sales, M. H. R., de Bruin, S., Souza, C., Herold, M. (2022): Land use and land cover area estimates from class membership probability of a random forest classification. – *IEEE Transactions on Geoscience and Remote Sensing* 60: 1-11.
- [31] Schölkopf, B., Smola, A. J. (2018): *Learning with Kernels*. – MIT Press, Cambridge, MA. <http://dx.doi.org/10.7551/mitpress/4175.001.0001>.
- [32] Sheykhmousa, M., Mahdianpari, M., Ghanbari, H., Mohammadimanesh, F., Ghamisi, P., Homayouni, S. (2020): Support vector machine versus random forest for remote sensing image classification: a meta-analysis and systematic review. – *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing* 13: 6308-6325. DOI: <https://doi.org/10.1109/jstars.2020.3026724>.
- [33] Sigurdsson, J., Armannsson, S. E., Ulfarsson, M. O., Sveinsson, J. R. (2022): Fusing Sentinel-2 and Landsat 8 satellite images using a model-based method. – *Remote Sensing* 14(13): 3224.
- [34] Singh, A., Harrison, A. (1985): Standardized principal components. – *International Journal of Remote Sensing* 6(6): 883-896.
- [35] Uddin, M. P., Mamun, M. A., Hossain, M. A. (2020): PCA-based feature reduction for hyperspectral remote sensing image classification. – *IETE Technical Review* 38(4): 377-396.
- [36] Xia, J., Falco, N., Benediktsson, J. A., Du, P., Chanussot, J. (2017): Hyperspectral image classification with rotation random forest via KPCA. – *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing* 10(4): 1601-1609.