



UNIVERSITÀ DEGLI STUDI DI SALERNO

IoT DATA ANALYTICS

MEASURING THE PERFORMANCES OF ONLINE ANALYSIS FOR
FOREST COVER TYPE CLASSIFICATION

Professors: Giuseppe Polese, Genoveffa Tortora

Project coordinator: Stefano Cirillo

Student: Michelangelo Esposito – 0522500982

Academic year 2020/2021

INDEX

1.	INTRODUCTION.....	3
2.	STATE OF THE ART FOR THE FOREST COVER TYPE PROBLEM.....	5
3.	CONTEXT OF THE PROJECT	6
4.	FOREST COVER TYPE DATASET	7
4.1	OVERVIEW	7
4.2	PREVIOUS DATASET USAGES.....	8
4.3	DATASET ANALYSIS	9
5.	EMPLOYED MODELS.....	15
5.1	DECISION TREE	15
5.2	RANDOM FOREST	16
5.3	HOEFFDING TREE	17
5.4	ADAPTIVE RANDOM FOREST	18
5.5	NAÏVE BAYES.....	18
6.	ONLINE DATASET ANALYSIS	19
6.1	MOA.....	19
6.2	SCIKIT MULTIFLOW.....	23
7.	BATCH LEARNING APPROACH.....	26
7.1	DATASET SPLITTING VERSUS CROSS VALIDATION.....	27
7.2	OVERFITTING AND UNDERFITTING EVALUATION.....	31
7.3	DATA BALANCING	33
7.4	CONFUSION MATRIX	34
8.	CONCLUSIONS	35
	REFERENCES	37

1. INTRODUCTION

Most of the times, machine learning is associated with batch data, meaning that a dataset is obtained from a source of some sort, and then it is used to train a classifier. As a consequence, most machine learning and data mining approaches assume that the elements of the dataset are independent, identically distributed and generated from a stationary distribution. This process is generally adaptable to most “static” circumstances and there are a variety of machine learning models and networks that employ it.

Sometimes, however, the data may not be available all at once from the start, or the hardware on which these machine learning processes run may not satisfy the requirements to run traditional batch learning. Nowadays, we are faced with a tremendous amount of distributed data that can be generated from the ever-increasing number of smart devices. In an IoT scenario, for example, dynamic streams of data are being constantly emitted from countless sensors, and resource-limited devices; in most cases, this data is transient and may be stored for a brief period of time before getting discarded. In this small fraction of time, we are often required to perform analysis operations on the data we receive; therefore, there is the need to adapt the learning process to a workflow that is able to satisfy limited hardware requirements and to dynamically adapt to the continuous flow of data.

Online learning is a method of machine learning in which data becomes available sequentially and is used to update the predictor for future data at each step. It can be data efficient since once an entry has been processed, it is no longer required; therefore, storing previously seen data is not necessary. This approach is also highly adaptable to evolving circumstances since it makes no assumptions about the distribution of the data: as the distribution morphs or drifts, the model can adapt on-the-fly to keep pace with trends in real-time. In order to do something similar with traditional offline learning, one would have to create a sliding window for the data and retrain it every time.

There are various reasons to use online learning over batch learning, among these:

- The whole dataset does not fit in memory.
- The distribution of the data is expected to morph over time.
- The data itself is a function of time (i.e., stock prices).

A very noticeable difference between online and batch learning resides in the training and testing phases of the classifier. Traditionally, there are various ways for evaluating a batch classifier, such as reserving a percentage of the data for testing, or performing k-fold cross validation, a technique by which the data is split into k groups and for each group an evaluation is performed by considering that group as the testing one and the others as the training ones.

In online learning there is no way of reserving a portion of the data for testing purposes, given the volatile nature of the stream, therefore a new approach becomes mandatory. In data stream classification, the most used evaluation technique is the prequential one, where instances are first used to test, and then to train the model.

Finally, something to be really aware of when performing classification over a data stream is concept drift; this is the condition due to which the statistical properties of the target variable, which the model is trying to predict, change over time in unforeseen ways. While this represents a problem in batch learning as well, it could pose a significant challenge in online learning since we do not even know the distribution of data to begin with.

In this project we focus on analysing the performances of machine learning classifiers on a dataset containing information about the forest cover type problem. Given a set of cartographic variables and seven different forest cover type classes, our goal is to train different classifiers, in both online and batch fashion, in order to achieve the best performance. In the next sections, we will discuss the state of the art for the problem in consideration and we will analyse the composition of the dataset in more detail, as well as our approach for the prediction.

2. STATE OF THE ART FOR THE FOREST COVER TYPE PROBLEM

The forest cover type problem concerns the ability to correctly classify the cover type (or the species of the trees) based on some kind of environmental variables (i.e., cartographic variables or images). This section will provide a summary for the literature analysis we performed for this problem.

First of all let us consider the approach discussed in [1], where starting from cartographic variables, the authors employ ANNs (Artificial Neural Networks) to classify the 11,340 instances of a dataset into seven different classes/cover types. They tuned the hyperparameters of their neural network, including the number of independent variables, the number of nodes in the network and the learning (momentum rate, in order to achieve an accuracy of 70.58%

The second study taken into consideration [2] was performed as an evolution of another study performed on Northern Europe's forests, with the goal of developing techniques for wall-to-wall mapping of forest variables. The authors describe a method based on the KNN (K-Nearest Neighbour) algorithm to improve the estimation and produce wall-to-wall basal area, volume and cover type maps in the context of USDA Forest Services' Forest Inventory and Analysis (FIA) monitoring system. Their dataset was made of variables extracted from 30x30 pixel images extracted from aerial thematic maps of various forest landscapes. According to their findings, the quality of cover type classification results is comparable with the quality obtained using parametric approaches, but with some advantages: most noticeably, the usage of KNN is very easily understandable, since it is solely based on the search techniques for similar units.

In [3], the authors provide a study of the cover type, among other things, in the Wisconsin region of the U.S. The study site contained a diverse mix of cover types typical of the region and the imagery making up the reference dataset was analysed by using both manual and software-assisted techniques. Focusing on the latter, the classification was performed using a variety of supervised and unsupervised classifiers, with a focus on maximum likelihood classification. First of all, three information classes were interpreted from the multispectral images and the accuracy results observed for them were 96% for the level 1 class, 91% for the level 2 class and 90% for the level three class.

An interesting approach used as an alternative to the analysis of traditional cartographic numerical variables is the usage of satellite thematic maps of forests, as discussed in [4]. With their approach, starting from AVHRR (Advanced very-high-resolution

radiometer) data obtained from satellites, the authors were able to apply statistical regression techniques to construct a map of predicted within-AVHRR-pixel percent forest cover type values and to provide a mechanism for selection of the optimal AVHRR dates and spectral channels for selection for cover type classification. They ended up developing separate models for each of the 15 physiographic regions on which they focused on in their paper. The classifications were then performed for each physiographic region on the same AVHRR data. Regarding the results, the mean and standard deviation of the difference variable between FIA and the RPA map estimations of 1.991 state percent forests were found to be 2.02 and 1.48 percent, respectively. The average agreements for the nine TM accuracy tests for forest and non-forests classes were 85% and 91%, based on average sample sizes of 1922 and 1421 TM pixels, respectively. The average of overall agreements was 89%.

3. CONTEXT OF THE PROJECT

The main focus of this project is to employ online machine learning techniques to analyse the performances of various classifiers on a data stream. This stream does not produce real-time data, but it is generated from a csv/arff file in order to use the online approach.

The platform chosen for the initial analysis is MOA, Massive Online Analysis. MOA is an open-source framework designed specifically for data stream mining with concept drift. It is written in Java and developed at the University of Waikato, New Zealand. It allows to build and run experiments of machine learning or data mining on evolving data streams; it includes a set of learners and stream generators that can be used from the Graphical User Interface (GUI), the command-line, and the Java API.

First of all, both the MOA GUI and the Java API were used to perform three parallel classification tasks on the dataset; these results acted as a baseline for any future run.

At this point, a first comparison among the three classifiers was made, in order to analyse their performances and costs. We then decided to repeat the same process in Python, to further elaborate on MOA's efficiency. Finally, a traditional batch approach was employed to additionally assess the potential gap with online learning.

4. FOREST COVER TYPE DATASET

4.1 OVERVIEW

The dataset used for this study is available for free on the UCI Machine Learning Repository website and its employment revolves around predicting forest cover type from cartographic variables only; no images or other kinds of data are present in this dataset. The actual forest cover type for a given observation (30 x 30 meters cell) was determined from US Forest Service (USFS) Region 2 Resource Information System (RIS) data. Independent variables were derived from data originally obtained from US Geological Survey (USGS) and USFS data. The data was not originally scaled and contains binary columns of data for qualitative independent variables, such as wilderness areas and soil types. A total of over 580,000 entries are present in the dataset, each with the following attributes:

- Elevation in meters.
- Aspect in degrees azimuth.
- Slope in degrees.
- Horizontal distance to nearest surface water.
- Vertical distance to nearest surface water.
- Horizontal distance to nearest roadway.
- Hill shade index at 9am, summer solstice.
- Hill shade index at noon, summer solstice.
- Hill shade index at 3am, summer solstice.
- Horizontal distance to nearest wildfire ignition points.
- Wilderness area designation.
- Soil type designation.
- Forest cover type designation (dependent attribute). The seven cover types, with the respective number of entries, are:
 - Spruce/Fir: 211840.
 - Lodgepole Pine: 283301.
 - Ponderosa Pine: 35754.
 - Cottonwood/Willow: 2747.
 - Aspen: 9493.
 - Douglas-fir: 17367.
 - Krummholz: 20510.

The primary reason for the collection of cartographic data pertains to terrain mapping. It is ultimately useful for applying topographic correction to satellite images in remote sensing or as background information in scientific studies. Topographic correction is necessary if, for example, we wish to identify materials on the Earth's surface by deriving empirical spectral signatures, or to compare images taken at different times with different Sun and satellite positions and angles. By applying the corrections, it is possible to transform the satellite-derived reflectance into their true reflectivity or radiance in horizontal conditions.[5]

4.2 PREVIOUS DATASET USAGES

The earliest use of the dataset was in a doctoral dissertation. Jock Blackard and Denis Dean at Colorado State University used Linear Discriminant Analysis to obtain a 58% accuracy and an Artificial Neural Network to achieve 70% accuracy. The paper concludes by saying that while the ANN does have its drawbacks, it is still more viable than traditional methods. In 2004, in *Systems, Man and Cybernetics* an IEEE International Conference publication, Xiaomei Liu, Kevin W. Bowyer of the University of Notre Dame use the forest cover dataset to conjecture a class of problems that are extremely difficult for FFBP NNs, but relatively simple for DTs.

Besides these, an open competition on Kaggle was also found, in which the community challenged themselves with using this dataset and studying the results.

No notable studies involving this dataset in an online fashion were found.

4.3 DATASET ANALYSIS

Before employing classification techniques on this dataset, it is necessary to analyse its composition, the distribution of the various features, its contained values, and scan for any outliers. Note that, for the online classification part, we will not consider any balancing or optimization techniques, since we will be assuming that the entries from the dataset are being streamed in real time from a source of some sort. On the other hand, for the batch learning approach we will try to obtain the best results by applying various pre-processing techniques.

First of all, the data appears to be well formatted and clean; there are no null values, and all features are numeric. As previously said, there are over 580,000 rows in the dataset and the following graph highlights their distribution of the seven cover type classes:

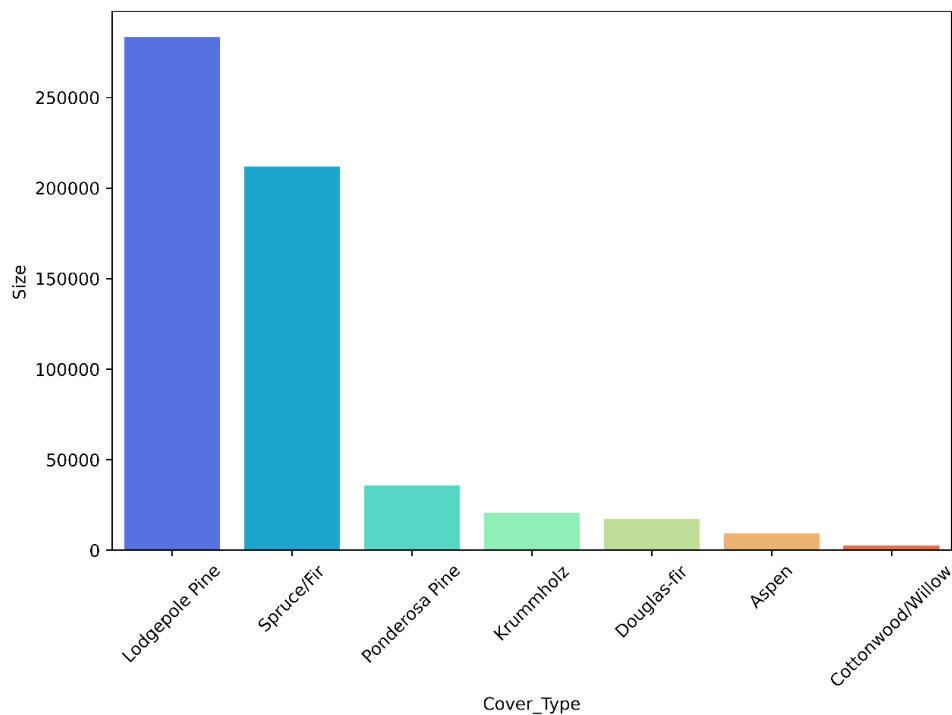


Figure 4.1 Dataset distribution according to the target attribute.

The respective percentages of each class are:

- Lodgepole Pine: 48.76%
- Spruce/Fir: 36.46%
- Ponderosa Pine: 6.15%
- Krummholz: 3.53%
- Douglas-fir: 2.99%
- Aspen: 1.63%
- Cottonwood/Willow: 0.47%

It is immediately noticeable how the balance across the classes is not quite optimal: while the first two classes, Lodgepole Pine and Spruce/Fir contain a relatively similar number of samples, the other five are underrepresented in the dataset; therefore, some data balancing techniques will be taken into consideration during the development of the project.

We can also observe that the dataset contains some negative values in the Vertical_Distance_To_Hidrology field: this indicates that the nearest source of water is below sea level, therefore these negative values are normal and should not cause any concern.

The next thing to look out for is skewness. Skewness is a quantifiable measure of how distorted a data sample is from the normal distribution. In normal distribution, the data is represented graphically in a bell-shaped curve, where the mean (average) and mode (maximum value in the data set) are equal; skewed values may give us an idea of which attributes to remove.

The skewness for a normal distribution is zero, and any symmetric data should have a skewness near zero. Negative values for the skewness indicate data that are skewed left and positive values for the skewness indicate data that are skewed right. By skewed left, it means that the left tail is long relative to the right tail. Similarly, skewed right means that the right tail is long relative to the left tail. After computing the skewness for our dataset, some variables appear to be heavily skewed, hence need to be corrected or transformed later. The next skewness representation shows us that, while most of the attribute present minimal skewness, some of the soil types suffer much more from this problem, especially “Soil_Type15”.

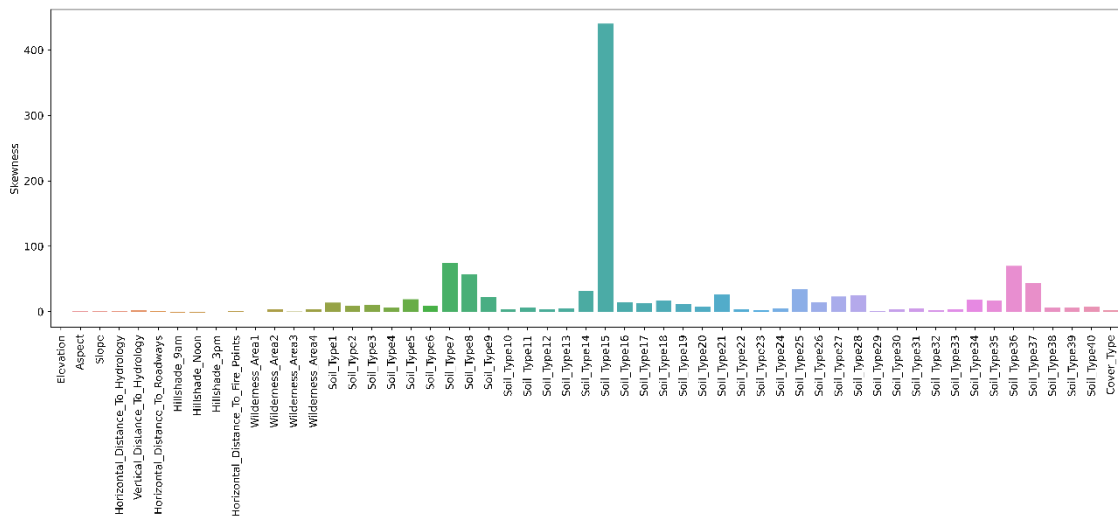
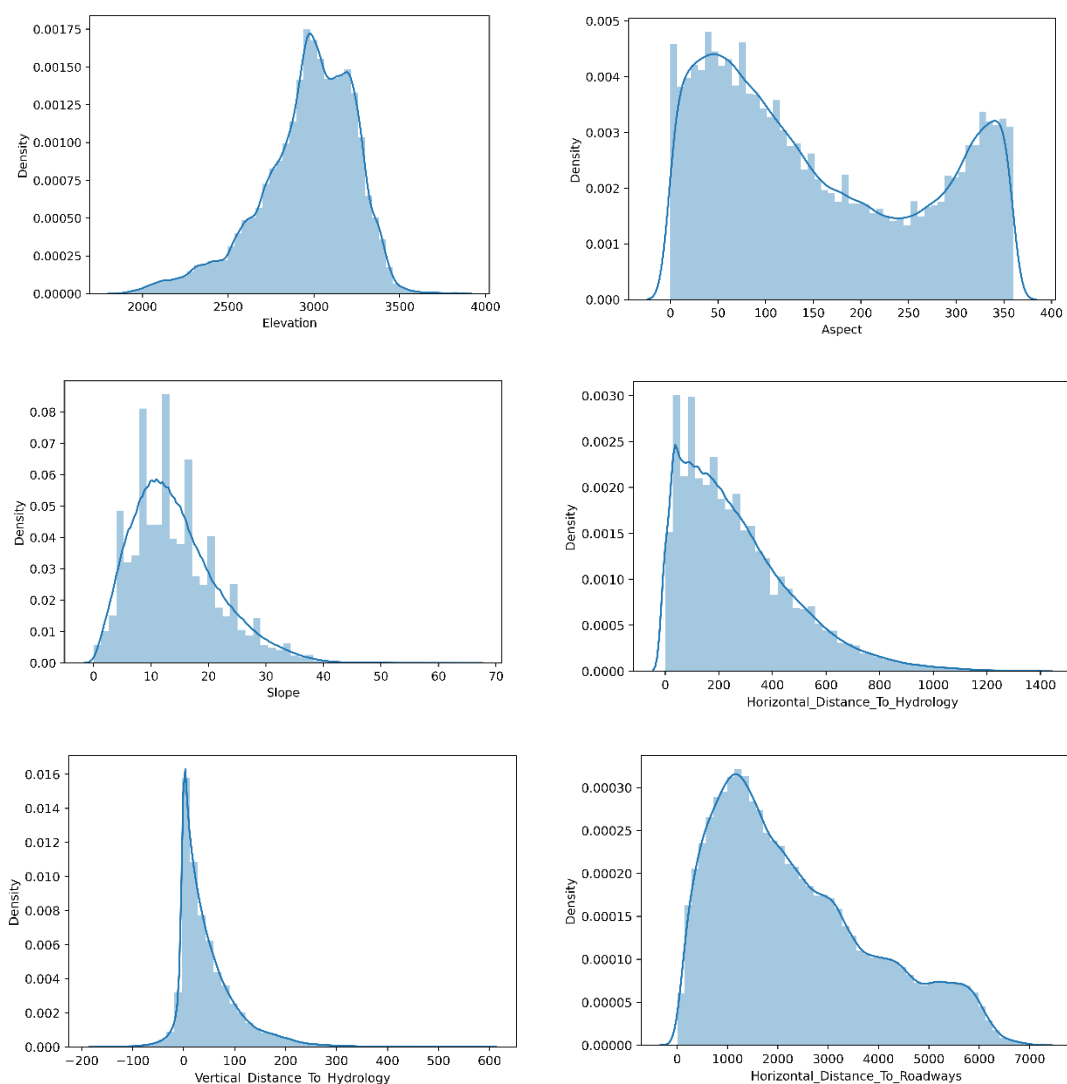


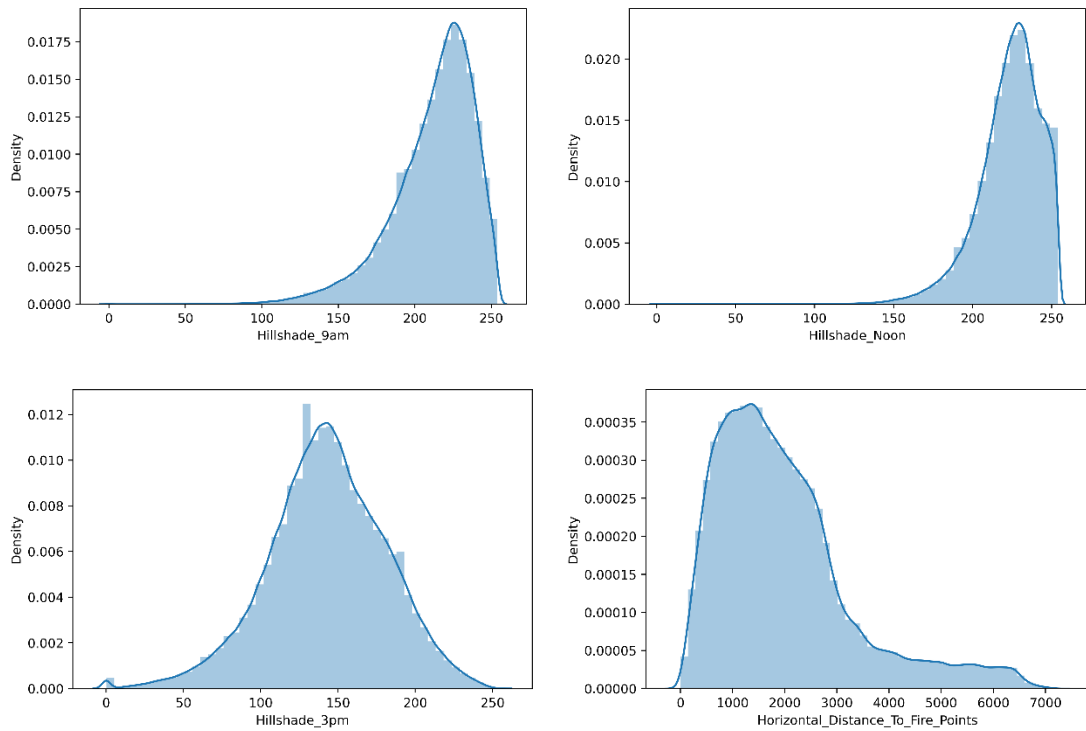
Figure 4.2 Skewness of each feature in the model.

Given the different nature of some of the features in the dataset (i.e., continuous, and binary features), it is worth splitting the dataset into four categories, in order to check the number of value counts within each feature; the considered subsets are:

- Continuous data: this includes elevation, aspect, slope, horizontal and vertical distance to hydrology, horizontal distance to roadways, the three hill shades and the horizontal distance to fire points.
- Binary data: the complete set of binary features, including the 4 wilderness areas and the 40 soil types.
- Wilderness areas data: the 4 wilderness areas alone.
- Soil type data: the 40 soil types alone.

This first set of graphs highlights the distribution of the continuous attributes of the dataset:





Another interesting feature to analyse is the wilderness area of the trees in the dataset, in relation to each cover type:

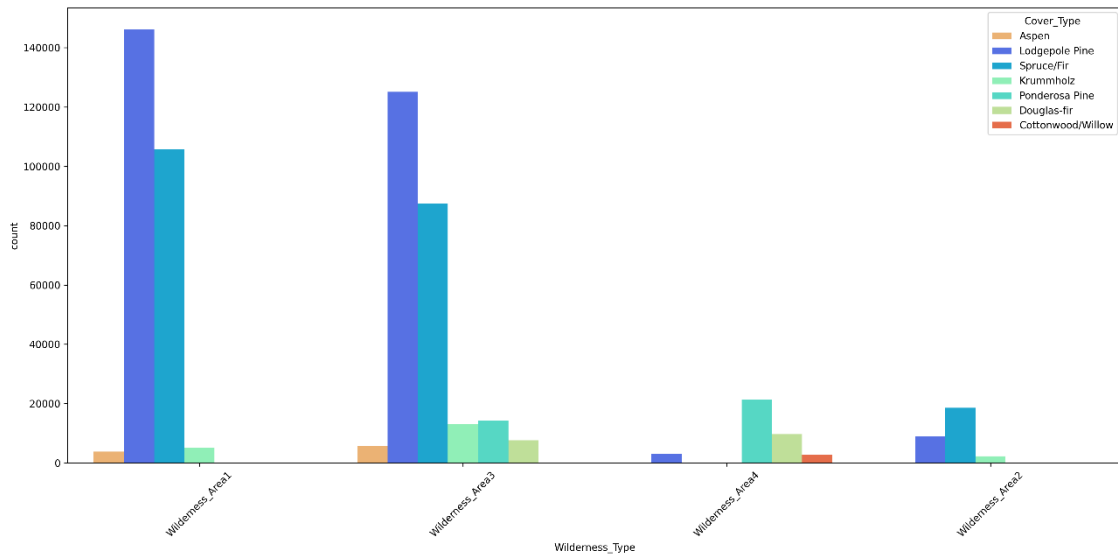


Figure 4.3 Distribution of trees among the four wilderness areas considered.

The distribution of the 40 soil types has also been analysed, finding an unbalanced distribution mostly in favour of the two most present tree types. We do not report the individual graphs for every soil type to avoid cluttering this document; however, the next graph will provide a panoramic view on their distribution:

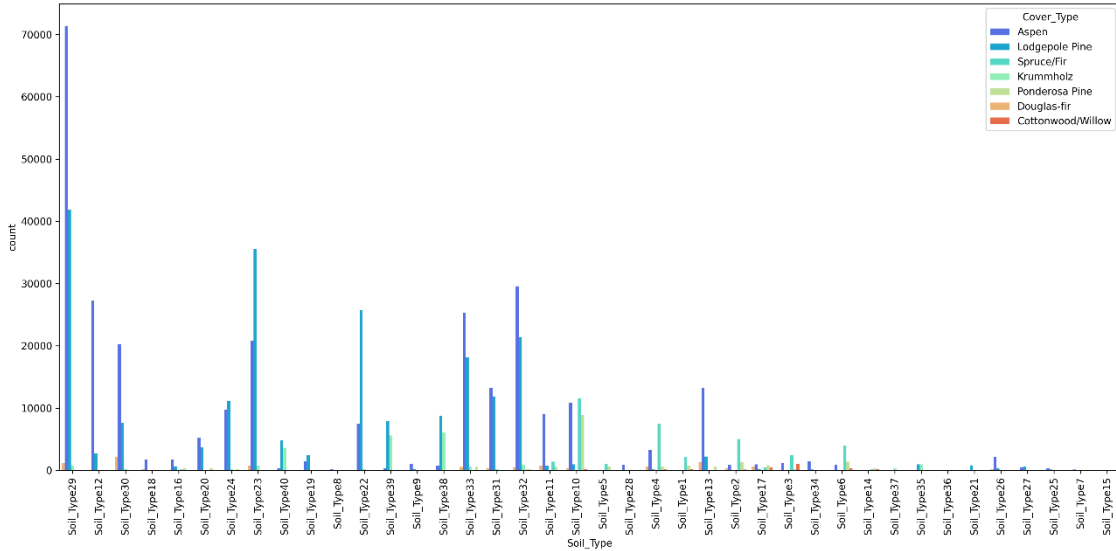


Figure 4.4 Distribution of soil types.

Some soil types are also present in a miniscule number of entries in the dataset; for example, soil_type15 is present in only 3 entries, soil_type36 is present in 119 entries and soil_type8 is present in 179 entries.

A final aspect worth considering is the correlation between the different features. Correlation explains how one or more features are related to each other; it gives us a general idea about the degree of the relationship of two features. Knowing the correlation of the feature of a dataset implies a few things:

- If two variables are closely correlated, then we can deduce one from the other, with a certain degree of precision.
- It plays a vital role in locating the important variables on which other variables depend.
- Proper correlation analysis leads to better understanding of data.

High correlation among values, however, can also be harmful for a classifier's performances. We define multicollinearity as a phenomenon in which one predictor variable can be linearly predicted from the others with a certain degree of accuracy. Multicollinearity does not reduce the overall predictive power of a model as a whole; it only affects calculations regarding individual predictors. As a result, we should remove the features which variables presents a higher-than-normal accuracy value.

The next table is the result of computing correlation values among the continuous attributes of our dataset.

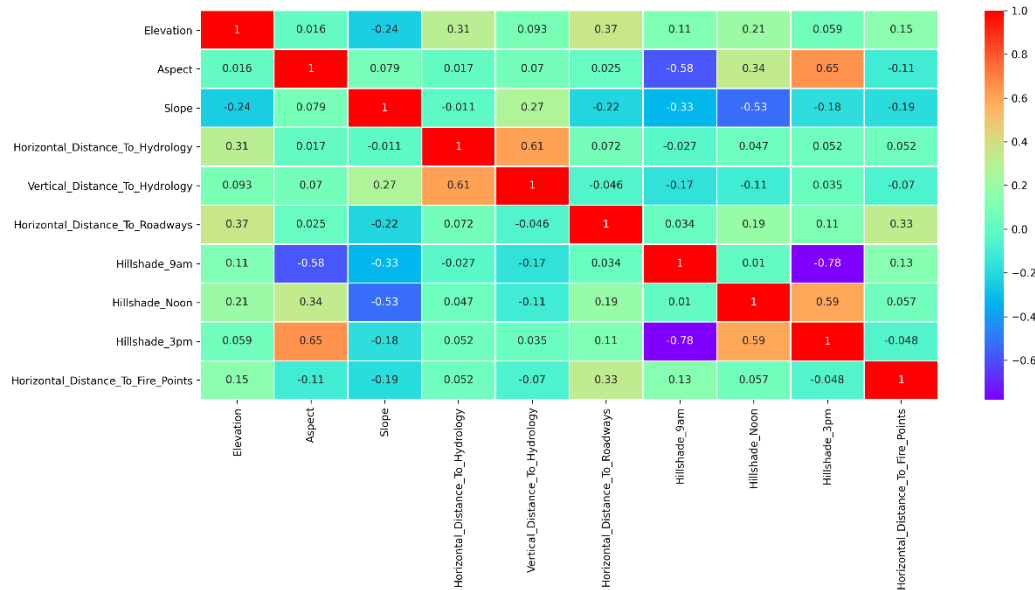


Figure 4.5 Correlation matrix of the dataset's attributes.

The most significant correlation scores appear between in the orange cells: Aspect is strongly correlated to Hillshade_3pm, Horizontal_Distance_To_Hydrology is strongly correlated to Vertical_Distance_To_Hydrology and Hillshade_3pm is strongly correlated to Hillshade_Noon. This can give us a general idea of what feature could be removed later on.

5. EMPLOYED MODELS

The following models have been employed on the forest cover type dataset:

Online learning:

- Adaptive Random Forest.
- Naïve Bayes
- Hoeffding Tree

Batch learning:

- Random Forest
- Decision Tree
- Naïve Bayes

In this section we will briefly explain how each of them works and what are their strengths and weaknesses.

5.1 DECISION TREE

First of all, a **Decision Tree** is a flowchart-like structure in which each internal node represents an evaluating expression on an attribute (in our case, an example of such test could be: “Is the altitude greater than 1500m?”), each branch represents the outcome of the test, and represents a class label (i.e., “Lodgepole Pine”). When they are being built, decision trees are constructed by recursively evaluating different features and using each node for the feature that best splits the data (in order to minimize the depth, and therefore the complexity of the algorithm traversing the tree).

Among decision support tools, decision trees have several advantages:

- They are very easy to understand and interpret.
- They can be combined with other decision techniques.
- They help determine worst, best and expected values for different scenarios.

Their disadvantages include:

- They are often relatively inaccurate. Many predictors perform better with similar data. This can be remedied by replacing a single decision tree with a random forest, as we will see soon.
- They can be unstable, meaning that a small change in the data can lead to a large change in the structure of the optimal decision tree.
- Calculations can get very complex, particularly if many values are uncertain and/or if many outcomes are linked.

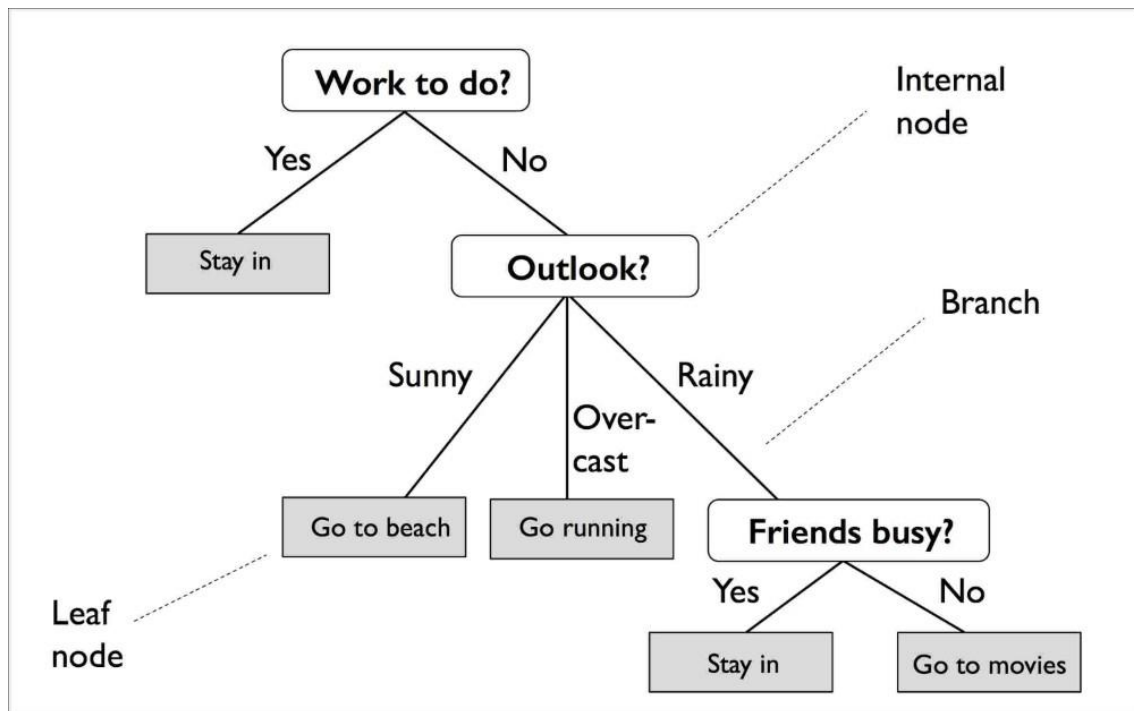


Figure 5.1 Decision Tree example.

5.2 RANDOM FOREST

Similarly, as the name implies, **Random Forests** consist of a large number of individual decision trees that operate as an ensemble: these kinds of models use multiple classifiers to obtain better predictive performance than could be obtained from any of the constituent learning algorithms alone. In the case of random forest, each individual tree in the forest guesses a class prediction and the class with the majority of the votes amongst the trees becomes the final prediction. This approach makes it so that the trees protect each other from their individual errors, as long as the majority does not constantly predict the wrong class; while some trees may be wrong, many other will be right, so as a group the trees are able to move towards the correct prediction.

An extension of the algorithm was developed in 2006 and makes use of bagging and random selection of features in order to construct a collection of decision trees with controlled variance. Decision trees are very sensitive to the data that are trained on, as previously stated; random forest takes advantage of this by allowing each individual tree to randomly sample from the dataset with replacement, resulting in different trees. This process is known as bagging, or bootstrap aggregation.

While generally performing better than singular decision trees, random forests lack their intrinsic interpretability which normally allows developers to confirm that the model

has learned realistic information from the data and allows end-users to have trust and confidence in the decisions made by the model. Furthermore, employing a random forest classifier with a non-trivial depth can significantly slow down the classification process.

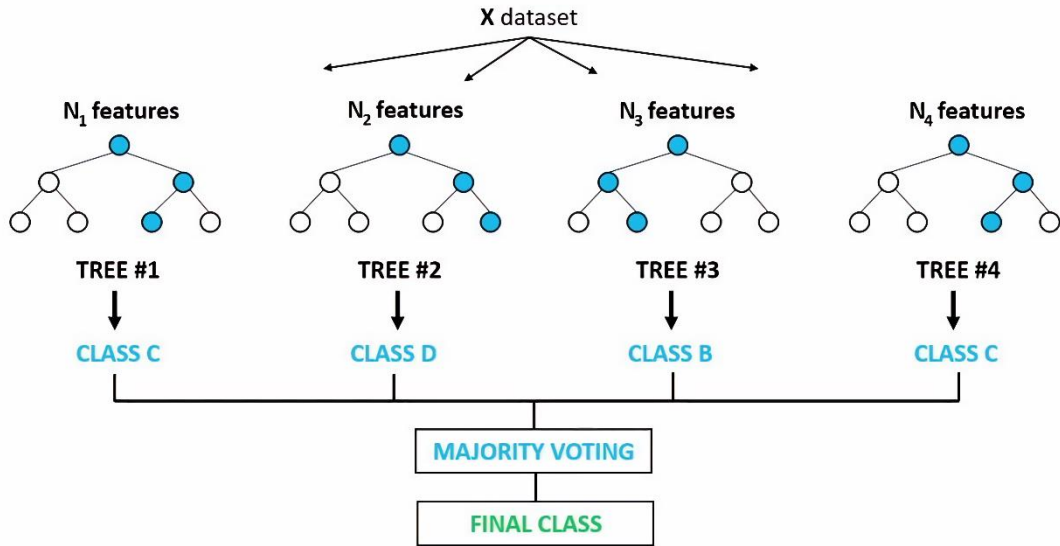


Figure 5.2 Basic decision process behind the Random Forest algorithm.

5.3 Hoeffding Tree

Considering the data stream context, where we cannot store all the data, the main problem of building a decision tree is the need to reuse instances to compute the best splitting attributes. The **Hoeffding Tree** is an incremental, anytime decision tree induction algorithm that is capable of learning from massive data streams assuming that the distribution generating examples does not change over time. It produces decision trees which are similar to the traditional batch learning method; Hoeffding trees and decision trees are asymptotically related. The HT algorithm is based on a simple idea that a small sample can be often sufficient to choose an optimal splitting attribute. The key point is that traditional batch learning methods also generate decision trees based on splitting attributes. A node is expanded as soon as there is sufficient statistical evidence that an optimal splitting feature exists, a decision based on the distribution-independent Hoeffding bound.

5.4 ADAPTIVE RANDOM FOREST

Remaining in the context of evolving data streams there is no random forest algorithm that can be considered state-of-the-art in comparison to bagging and boosting based algorithms [6]. In contrast to the traditional approach, **Adaptive Random Forest** includes an effective resampling method and adaptive operators that can cope with different types of concept drift without complex optimizations for different datasets (hence Adaptive). The adaptation process depends on an appropriate bootstrap aggregation process and limiting each leaf split decision to a subset of features. The second requirement is achieved by modifying the base tree induction algorithm, effectively by restricting the set of features considered for further splits to a random subset of size m , where $m < M$ and M corresponds to the total number of features. To address the first requirement, given the Online Bagging algorithm, proposed by the authors in [7], which approximates the original random sampling with replacement by weighting instances according to a Poisson ($\lambda = 1$) distribution, and switching the λ to 6, it is possible to “leverage” resampling in order to increase the probability of assigning higher weights to instances while training the base model.

5.5 NAÏVE BAYES

Finally, **Naïve Bayes** is a classification technique based on Bayes’ Theorem with an assumption of independence among predictors; this means that a Naïve Bayes classifier assumes that the presence of a particular feature in a class is unrelated to the presence of any other feature.

Bayes’ Theorem provides a way of calculating posterior probability $P(A/B)$ from $P(A)$, $P(B)$ and $P(B/A)$, as:

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

Here, A and B are two independent events, while $P(A)$ and $P(B)$ are their respective probabilities; $P(A/B)$ and $P(B/A)$ are the compound probabilities of A given B and B given A , respectively. We cannot obtain $P(A/B)$ directly, however we can get $P(B/A)$ and $P(A)$ from the training data: B is our set of features in the dataset, while A is the set of all predictable classes. Theoretically, finding $P(A|B)$ is not harder, however it becomes much more difficult as the number of features grows.

The advantages of using Naïve Bayes include:

- It is relatively simple to understand and build.

- It is easily trained, even with small datasets.
- It is extremely fast compared to other models.
- It is not sensitive to irrelevant features.

However, its main strength is also its weakness: the approach assumes that every feature is independent, which is not always the case.

6. *ONLINE DATASET ANALYSIS*

In order to use the dataset in online fashion, and simulate real-time learning, its entries are streamed from the arff file to MOA. Later on we also employed scikit multiflow to check whether or not the two techniques would produce similar results.

6.1 *MOA*

Three algorithms were used for the MOA analysis: Adaptive Random Forest, Hoeffding Tree, and Naïve Bayes. Amongst the three, the Adaptive Random Forest produces the best result, however it managed to do so with a significantly slower execution time. Each of the three runs produced a set of parameters on which to evaluate the model; furthermore, every 100,000 processed instances, intermediary results were logged, in order to measure the growth of the predictor over time.

For each task run on the dataset, the following attributes have been measured:

- **Classified instances:** the number of entries in the dataset that has been processed so far by the classifier.
- **Time:** the execution time, in seconds, necessary for the classification.
- **Correct predictions (accuracy):** percentage of entries successfully classified by the model so far.

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

Where, given a positive and a negative class we define:

TP – True Positives: the number of outcomes where the model correctly predicts the positive class.

TN – True Negatives: the number of outcomes where the model incorrectly predicts the positive class.

FP – False Positives: the number of outcomes where the model correctly predicts the negative class.

FN – False Negatives: the number of outcomes where the model incorrectly predicts the negative class.

- **Kappa score:** a statistic that is used to measure inter-rater reliability, a score of how much homogeneity exists in the classification.
- **Precision:** the ratio of correctly predicted positive observations to all positive observations in the positive class

$$\text{Precision} = \frac{TP}{TP + FP}$$

- **Recall:** the ratio of correctly predicted positive observations to all observations in the positive class

$$\text{Recall} = \frac{TP}{TP + FN}$$

- **F1-score (F-Measure):** a weighted average of precision and recall. It takes into account both false positives and false negatives.

$$\text{F-Measure} = 2 * \frac{\text{precision} * \text{recall}}{\text{precision} + \text{recall}}$$

The following tables and graphs display the progression of the three algorithms during their execution:

Adaptive Random Forest

Classified instances	Accuracy	Kappa	Kappa temp.	Kappa M	Precision	Recall	F1-score	Time (s)
100,000	0.938	0.938	-58	0.938	0.721	0.715	0.718	347
200,000	0.982	0.982	43	0.982	0.809	0.792	0.801	732
300,000	0.938	0.938	-40	0.938	0.716	0.703	0.709	1135
400,000	0.912	0.912	-95	0.912	0.702	0.699	0.700	1578
500,000	0.943	0.943	-19	0.943	0.723	0.709	0.716	1966
581,012	0.984	0.984	20	0.984	0.821	0.786	0.805	2232

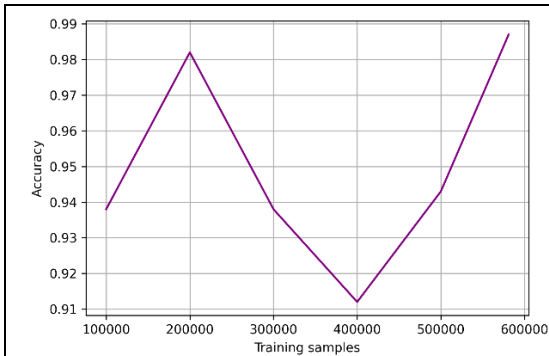


Figure 6.1 Adaptive Random Forest accuracy evolution.

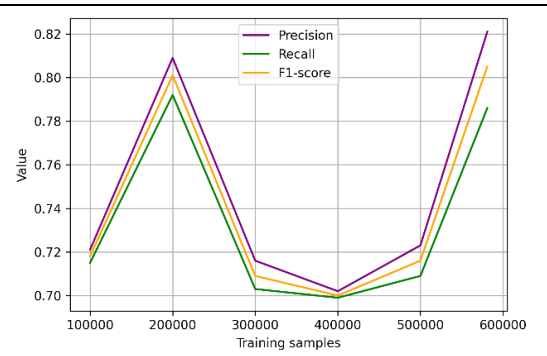


Figure 6.2 Adaptive Random Forest precision, recall and F1-score evolution.

This classifier is the one that produced the highest accuracy amongst the three, however its execution time was significantly slower compared to both Hoeffding Tree and Naïve Bayes. Its accuracy reaches the lowest value at around 400,000 classified instances but manages to regrow by the end of the process.

Hoeffding Tree

Classified instances	Accuracy	Kappa	Kappa temp.	Kappa M	Precision	Recall	F1-score	Time (s)
100,000	0.909	0.909	-133	0.909	0.723	0.718	0.720	1.92
200,000	0.832	0.832	-424	0.832	0.642	0.606	0.623	3.45
300,000	0.752	0.752	-463	0.752	0.584	0.532	0.556	5.09
400,000	0.685	0.685	-599	0.685	0.503	0.498	0.500	6.82
500,000	0.756	0.756	-419	0.756	0.579	0.529	0.552	8.67
581,012	0.673	0.673	-1534	0.673	0.499	0.472	0.485	10.18

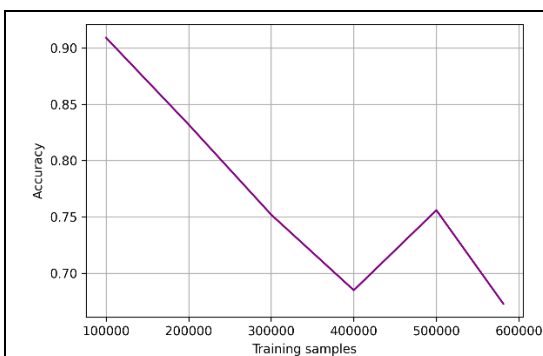


Figure 6.3 Hoeffding Tree accuracy evolution.

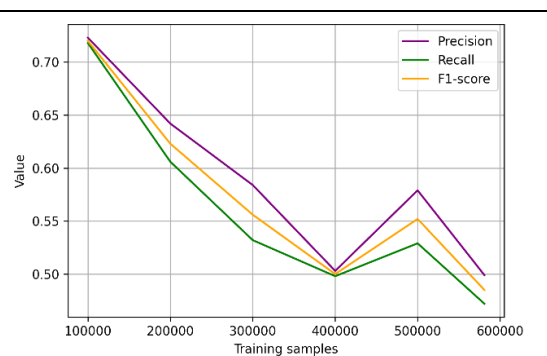


Figure 6.4 Hoeffding Tree precision, recall and F1-score evolution.

This classifier proved to be to be a very good compromise between performance and execution time. After the first 300,000 instances, we can see the accuracy oscillating; this suggests us that the classifier might have needed more samples before stabilizing its accuracy value.

Naïve Bayes

Classified instances	Accuracy	Kappa	Kappa temp.	Kappa M	Precision	Recall	F1-score	Time (s)
100,000	0.789	0.789	-441	0.789	0.589	0.472	0.524	1.75
200,000	0.526	0.526	-1381	0.526	0.392	0.314	0.348	3.46
300,000	0.643	0.643	-711	0.643	0.467	0.399	0.435	5.12
400,000	0.315	0.310	-1433	0.310	0.223	0.204	0.215	6.79
500,000	0.411	0.411	-1155	0.410	0.392	0.315	0.353	8.46
581,012	0.638	0.638	-1709	0.638	0.445	0.405	0.425	9.85

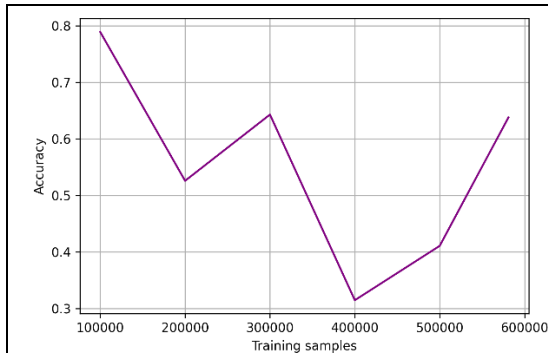


Figure 6.5 Naïve Bayes accuracy evolution.

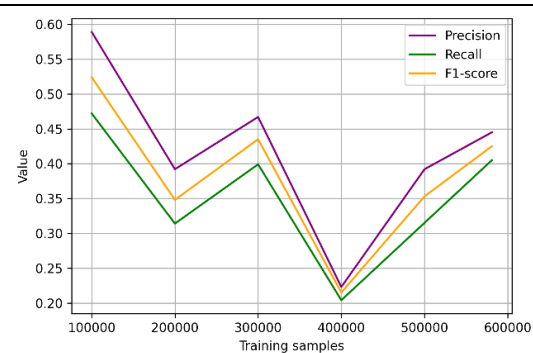


Figure 6.6 Naïve Bayes precision, recall and F1-score evolution.

This classifier performs a classic Bayesian prediction while making naïve assumption that all inputs are independent. It is a very simple algorithm with low computational costs, and this is reflected in its execution time, but also in its underwhelming accuracy when it comes to the considered problem.

By performing an early comparison between the three, we can easily notice how the Adaptive Random Forest algorithm clearly outclasses the other two in every observed metric. Therefore, if the extra training time is not an issue, the best choice easily points towards this tree classifier. Regarding the other two, they manage to stay relatively close, however the Naïve Bayes classifier seemed more “unstable” during the test runs: as we can see, at one point its accuracy drops to 31%, while the Hoeffding Tree was much more consistent in its predictions.

An interesting thing observed was how all three classifiers reach their lowest accuracy value at around 400,000 classified instances; this can be due to the fact that at that point in time, the algorithms finally encounter more samples of the classes with a lower number of entries, therefore resulting in never-seen-before patterns.

The final graph in this section directly compares the accuracy growth of these three classifiers:

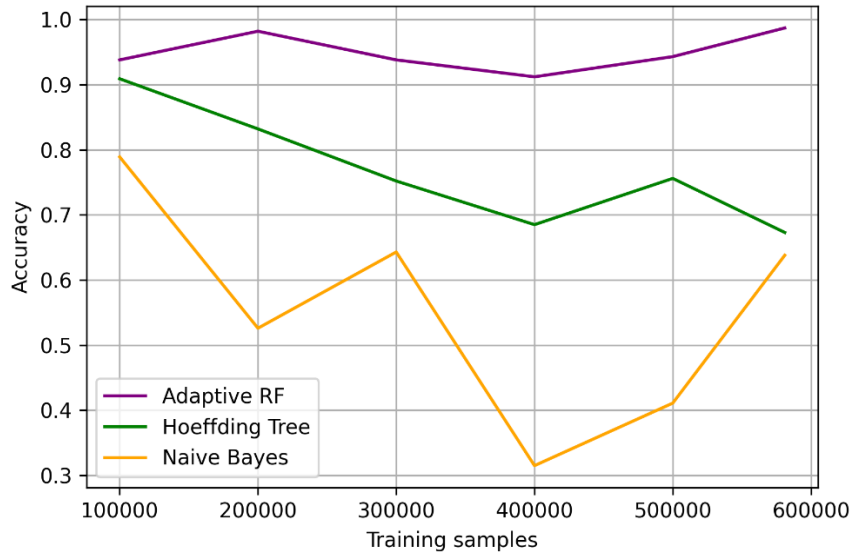


Figure 6.7 Accuracy comparison between the MOA classifiers.

6.2 SCIKIT MULTIFLOW

Just for the sake of comparison, we wanted to put MOA side-to-side with another powerful online machine learning library: Scikit Multiflow. Similarly to what was done in MOA, we took snapshots of the state of the algorithms during their execution. The results are reported ahead:

Adaptive Random Forest

#Instances	Accuracy	Kappa	Precision	Recall	F1-score	Time (s)
100,000	0.923	0.885	0.785	0.753	0.769	1221
200,000	0.931	0.889	0.790	0.754	0.772	2463
300,000	0.931	0.893	0.786	0.752	0.769	3743
400,000	0.933	0.895	0.799	0.755	0.777	5085
500,000	0.936	0.902	0.793	0.755	0.776	6291
581,012	0.940	0.902	0.798	0.760	0.777	7386

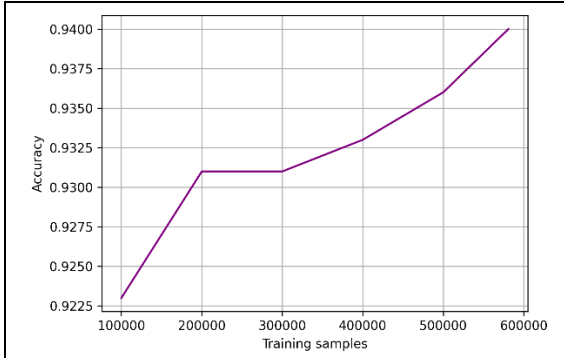


Figure 6.8 Adaptive Random Forest accuracy evolution.

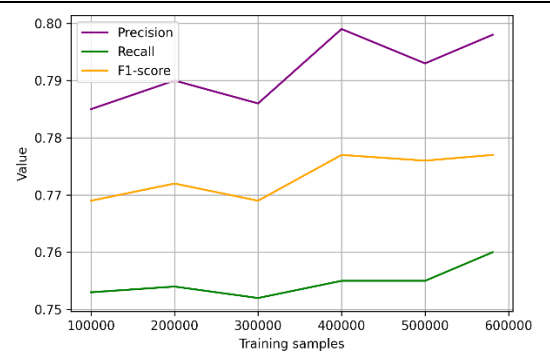


Figure 6.9 Adaptive Random Forest precision, recall and F1-score evolution.

Hoeffding Tree

#Instances	Accuracy	Kappa	Precision	Recall	F1-score	Time (s)
100,000	0.837	0.657	0.658	0.587	0.615	257
200,000	0.848	0.680	0.646	0.612	0.628	398
300,000	0.836	0.720	0.643	0.635	0.639	586
400,000	0.831	0.720	0.645	0.636	0.640	770
500,000	0.826	0.717	0.640	0.644	0.642	940
581,012	0.820	0.710	0.640	0.648	0.643	1122

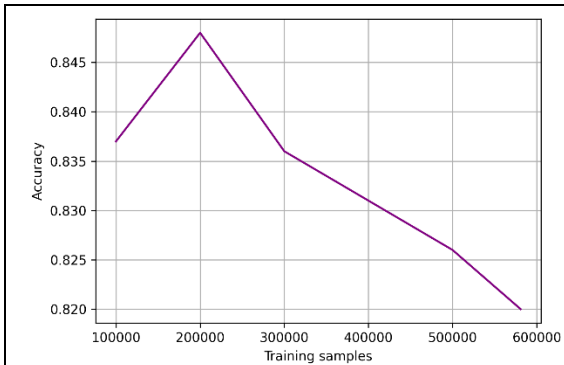


Figure 6.10 Hoeffding Tree accuracy evolution.

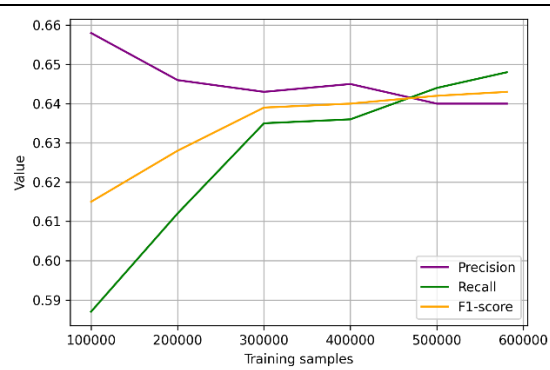
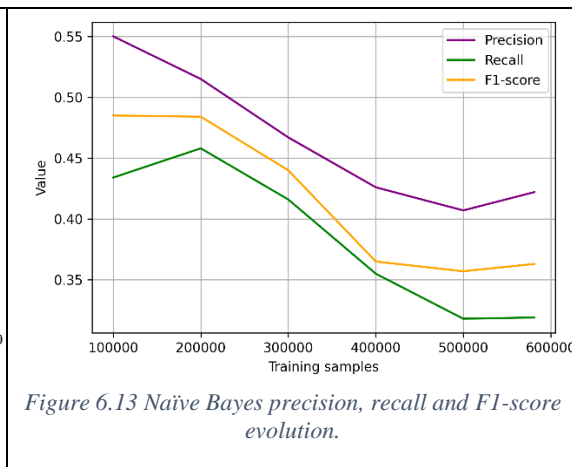
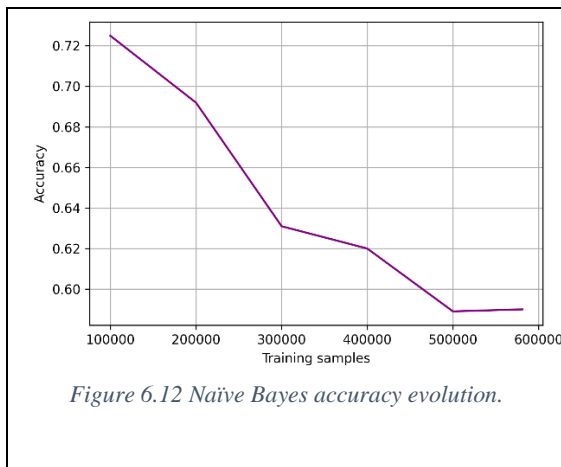


Figure 6.11 Hoeffding Tree precision, recall and F1-score evolution.

Naïve Bayes

#Instances	Accuracy	Kappa	Precision	Recall	F1-score	Time (s)
100,000	0.725	0.314	0.550	0.434	0.485	284
200,000	0.692	0.216	0.515	0.458	0.484	504
300,000	0.631	0.309	0.467	0.416	0.440	723
400,000	0.620	0.324	0.426	0.355	0.365	943
500,000	0.589	0.280	0.407	0.318	0.357	1160
581,012	0.590	0.301	0.422	0.319	0.363	1408



After a first analysis between the two online approaches we can notice that the three classifiers display the same relative behaviour relatively to each other; in particular, the Random Forest Algorithm still ranks on top, while the Naïve Bayes still performs the worse, as expected.

Regarding the differences between MOA and Scikit Multiflow there are some things worth mentioning:

- The Adaptive Random Forest was worse compared to its MOA counterpart in terms of accuracy; however precision and recall were much closer. We are not sure about the reasons of this extreme slowdown, since the parameters used for both classifiers were very similar (i.e., both used 100 estimators, and this parameter is among the ones contributing more to classification speed).
- The performance of the Hoeffding Tree algorithm in SM does not converge towards the Naïve Bayes scores, as seen in the MOA evaluation; instead, it performs much better, in terms of accuracy, precision and recall.
- The Naïve Bayes classifier also performed considerably worse compared to its MOA counterpart, especially in terms of recall.

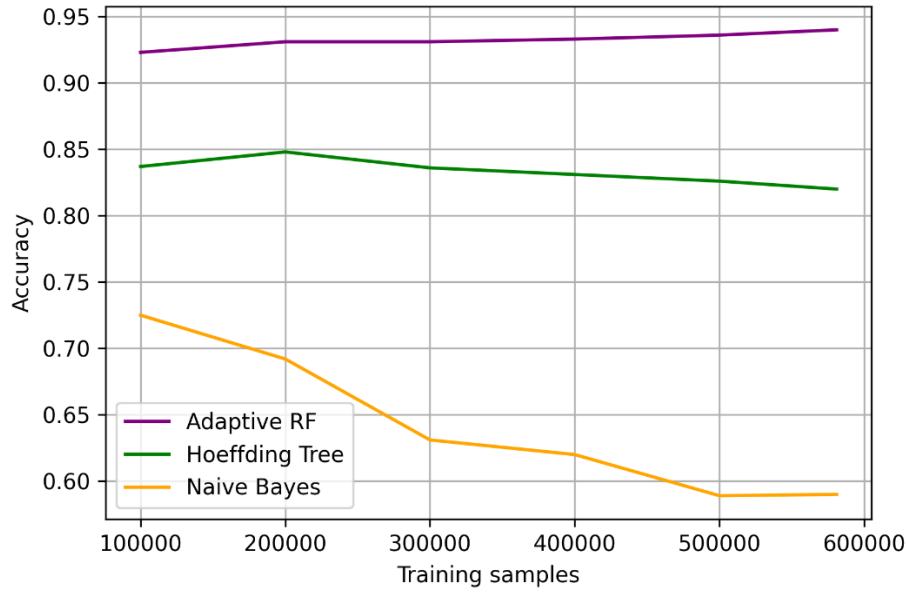


Figure 6.14 Accuracy comparison between the Scikit Multiflow classifiers.

7. BATCH LEARNING APPROACH

To further expand on the analysis of the goodness of the results obtained in MOA and Scikit Multiflow, we decided to compare these findings with other analysis techniques using batch learning. In particular, the focus shifted towards one of the most widely used libraries for machine learning, python's Scikit Learn.

In this section we will examine the results obtained by employing different batch machine learning solutions, such as traditional Random Forest, Decision Tree and Naïve Bayes. Furthermore, we will discuss on the performance results obtained by of splitting the dataset into training and test sets vs. employing 10-fold cross validation for the evaluation of the performances and we will analyse overfitting and underfitting in our results.

Finally, as a result of the previous analysis, we will apply some data balancing techniques, in order to reduce the gap between the classes, and we will remove correlated attributes and skewed attributes.

The removed attributes are:

- Hillshade_3pm.
- Vertical_Distance_To_Hidrology.
- The following soil types: 7, 8,14,15,21,25,28, 36 and 37.

7.1 DATASET SPLITTING VERSUS CROSS VALIDATION

Learning the parameters of a prediction function and testing it on the same data is a methodological mistake: a model that would just repeat the labels of the samples that it has just seen would have a perfect score but would fail to predict anything useful on unseen data. This condition is called overfitting, as we will see in the next section. To avoid this, it is common practice when performing a supervised machine learning task to hold out part of the available data as test data. When evaluating different settings for estimators (their hyperparameters) there is still a risk of overfitting on the test set because the parameters can be tweaked until the estimator performs optimally. This way knowledge about the data can “leak” into the model and evaluation metrics no longer report on generalization performance. To solve this problem, yet another part of the dataset can be held as a “validation set”: training proceeds on the train set, after which evaluation is done on the validation set, and when the experiment seems to be successful, the final evaluation on the test set is performed. However, by partitioning the available data into three sets, we drastically reduce the number of samples which can be used for learning the model, and the results can depend on a particular random choice for the pair of sets.

A solution to this problem is a procedure called cross-validation (CV). A test set should still be held out for final evaluation, but the validation set is no longer required when performing CV. In the basic approach, called k-fold CV, the training set is split into k smaller sets; then, for each of the k “folds”, the next procedure is followed:

- A model is trained using k-1 of the folds as training data.
- The resulting model is validated on the remaining part of the data.

Finally, the performance measure reported by the k-fold cross-validation is the average of the values computed in the loop. This approach can be computationally expensive, but does not waste as much data, which is a major advantage in most problems.

The following tables and graphs report the improvement obtained with 10-fold cross-validation:

Random Forest

Hyperparameters:

- Number of estimators: 100.
- Maximum depth: 23.
- Maximum number of features: “auto”.
- Minimum number of samples required to be at a leaf node: 1.
- Minimum number of samples for splitting a node: 3.

Train-test split:

Accuracy	Kappa	Precision	Recall	F1-Score	MSE	Time(s)
0.931	0.888	0.943	0.868	0.899	0.285	36.8

10-fold cross validation:

	Accuracy	Kappa	Precision	Recall	F1-Score	MSE
Fold 1	0.928	0.883	0.940	0.867	0.897	0.330
Fold 2	0.927	0.881	0.942	0.863	0.895	0.305
Fold 3	0.926	0.880	0.940	0.856	0.890	0.311
Fold 4	0.927	0.882	0.943	0.864	0.896	0.305
Fold 5	0.927	0.883	0.946	0.855	0.892	0.318
Fold 6	0.929	0.885	0.940	0.866	0.895	0.308
Fold 7	0.928	0.883	0.945	0.863	0.895	0.315
Fold 8	0.928	0.886	0.944	0.866	0.898	0.300
Fold 9	0.927	0.882	0.935	0.862	0.890	0.305
Fold 10	0.931	0.888	0.940	0.866	0.896	0.298

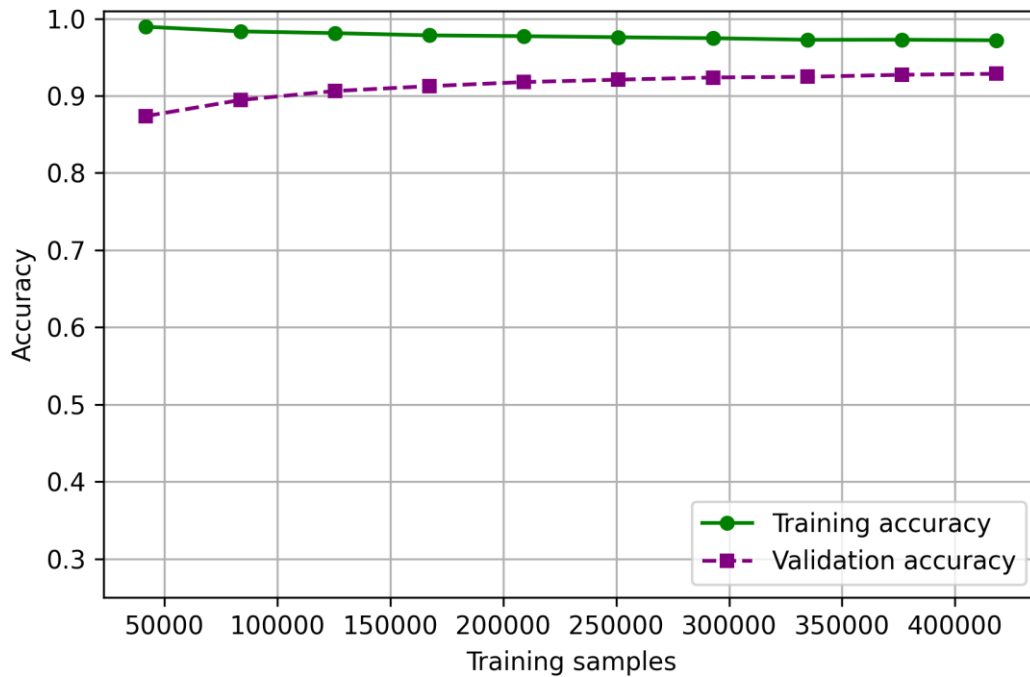


Figure 7.1 Mean accuracy of the 10-fold CV for the Random Forest classifier.

Decision Tree

Hyperparameters:

- Maximum depth: 23.

Train-test split:

Accuracy	Kappa	Precision	Recall	F1-Score	MSE	Time(s)
0.921	0.873	0.898	0.849	0.870	0.375	5.36

10-fold cross validation:

	Accuracy	Kappa	Precision	Recall	F1-Score	MSE
Fold 1	0.917	0.866	0.887	0.863	0.874	0.409
Fold 2	0.916	0.865	0.890	0.849	0.867	0.418
Fold 3	0.918	0.868	0.888	0.840	0.860	0.385
Fold 4	0.918	0.868	0.882	0.846	0.862	0.401
Fold 5	0.916	0.866	0.885	0.845	0.863	0.416
Fold 6	0.916	0.865	0.878	0.857	0.867	0.423
Fold 7	0.918	0.868	0.887	0.860	0.871	0.415
Fold 8	0.919	0.870	0.883	0.859	0.870	0.404
Fold 9	0.915	0.863	0.884	0.845	0.861	0.405
Fold 10	0.921	0.872	0.887	0.850	0.867	0.390

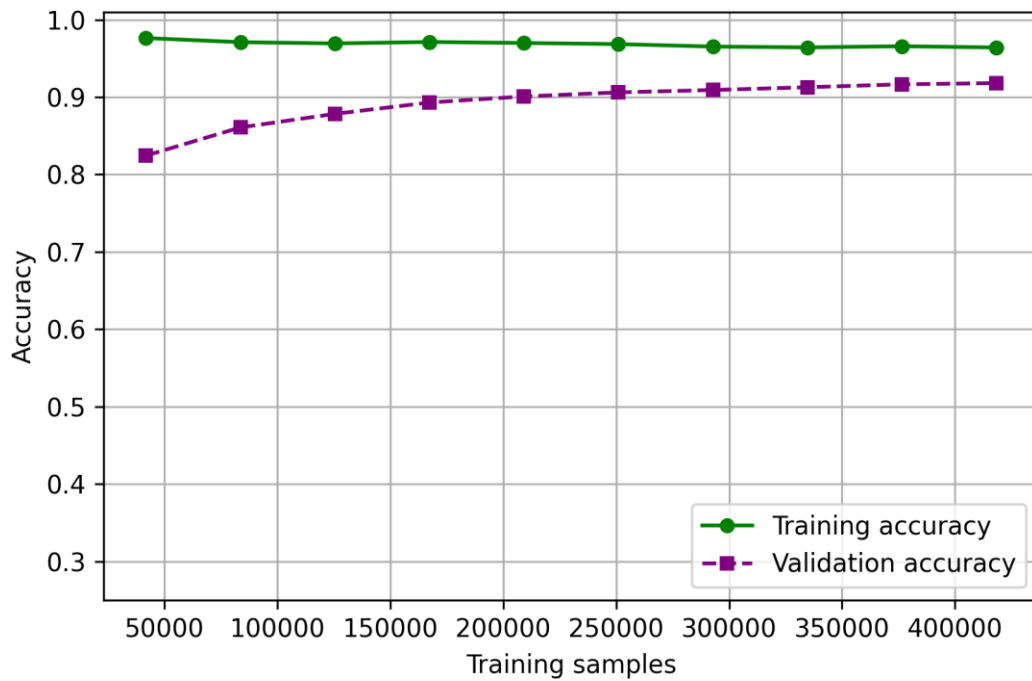


Figure 7.2 Mean accuracy of the 10-fold CV for the Decision Tree classifier.

Naïve Bayes

Train-test split:

Accuracy	Kappa	Precision	Recall	F1-Score	MSE	Time(s)
0.453	0.258	0.397	0.594	0.367	3.814	0.866

10-fold cross validation:

	Accuracy	Kappa	Precision	Recall	F1-Score	MSE
Fold 1	0.457	0.260	0.392	0.591	0.363	3.81
Fold 2	0.454	0.256	0.395	0.595	0.361	3.76
Fold 3	0.457	0.261	0.389	0.588	0.364	3.81
Fold 4	0.457	0.258	0.389	0.592	0.361	3.76
Fold 5	0.457	0.260	0.395	0.587	0.361	3.79
Fold 6	0.456	0.260	0.394	0.584	0.363	3.78
Fold 7	0.455	0.258	0.387	0.588	0.357	3.80
Fold 8	0.451	0.256	0.389	0.577	0.353	3.92
Fold 9	0.456	0.259	0.395	0.582	0.359	3.86
Fold 10	0.455	0.259	0.389	0.584	0.361	3.81

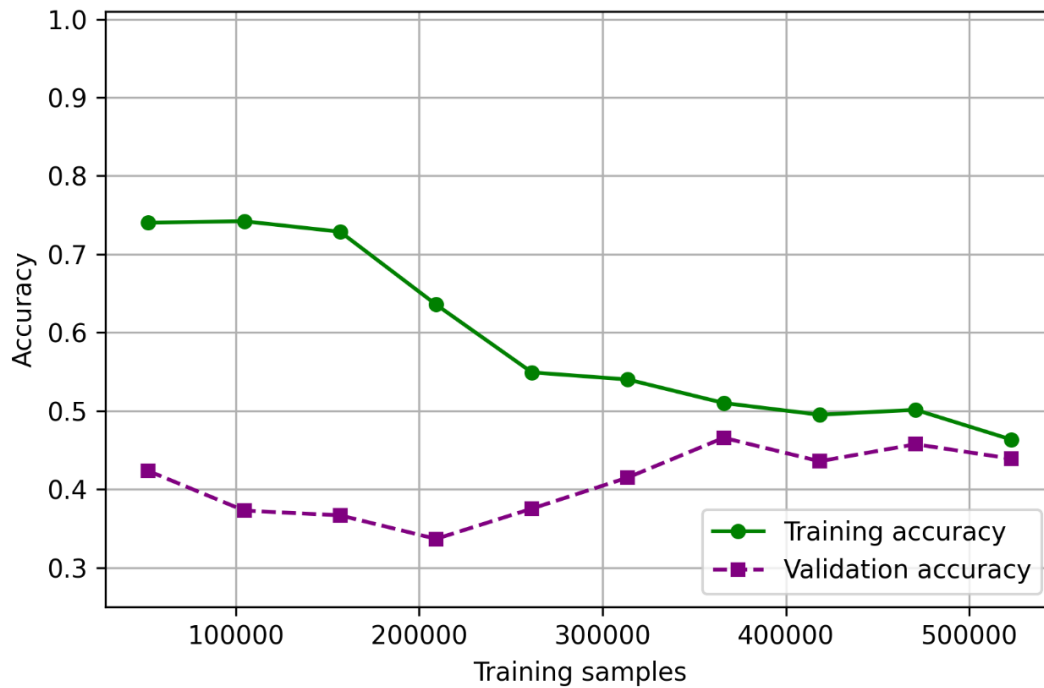


Figure 7.3 Mean accuracy of the 10-fold CV for the Naïve Bayes classifier.

7.2 OVERFITTING AND UNDERFITTING EVALUATION

Two particularly important aspects worth considering when evaluating the performance of batch classifiers are over- and underfitting. Overfitting refers to the condition by which a model learns too much from the training dataset, thus resulting in poor generalization capabilities. This can happen if the noise or random fluctuations in the training dataset are picked up and learned as concepts by the model; the problem is that these concepts do not apply to new data. Overfitting is more likely with non-parametric and non-linear models that have more flexibility when learning a target function. As such, many non-parametric machine learning algorithms also include parameters or techniques to limit and constrain how much detail the model learns.

As an example, decision trees and random forests are two kinds of machine learning algorithms that are very flexible and thus are often subject to overfitting on data. This problem can be addressed by pruning a tree after it has learned in order to remove some of the details it has picked up.

Underfitting, on the other hand, refers to a model that can neither fit training data nor generalize to new, unseen, data. An underfit machine learning model is not a suitable model and will be obvious as it will have poor performances on the training data. Underfitting is often not discussed as it is easily detectable given a good performance

metric. The remedy is to move on try alternate machine learning algorithms; nevertheless, it does provide a good contrast to the problem overfitting.

In our case, when tuning the hyperparameters of our two tree classifiers we settled for a maximum depth of 22 and 23 for the Random Forest and Decision Tree specifically. Why such a small number? Simply, by using a higher depth for the trees, we would quickly encounter an overfit model. The following two graphs highlight the accuracy curves of the two classifiers when the max depth of the trees is left untouched:

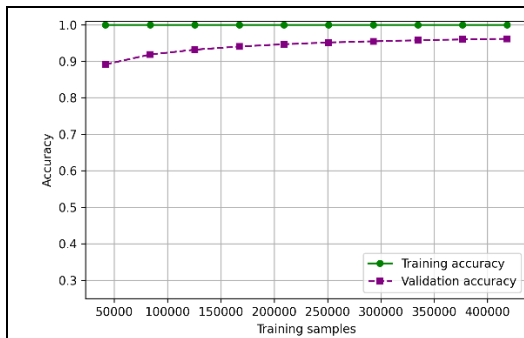


Figure 7.4 Overfitting in the Random Forest classifier.

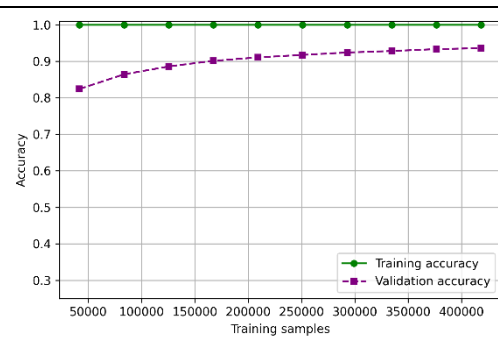


Figure 7.5 Overfitting in the Decision Tree classifier.

We can clearly see that, while the validation accuracy increases as it should, the training accuracy remains constant from the start at a value of 1. So, why does this happen? By default, the trees in these classifiers are not pruned; thus, roughly around $(1 - 1/e)$ of the trees will contain the training case and will be able to lookup the correct solution, resulting in a higher apparent accuracy. Even if the $1/e$ trees which were not trained with this case would always predict wrongly, they would always be outvoted by the others.

7.3 DATA BALANCING

A common reaction when facing an imbalanced dataset is to consider that the data are not representative of the reality: if so, we assume that real data are almost balanced but that there is a proportions bias in the collected data (i.e., due to the gathering method). The challenge of working with unbalanced datasets is that most machine learning techniques will ignore, and in turn have poor performance on, the minority class.

There are three main ways of dealing with unbalanced datasets:

- **Oversampling:** we replicate some points from the minority classes in order to create new similar instances and increase their cardinality.
- **Undersampling:** we collect a reduced number of samples from the majority classes in order to reduce their number of instances.
- **Generation of synthetic data:** we create new synthetic points from the minority classes by using an approach such as SMOTE (Synthetic Minority Oversampling Technique) to generate more instances.

In this study we focused on oversampling the minority classes, by randomly picking samples with replacement in the minority class and bringing the total number of instances to 861,566, thus generating effectively over 300,000 artificial instances.

The results obtained with this technique are reported in the following table:

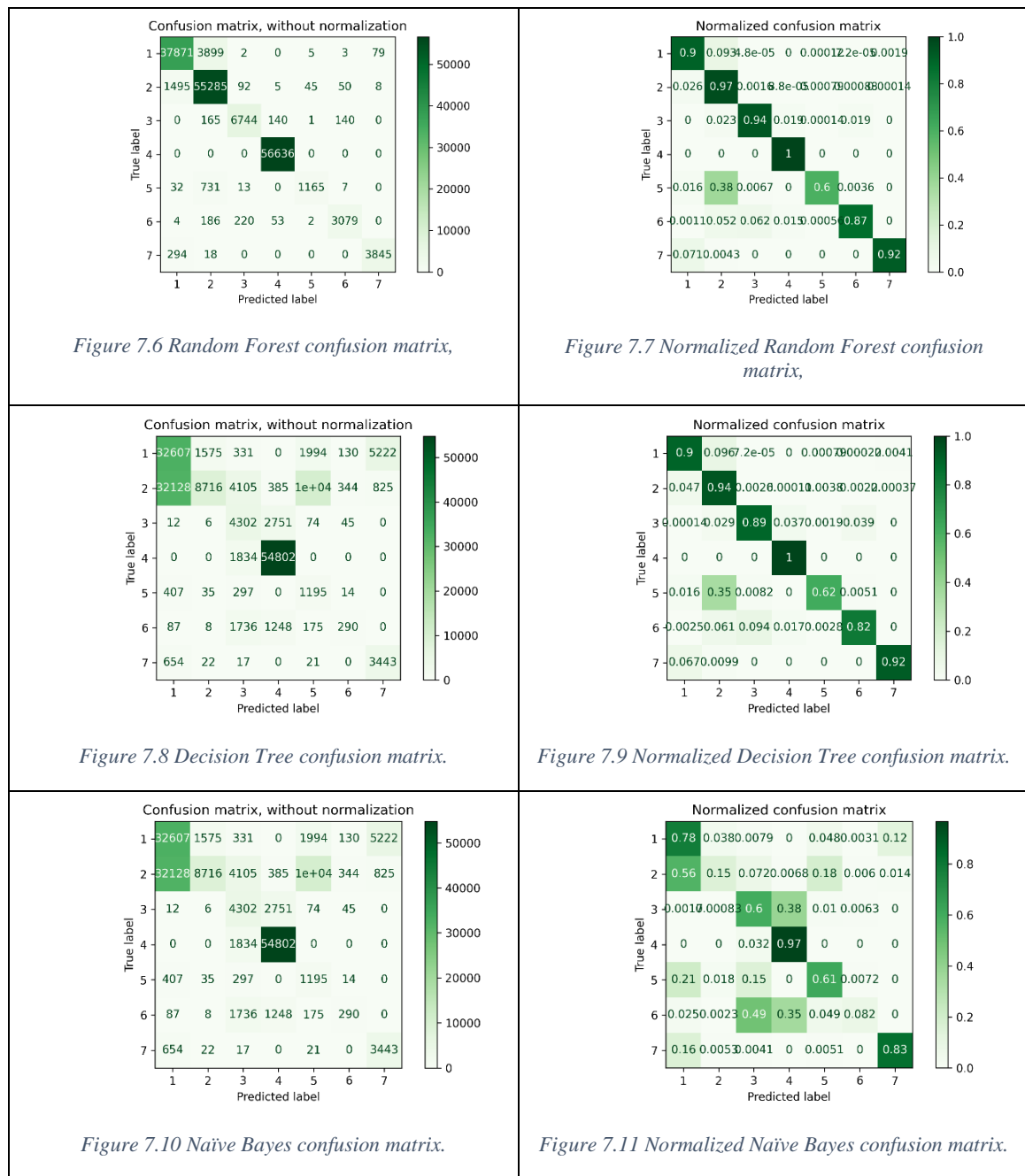
	Accuracy	Kappa	Precision	Recall	F1-Score	MSE
Random Forest	0.955	0.937	0.956	0.886	0.914	0.209
Decision Tree	0.942	0.920	0.914	0.877	0.891	0.269
Naïve Bayes	0.611	0.494	0.486	0.574	0.433	2.58

Overall the improvement is very noticeable, especially in the Naïve Bayes classifier, where we see an improvement of almost 35% in accuracy, 90% in precision, 45% in recall and 98% in mean squared error. The other two classifiers saw a growth of around 2.5% for accuracy, precision, and recall, while the mean squared error improved about 36% and 39% for the Random Forest and Decision Tree respectively..

7.4 CONFUSION MATRIX

Another interesting way to measure our results is the usage of a confusion matrix. A confusion matrix is a table layout that allows the visualization of the performance of the employed classifiers. Each row of the matrix represents the instances in the actual class, while each column represents the instances in a predicted class, or vice versa.

In this section we report the two versions, standard and normalized, for each classifier. The following table reports all of the matrices:



8. CONCLUSIONS

The forest cover type dataset is a relatively large dataset containing cartographic information about different species of three across four wilderness areas. As we saw in this study, employing machine learning techniques, in both online and batch fashion, can definitely be a viable option to build a powerful predictor for classifying the various trees in the dataset. More specifically, if we consider a scenario in which a cyberphysical system, made up of countless sensors, continuously sends data collected from the soil to a remote source in order for it to be processed, with the employed model we could be able to analyse and make predictions on the measurements in near real time.

While the initial dataset was quite unbalanced at the start, by applying various pre-processing steps to it we were successfully to alter the distribution of the data in favour of stronger prediction performances. Furthermore, different evaluation techniques were analysed in order to find the best approach, and the overfitting of our model was studied to better generalize our model on new data.

This last table provides the final comparison between all the different employed techniques and classifiers:

	Accuracy	Kappa	Precision	Recall	F1-score
<u>MOA</u>					
Adaptive Random Forest	0.984	0.984	0.821	0.786	0.805
Hoeffding Tree	0.673	0.673	0.499	0.472	0.485
Naïve Bayes	0.638	0.638	0.445	0.405	0.425
<u>Scikit Multiflow</u>					
Adaptive Random Forest	0.940	0.902	0.798	0.760	0.777
Hoeffding Tree	0.820	0.710	0.640	0.648	0.624
Naïve Bayes	0.590	0.301	0.422	0.319	0.327
<u>Scikit Learn</u>					
Random Forest	0.955	0.937	0.956	0.886	0.914
Decision Tree	0.942	0.920	0.914	0.877	0.891
Naïve Bayes	0.611	0.494	0.486	0.574	0.433

Given the nature of online learning algorithms, its more likely that they overfit on the training data. Since most online learners use a greedy, hill-climbing search in the space of models, they are much more prone to overfitting, local maxima, etc. This is the reason why the online measurements are better in the case of Random Forest and Naïve Bayes.

On the other hand, in the case of the Decision Tree algorithms, and its online counterpart, given the simpler structure of the single tree, it is much less likely to overfit; this is why the batch results are superior.

REFERENCES

- [1] Comparative accuracies of artificial neural networks and discriminant analysis in predicting forest cover type from cartographic variables, J. A. Blackard and D. J. Dean, Computers and Electronics in Agriculture, Vol 24, 1999.
- [2] Estimation and mapping of forest stand density, volume, and cover type using the k-nearest neighbors method, Hector Franco-Lopez et al., Remote Sensing of Environment, Vol.77, 2001.
- [3] Forest Cover Type Mapping and Spruce Budworm Defoliation Detection Using Simulated SPOT Imagery, Martin P. et al. Environmental Remote Sensing Centre, University of Wisconsin-Madison, Madison, 1985.
- [4] U.S. Forest Types and Predicted Percent Forest Cover from AVHRR Data, Z. Hiliang Zhu and David L. Evans, Photogrammetric Engineering & Remote Sensing, Vol.60, 1994.
- [5] Forest Cover Type Classification Study, Thomas Kolasa and Aravind Kolumum Raaja, Amazon AWS, 2016.
- [6] Adaptive random forests for evolving data stream classification, Heitor M. Gomes et al., 2017.
- [7] A new ensemble approach for dealing with concept drift, Leandro L. et al., IEEE Transactions on Knowledge and Data Engineering, 2012.