



FOREST

MONITOR

REPORT

2025 - 2

- Methodology to obtain the potential location of deciduous forest with high conservation values in southern Sweden



Report 2025 - 2

FOREST MONITOR

- Methodology to obtain the potential location of deciduous forest with high conservation values in southern Sweden

Author

Hainner Aparicio

Proofreading

Jon Andersson

Graphics

Hainner Aparicio & Jon Andersson

Layout

Jon Andersson

Acknowledgments

We want to express our greatest appreciation to the Postcode Lottery Foundation, to Naturkompaniet and to our private donors!

Cover photo

Jon Andersson

INTRODUCTION

Sweden is home to a significant part of the EU's natural heritage not only in the form of valuable coniferous forests, but also older deciduous forests (DF) with high conservation values (HCV). Older deciduous forests often have unique flora, fungi, and fauna associated with them. To meet national and international environmental goals and EU laws and regulations, Sweden must protect all forests with high natural values and restore strategic areas around them.

An important step in the work to protect and restore vital forest ecosystems is to locate known and potential conservation values. This greatly streamlines inventory efforts and contributes to a unique overview of the values in forest landscapes.

The map tool Skogsmonitor.se presents data on known and potential conservation values for the entire country, linked to older forests and so-called continuity forests, forests that have not previously been clear-cut. To find potential older forests and continuity forests, Skogsmonitor has used historical aerial photographs and satellite images covering the entire country and then classified these using inventory data and algorithms that utilize findings of species of conservation interest.

There are challenges in using classic analysis of aerial and satellite images, as these do not always capture deciduous forests with potentially high conservation values (DFs with potential HCV) in a satisfactory way. New tools, like machine learning, offer a great opportunity by utilizing computational power to analyze massive loads of data from different sources in a consolidated/rapid yet detailed and robust manner.

We have used a machine learning approach to produce a map layer with predictions of where the most valuable deciduous forests can be found in southern Sweden. The map layer of deciduous forests with potential conservation values was created using a random forest classifier trained on satellite images, vegetation maps, terrain data, and data from tree and forest surveys.

The purposeful methodology attempted in this study and its results help locate these important forests within reasonable accuracy of over 90 percent (Accuracy > 0.9; Kappa > 0.8; RMSE < 0.3).

In this report we describe the methods used for the prediction of locations of DFs with potential HCV. We describe the datasets construction, pre-processing steps, modelling and the post processing steps that were carried out.

METHODOLOGY

Study area and data sources

The study area covers the southern part of Sweden except for the counties Värmland, Dalarna and Gotland, see **Figure 1**.



Figure 1. Map showing the county-wise coverage of the study area.

Reference data

The rate of success or failure in detecting various forest features with any kind of algorithm relies heavily on the construction of a solid reference dataset and well-chosen variables for modelling. We will give an exhaustive explanation of how we went about choosing the data and how we prepared the data sets for the analysis.

We used two different open data sources, where we downloaded datasets which we later mined, filtered and cleared to select specific places where the deciduous forests with high conservation value have been registered/reported. These forests correspond in ecological characteristics with our targeted forests (DFs with HCV). Therefore, they served as reference locations for creating a consistent reference dataset and further training and testing of the machine learning (ML) model.

Tabel 1. Data sources that were used to supervise the training procedure.

Data source	Forest type	Relevance
Naturvarden (Skogsstyrelsen)	Alsumpskog (ALSUMP)	High
	Aspskog (ASPSKOG)	High
	Bokskog (BOKSKOG_ONV)	High
	Hassellund (HASSLUND)	Medium
	Hedädellövskog (HEDÄDEL)	High
	Kalklövskog (KALKLÖV)	High
	Lövnaturskog (LÖVSKOG)	High
	Lövskog (LÖVSKOG_ONV)	High
	Lövskogslund (LÖVLUND)	Medium
	Lövskogslund/Hagmarksskog (LÖVLUND_ONV)	High
	Lövsumpskog (LÖVSUMP)	High
	Sekundär lövnaturskog (SEKNSKOG)	High
	Sekundär ädellövnaturskog (SEKÄDEL)	High
	Ädellövnaturskog (ÄDELLÖV)	High
	Ädellövskog (ÄDELBEST)	High
	Ädellövskog (ÄDELLÖV_ONV)	High
	Ädellövskog (ÄDELSKOG)	High
	Ädellövsumpskog (ÄDELSUMP)	High
	Örtrik allund (ALLUND)	High
	Övriga lövträd (ÖLÖVTRÄD)	High
Biotoskydd (Skogsstyrelsen)	Alkärr	High
	Hassellundar och hasselrika skogar	Medium
	Örtrika allundar	Medium
Nyckelbiotoper (Skogsstyrelsen)	Alsumpskog (ALSUMP)	High
	Aspskog (ASPSKOG)	High
	Bokskog (BOKSKOG_ONV)	High
	Hassellund (HASSLUND)	Medium
	Hedädellövskog (HEDÄDEL)	High
	Kalklövskog (KALKLÖV)	High
	Lövnaturskog (LÖVSKOG)	High
	Lövskog (LÖVSKOG_ONV)	High
	Lövskogslund (LÖVLUND)	Medium
	Lövskogslund/Hagmarksskog (LÖVLUND_ONV)	High
	Lövsumpskog (LÖVSUMP)	High
	Sekundär lövnaturskog (SEKNSKOG)	High
	Sekundär ädellövnaturskog (SEKÄDEL)	High
	Ädellövnaturskog (ÄDELLÖV)	High

Tabel 1, continued. Data sources that were used to supervise the training procedure.

Data source	Habitat type	Relevance
	Ädellövskog (ÄDELBEST)	High
	Ädellövskog (ÄDELSKOG)	High
	Ädellövsumpskog (ÄDELSUMP)	High
	Ädellövträäd (ÄDELTRÄD)	High
	Örtrik allund (ALLUND)	High
	Övriga lövträäd (ÖLÖVTRÄD)	High
Naturvardsavtal (Skogsstyrelsen)	Lövsumpskog	High
	Triviallövskog asp	High
	Triviallövskog björk	High
	Triviallövskog med ädellövinslag	High
	Triviallövskog övrigt	High
	Ädellövskog bok	High
	Ädellövskog ek	High
	Ädellövskog övrigt	High
Naturtypskartan – NNK (Naturvårdsverket)	9080 - Lövsumpskog	High
	9110 - Näringsfattig bokskog	High
	9130 - Näringsrik bokskog	High
	9160 - Näringsrik ekskog	High
	9161 - Näringsrik ekskog - Ek-avenbokskog	High
	9162 - Näringsrik ekskog - Ek-hassellund	High
	9170 - Torr ekskog	High
	9180 - Ädellövskog i branter	High
	9190 - Näringsfattig ekskog	High
	9750 - Svämlövskog (91E0)	High
	9760 - Svämädellövskog (91F0)	High
	9820 - Obestämd ädellövskog (9020, 9850, 9860)	High

For details on the specific data sources that were used to build the reference dataset, see **Table 1** below.

The reference datasets in **Table 1** were compiled and then prepared for the ML modelling. We used the Random Forest algorithm which is built on an assemblage of multiple decision trees. However, decision trees are sensitive to class imbalance and the classifier shows better results when it is trained with datasets where the number of data points per class is more balanced. Therefore, a balancing procedure was required.

The landscape was sampled on target areas (labelled class 1 = protected DF) against all other types of forests (labelled class 0). Finally, to avoid the influence of edge effects, sample points within a 100 meter buffer from boundary lines were subtracted. Note that forest masking was performed throughout all sampling implementations. Consequently, sampling of non-forest was avoided. Reasonable balance was attained once the difference between number of sampled points per class reached under 25%.

Once the sampling scheme was implemented and the geodata sets with the modelling variables ready (this last step will be explained in detail in the next section), variable's values were extracted for each sampled-labelled point. Then our reference dataset was ready. Depending on the size of the county, between 400,000 and 1,000,000 points were sampled per county.

Modelling data

The characteristics that make our target forest type valuable are influenced by factors other than the tree species assemblage alone, e.g. tree age, biomass, land use, terrain etc. Therefore, we explored and mined all the available open-source data that could help us understand the interactions between these factors as a predictor of DFs with HCV.

Multispectral satellite imagery with its different bands and their correlation to vegetation signatures allowed us to infer vegetation cover and the state of the vegetation. Topographic factors such as elevation and slope influence soil profiles and water fluxes, and therefore they could help us discern where our target forests thrive. Several high-quality studies from the Swedish Forest Agency, SFA and the Swedish University of Agricultural Sciences, SLU have been carried out to gather data on the status of Swedish forests. Such data on forest status could also be of value to decipher the presence of our targeted forests.

For our methodology implementation, we also used data on general dendrometry on biomass, tree height and soil moisture from the SFA, and the tree volume of various deciduous tree species from the SLU. Data from the Swedish National Land Cover Database (NMD) was used to make a visual investigation of the ML resulting rasters. However, the NMD data was not used to train our models. In **Table 2** below, the specific data sources.

Table 2. Predictors used for modelling.

Data type		Source	Unit	Notes
Satellite data	Sentinel 2 – band 2 - 8 (8 & 8a)	Copernicus hub	DN	Sentinel 2 L2A* Resolution between 10 – 20 m Captured 2023-2024, cloud free scenes from late autumn. Accessed via: https://dataspace.copernicus.eu/browser/
Topographic data	Elevation Slope	The Swedish Geographical Survey	meters degrees	Point cloud data were accessed via the Swedish Geographical Survey’s web portal, Geotorget: https://geotorget.lantmateriet.se/ The data were processed to a DEM with 10 meter resolution.
Dendrometry	BEECHVOL_XX_P_15 BIRCHVOL_XX_P_15 DECIDUOUSVOL_XX_P_15 OAKVOL_XX_P_15 Biomass Tree height	The Swedish University of Agricultural Sciences (SLU) The Swedish Forest Agency	m3/ha Ton/ha dm	Accessed via SLU’s open data service: https://gis.slu.se/data/slu_forest_map/ Resolution 10 m Accessed via the Swedish Forest Agency’s open data service: https://www.skogsstyrelsen.se/e-tjanster-
Soil type	Soil moisture		wtd	

* Level-2A provides Bottom of Atmosphere (BOA) reflectance, which are corrected for atmospheric effects. Level-2A products are useful for applications requiring accurate surface reflectance data, such as vegetation monitoring and land cover classification.

The Sentinel 2 - L2A band sets and geospatial data processing was carried out using Python 3.7, QGIS and ArcGIS.

Machine learning approach

The accurate detection of DFs with HCV using remote sensing is a complex task. As mentioned before, the particular characteristics that make this and other forest types valuable for nature conservation are influenced by multiple factors. Methods that can capture problems on relatively large areas are needed for an analysis on the spatial scale of this study. Supervised machine learning (ML) is a suitable and successfully tested approach in studies of this spatial scale. It is a method to investigate the relevance of many factors from a variety of data sources in a systematic manner. Several ML-algorithms have already been successfully used to predict forest types, but so far few have been made with the aim of detecting DFs with HCV.

We used the “scikit-learn” library in Python to implement the ML approach and investigated the relationship between variables with a cross-correlation and principal component analysis (PCA) – examples can be ob-

served in **Figure A - C**. Afterwards, we used all derived variables as predictor information for the classification with supervised machine learning.

All labelled data (class 1 = protected DF and 0 = other forest) was used to train several machine learning algorithms. With the library “lazy predict” many basic ML models were built to determine which models would work better without any parameter tuning. This included ML algorithms like Random Forest classifier, Ada-Boost classifier and logistic regression. After a thorough comparison of these algorithms, the Random Forest

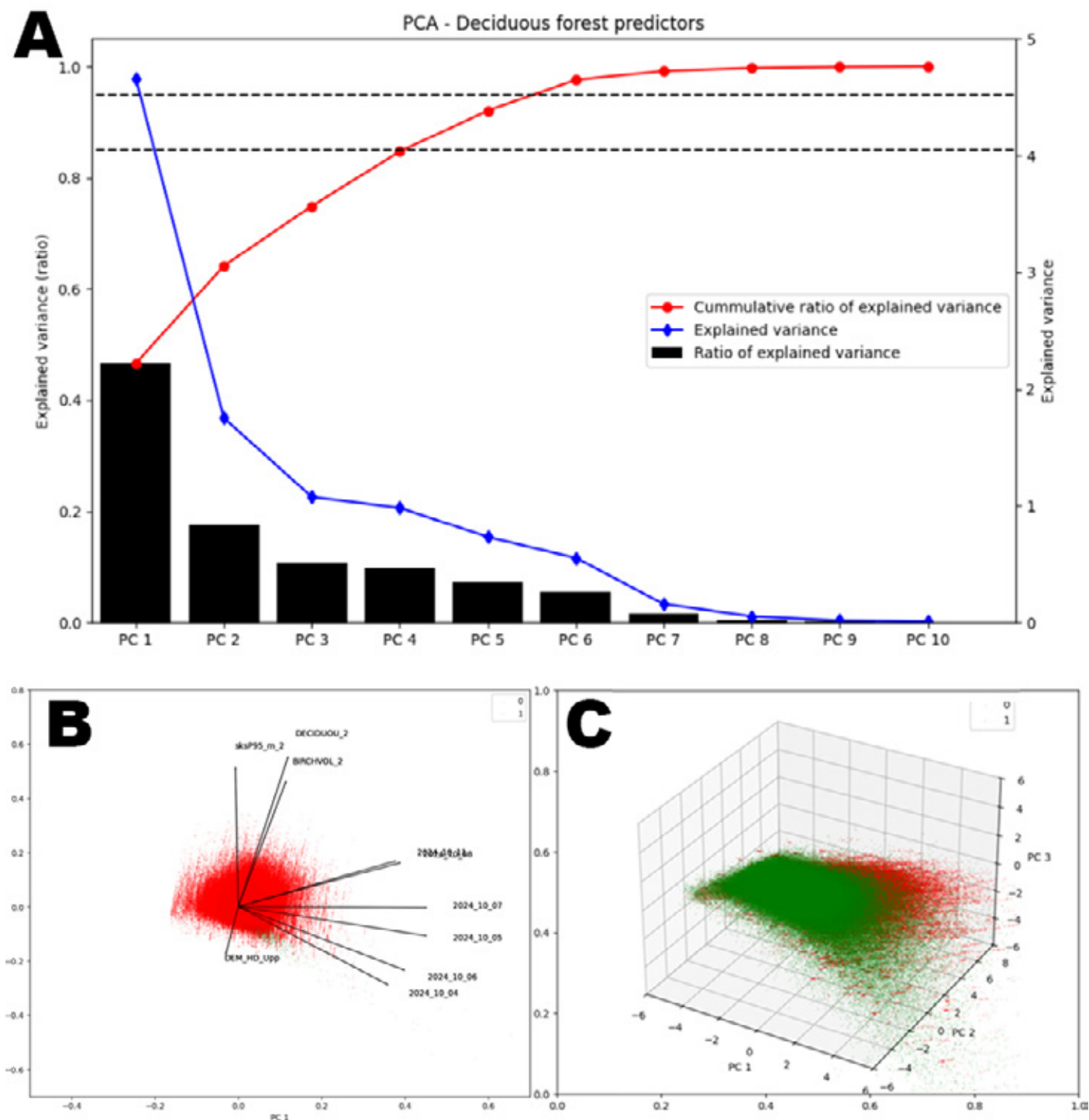


Figure 2 A - C. In A, an example of a screen plot from the PCA on data for one of the 13 county sections that were mapped in this study. In B and C, the 2D (left) and 3D (right) projection of the multidimensional vectors (predictors) and the achieved segmentation of sampled points by class (prediction).

classifier came out as the best choice. Based on this result, we chose to use this classifier.

The random forest classifier was trained (0.7 train: 0.3 test partitioned) and its performance was evaluated. To optimize the classifier, we tuned the hyper parameters and then trained the algorithm again. We used a random grid search with cross validation to define a grid of hyperparameter ranges, and randomly sampled from the grid, performing K-Fold cross validation with each combination of values. The best parameter combination was used in the final model, see **Appendix 1 - 4** for examples.

The performance of the predictions was assessed based on the testing partition of the dataset through the following metrics: overall accuracy, balanced accuracy, F1 score, recall, precision, R2, RMSE and Kappa-values. To evaluate the predictor's relevance in the results, the importance of each predictor was also calculated.

For a full description of the workflow of the study, see **Figure 3** below.

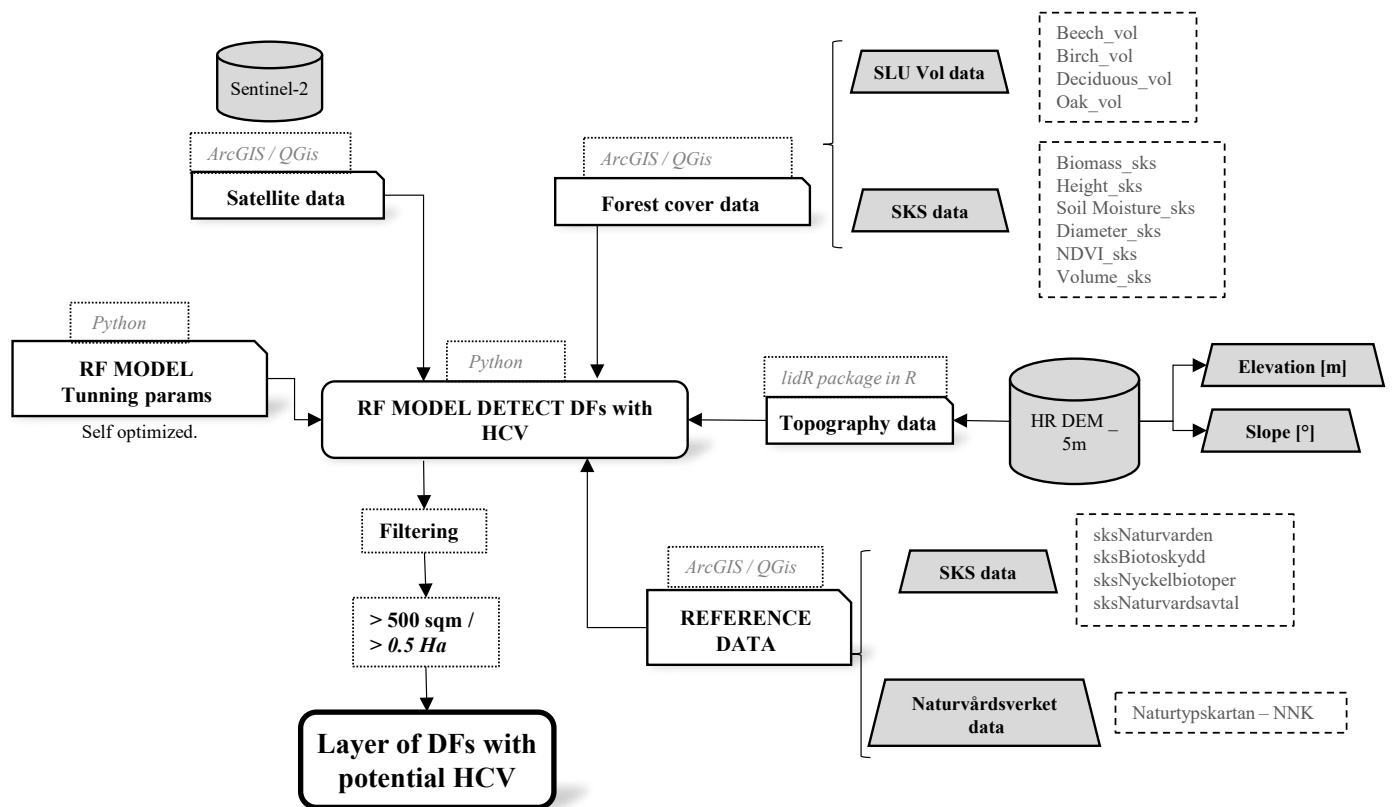


Figure 3. Workflow diagram that illustrates the general approach of our study.

RESULTS

Once the accuracy metrics were calculated, we could conclude that we achieved good accuracy values and acceptable to good error ranges (Accuracy > 0.9; Kappa > 0.8; RMSE < 0.3). According to our results, approximately a little more than 289,000 hectares of DFs with potential HCV lay in the studied area. From those, a mere 18.5% (53,519 hectares) lay inside protected areas.

Skogsmonitor would like to emphasize that it is important for everyone who uses the new deciduous forest layer to understand that this is a prediction, or forecast, of where in the landscape there could be DFs with potential HCV.

The degree of accuracy determines how valuable these types of predictions are for nature conservation work,

environmental monitoring, and in terms of streamlining the prioritization of what should be inventoried in the field.

There is no scientific consensus on what constitutes “good” accuracy in forest classification studies using machine learning. However, similar studies in this field vary from acceptable and valuable results above 70% to excellent accuracy, which usually exceeds 90–95%. High accuracy depends on factors such as the complexity of the classification task and the quality and type of data used.

Our results have a high degree of accuracy. Although they are usually correct, we have also seen examples of accumulations of standing dead conifers being detected as defoliated large deciduous trees in some cases.

For a comparison between known areas with HCV DFs and predicted presense of such forests, see **Figure 4 A-B**.

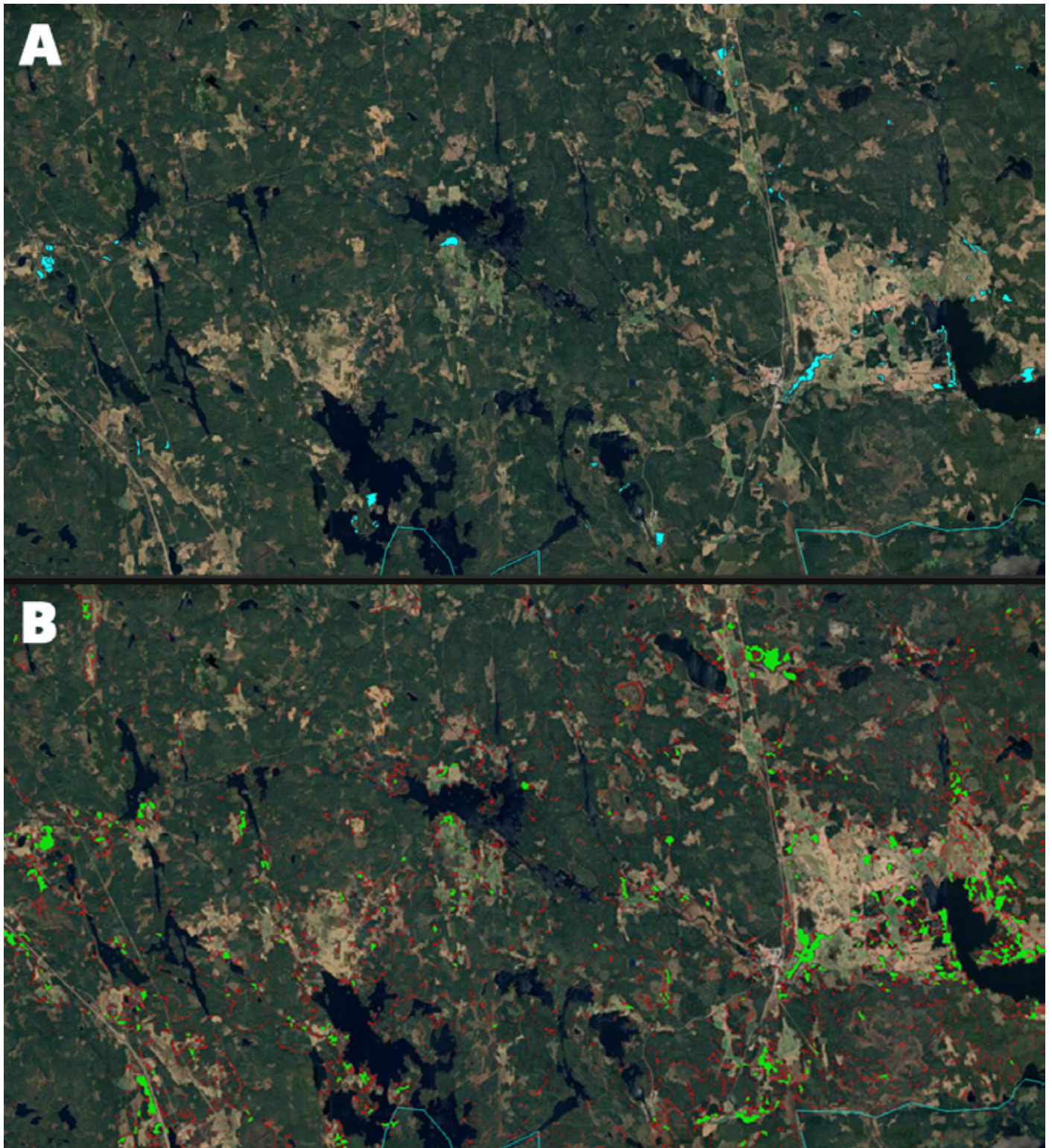


Figure 4 A - B. In A, formally protected deciduous forests in light blue and in B, the ML predictions in light green. The area covers a section of southeast Örebro county.

VALIDATION

As part of the validation process, two complementary analyses were conducted:

- An overlay analysis of predicted rasters with stand data from the Swedish Forest Inventory (SFI).
- An overlay analysis of predicted rasters with land cover data derived from Nationella Marktäckedata 2023 – Basskikt NMD2023 version 2.0, released in July 2025.

The analysis using SFI data revealed that 89% of the predicted DF–HCV areas were located within stands characterized by a basal area-weighted average stand age over 40 years. This finding demonstrates a strong coherence between the predicted patterns and field-based observations made by the SFI. The corresponding results are presented in **Figure 5**. Furthermore, about 82% of the plots found inside of predicted DF–HCV areas were dominated by deciduous trees. Here, note that the SFI is a plot survey, and hence plots may have ended up in parts of forests that are not representative for the forest as a whole.

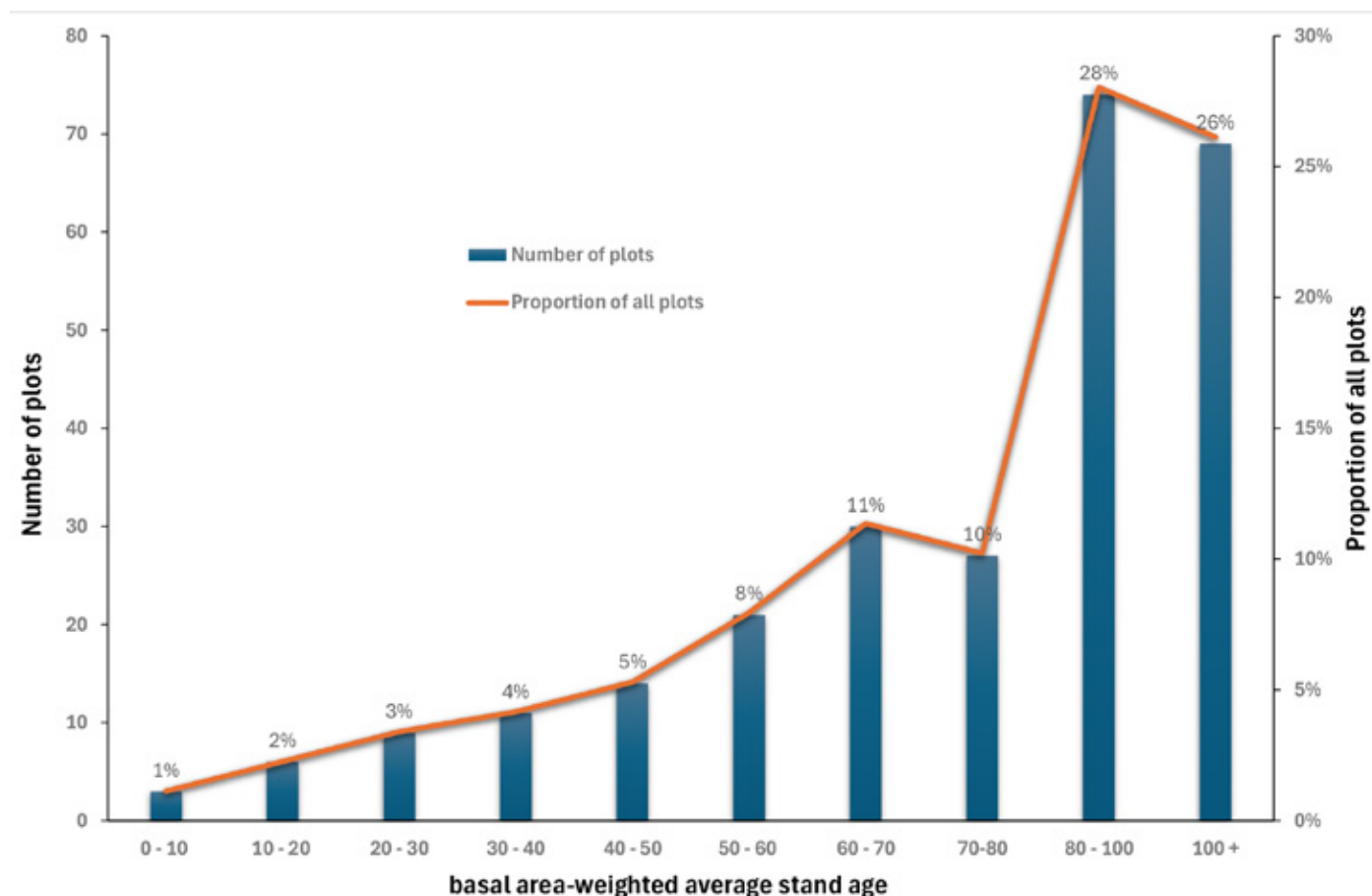


Figure 5. The number of SFI-plots that fell inside of the mapping of DF-HCV (left y-axis) as a function of forest age presented as 10-year intervals from zero years to 100+ years. The proportion of plots within each age interval is shown on the orange line and the proportions correspond to the right y-axis.

The results indicate a very high correspondence between the predicted deciduous forest (DF) areas with high conservation value (HCV) and the NMD dataset. Approximately 90% of the predicted DF–HCV areas overlapped with the corresponding land cover categories in the NMD2023 mapping. The detailed results of this comparison are presented in **Table 3**.

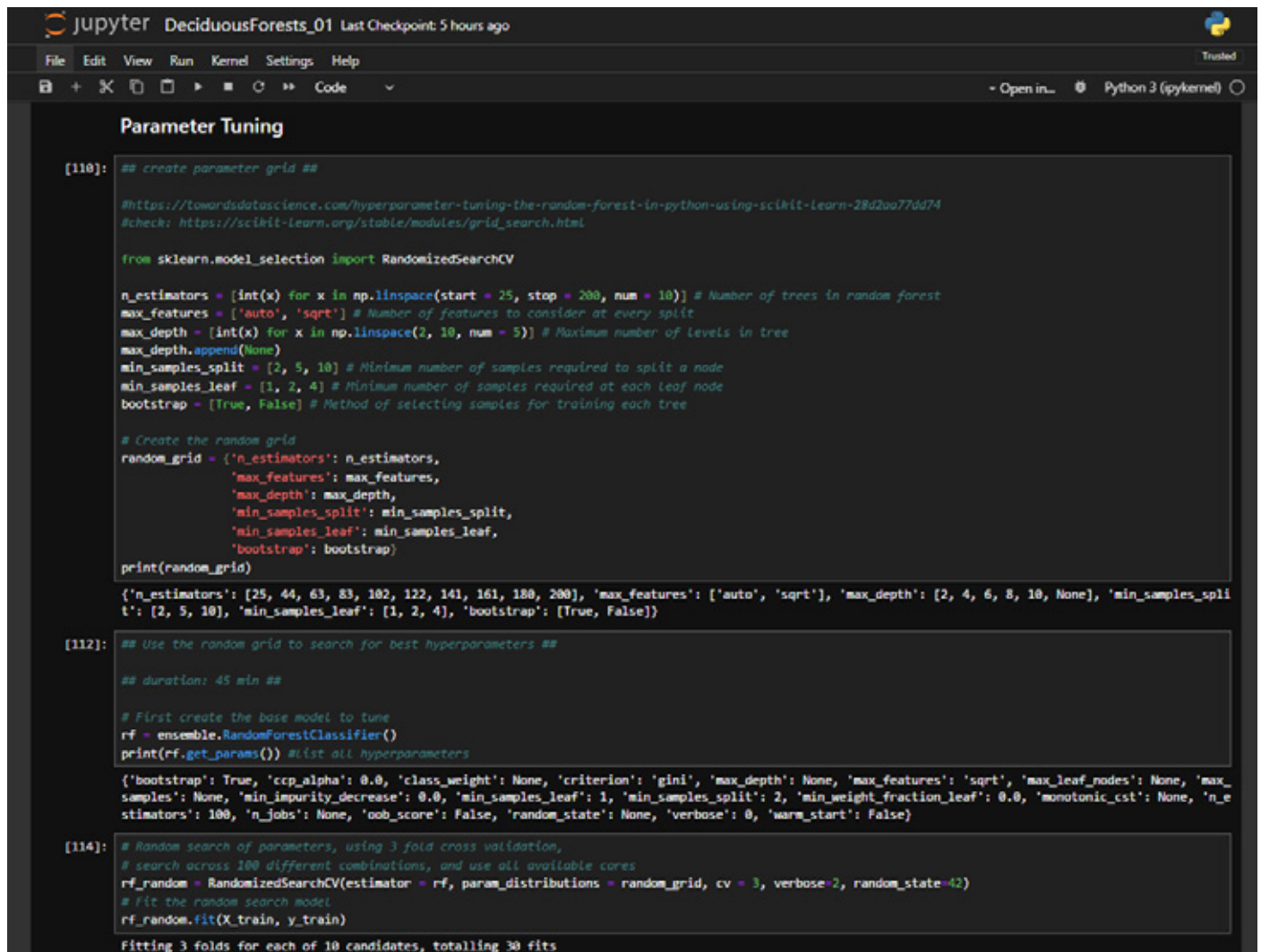
Overall, the analyses suggest that approximately only 284,000 hectares of deciduous forest with potential natural values remain within the study area. Of this total, a mere 18.8% (equivalent to 53,371 hectares) are situated within formally protected areas.

Tabel 3. Predicted deciduous forests with high conservation values over NMD 2023 data released 2025.

NMD Klass	pixel count	Proportion
Åkermark	49958	0,2%
Byggnad	10641	0,0%
Anlagd mark, ej byggnad eller väg/järnväg	1589	0,0%
Väg eller järnväg	77637	0,3%
Torvtäkt	12	0,0%
Inlandsvatten	19811	0,1%
Hav	731	0,0%
Tallskog på fastmark	696274	2,5%
Granskog på fastmark	791639	2,8%
Barrblandskog på fastmark	530023	1,9%
Lövblandad barrskog på fastmark*	1489288	5,2%
Triviallövskog på fastmark*	6651469	23,4%
Ädellövskog på fastmark*	11484792	40,4%
Triviallövskog med ädellövinslag på fastmark*	3569780	12,6%
Temporärt ej skog på fastmark	259032	0,9%
Tallskog på våtmark	104533	0,4%
Granskog på våtmark	59180	0,2%
Barrblandskog på våtmark	48983	0,2%
Lövblandad barrskog på våtmark*	153600	0,5%
Triviallövskog på våtmark*	1951093	6,9%
Ädellövskog på våtmark*	15171	0,1%
Triviallövskog med ädellövinslag på våtmark*	42620	0,2%
Temporärt ej skog på våtmark	19671	0,1%
Öppen våtmark (underindelning saknas)	26404	0,1%
Buskmyr	2514	0,0%
Ristuvemyr	34	0,0%
Fastmattemyr, mager	942	0,0%
Fastmattemyr, frodig	2088	0,0%
Sumpkärr	898	0,0%
Mjukmattemyr	34	0,0%
Lösbottenmyr	97	0,0%
Våtmark med buskar	3942	0,0%
Risdominerad våtmark	25	0,0%
Gräsdominerad våtmark, mager	1520	0,0%
Gräsdominerad våtmark, frodvuxen	2736	0,0%
Gräsdominerad våtmark, högvuxen	2188	0,0%
Mossdominerad våtmark	23	0,0%
Våtmark utan växttäck	204	0,0%
Öppen fastmark utan vegetation (ej glaciär eller varaktigt snöfält)	2465	0,0%
Torr buskdominerad mark	13152	0,0%
Frisk buskdominerad mark	14701	0,1%
Frisk-fuktig buskdominerad mark	9500	0,0%
Torr risdominerad mark	25475	0,1%
Frisk risdominerad mark	34044	0,1%
Frisk-fuktig risdominerad mark	21522	0,1%
Torr gräsdominerad mark	47817	0,2%
Frisk gräsdominerad mark	100375	0,4%
Frisk-fuktig gräsdominerad mark	61497	0,2%

* relevant categories

APPENDIX



```
[110]: ## create parameter grid ##

#https://towardsdatascience.com/hyperparameter-tuning-the-random-forest-in-python-using-scikit-learn-28d2aa77dd74
#check: https://scikit-learn.org/stable/modules/grid_search.html

from sklearn.model_selection import RandomizedSearchCV

n_estimators = [int(x) for x in np.linspace(start = 25, stop = 200, num = 10)] # Number of trees in random forest
max_features = ['auto', 'sqrt'] # Number of features to consider at every split
max_depth = [int(x) for x in np.linspace(2, 10, num = 5)] # Maximum number of levels in tree
max_depth.append(None)
min_samples_split = [2, 5, 10] # Minimum number of samples required to split a node
min_samples_leaf = [1, 2, 4] # Minimum number of samples required at each leaf node
bootstrap = [True, False] # Method of selecting samples for training each tree

# Create the random grid
random_grid = {'n_estimators': n_estimators,
               'max_features': max_features,
               'max_depth': max_depth,
               'min_samples_split': min_samples_split,
               'min_samples_leaf': min_samples_leaf,
               'bootstrap': bootstrap}

print(random_grid)

{'n_estimators': [25, 44, 63, 83, 102, 122, 141, 161, 180, 200], 'max_features': ['auto', 'sqrt'], 'max_depth': [2, 4, 6, 8, 10, None], 'min_samples_split': [2, 5, 10], 'min_samples_leaf': [1, 2, 4], 'bootstrap': [True, False]}

[112]: ## Use the random grid to search for best hyperparameters ##

## duration: 45 min ##

# First create the base model to tune
rf = ensemble.RandomForestClassifier()
print(rf.get_params()) #list all hyperparameters

{'bootstrap': True, 'ccp_alpha': 0.0, 'class_weight': None, 'criterion': 'gini', 'max_depth': None, 'max_features': 'sqrt', 'max_leaf_nodes': None, 'max_samples': None, 'min_impurity_decrease': 0.0, 'min_samples_leaf': 1, 'min_samples_split': 2, 'min_weight_fraction_leaf': 0.0, 'monotonic_cst': None, 'n_estimators': 100, 'n_jobs': None, 'oob_score': False, 'random_state': None, 'verbose': 0, 'warm_start': False}

[114]: # Random search of parameters, using 3 fold cross validation,
# search across 100 different combinations, and use all available cores
rf_random = RandomizedSearchCV(estimator = rf, param_distributions = random_grid, cv = 3, verbose=2, random_state=42)
# Fit the random search model
rf_random.fit(X_train, y_train)

Fitting 3 folds for each of 10 candidates, totalling 30 fits
```

Appendix 1. Screenshot showing the script run to find out the best parameter's tuning of the RF model.

```
[114]: RandomizedSearchCV
>
> best_estimator_: RandomForestClassifier
- RandomForestClassifier
RandomForestClassifier(bootstrap=False, min_samples_leaf=4, n_estimators=122)

print best parameter combination

[116]: rf_random.best_params_

[116]: {'n_estimators': 122,
'min_samples_split': 2,
'min_samples_leaf': 4,
'max_features': 'sqrt',
'max_depth': None,
'bootstrap': False}

[118]: def evaluate(model, X_test, y_test):
    predictions = model.predict(X_test)
    errors = abs(predictions - y_test)
    accuracy = 100 - 100*np.mean(errors)
    print('Model Performance:')
    print('Average Error: {:.4f}'.format(np.mean(errors)))
    print('Accuracy = {:.2f}%'.format(accuracy))

    return accuracy

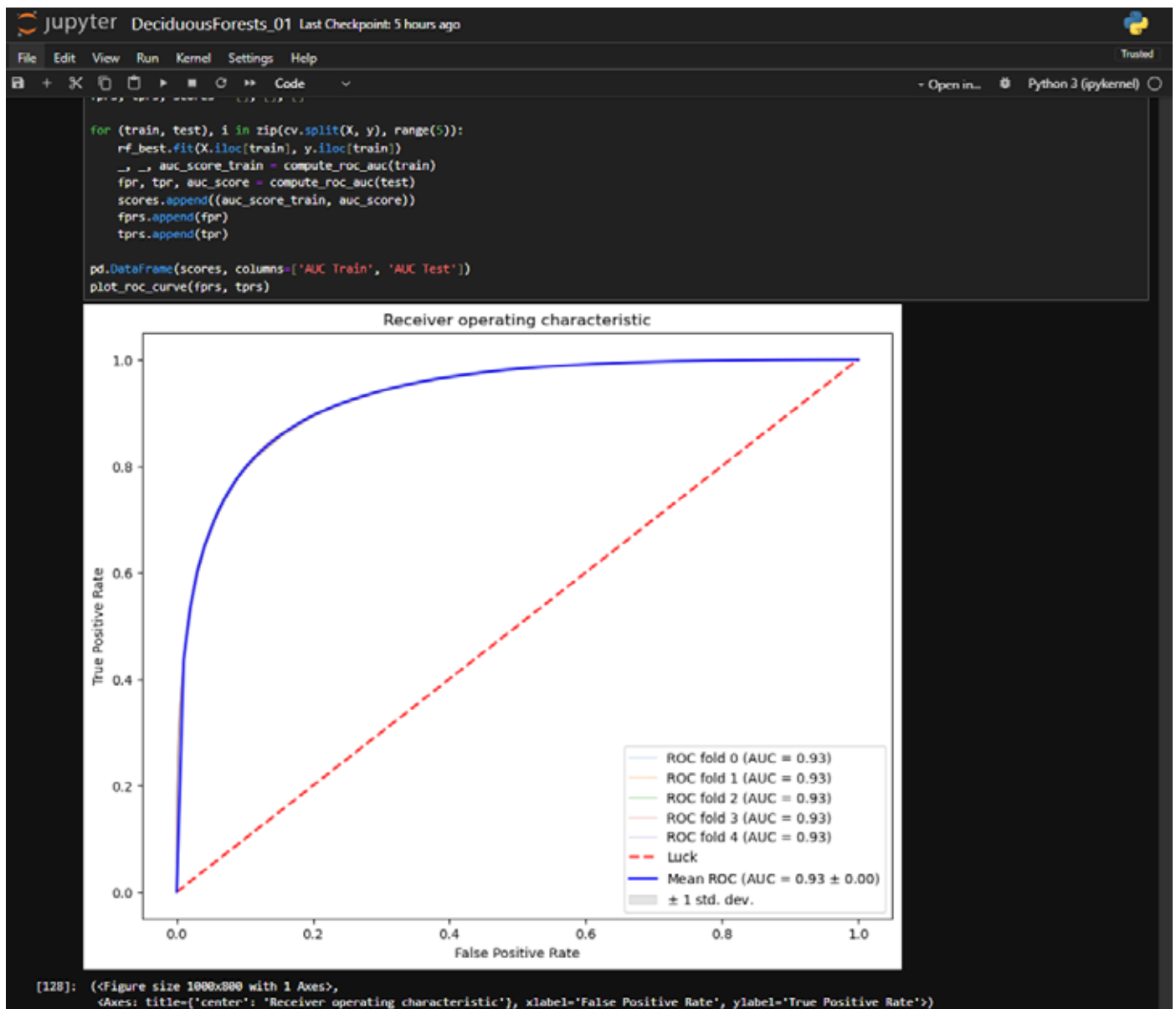
base_model = ensemble.RandomForestClassifier(n_estimators = 122, random_state = 42)
base_model.fit(X_train, y_train)
base_accuracy = evaluate(base_model, X_test, y_test)

best_random = rf_random.best_estimator_
random_accuracy = evaluate(best_random, X_test, y_test)

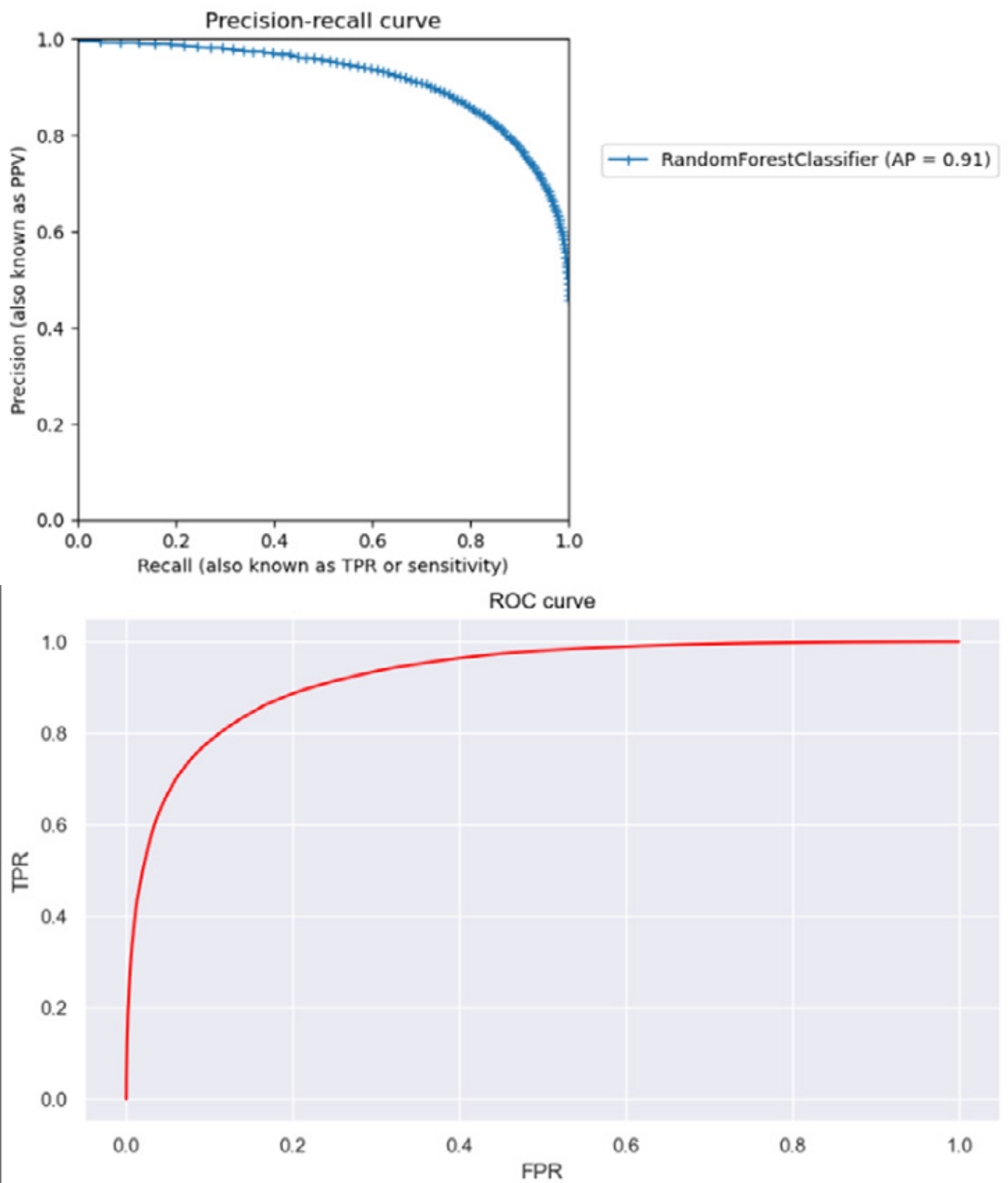
print('Improvement of {:.2f}%'.format( 100 * (random_accuracy - base_accuracy) / base_accuracy))

Model Performance:
Average Error: 0.1508.
Accuracy = 84.92%.
Model Performance:
Average Error: 0.1500.
Accuracy = 85.00%.
```

Appendix 2. Screenshot showing the ML results of best classifier, best parameter combination and its performance.



Appendix 3. Screenshot showing the ROC curve of the RF model results for one county section.



Appendix 4. Precision-recall curve and ROC curve.



www.skyddaskogen.se

CONTACT

US



skogsmonitor@skyddaskogen.se

www.skogsmonitor.se

