

Article



Flood Early Warning Systems Using Machine Learning Techniques: The Case of the Tomebamba Catchment at the Southern Andes of Ecuador

Paul Muñoz ^{1,2,*}, Johanna Orellana-Alvear ^{1,2}, Jörg Bendix ³, Jan Feyen ⁴ and Rolando Célleri ^{1,2}

- ¹ Departamento de Recursos Hídricos y Ciencias Ambientales, Universidad de Cuenca,
- Cuenca 010150, Ecuador; johanna.orellana@ucuenca.edu.ec (J.O.-A.); rolando.celleri@ucuenca.edu.ec (R.C.) ² Facultad de Ingeniería, Universidad de Cuenca, Cuenca 010150, Ecuador
- ³ Laboratory for Climatology and Remote Sensing, Faculty of Geography, University of Marburg, 35032 Marburg, Germany; bendix@mailer.uni-marburg.de
- Faculty of Bioscience Engineering, Catholic University of Leuven, 3001 Leuven, Belgium; jan.feyen@kuleuven.be
- * Correspondence: paul.munozp@ucuenca.edu.ec

Abstract: Worldwide, machine learning (ML) is increasingly being used for developing flood early warning systems (FEWSs). However, previous studies have not focused on establishing a methodology for determining the most efficient ML technique. We assessed FEWSs with three river states, *No-alert, Pre-alert* and *Alert* for flooding, for lead times between 1 to 12 h using the most common ML techniques, such as multi-layer perceptron (MLP), logistic regression (LR), K-nearest neighbors (KNN), naive Bayes (NB), and random forest (RF). The Tomebamba catchment in the tropical Andes of Ecuador was selected as a case study. For all lead times, MLP models achieve the highest performance followed by LR, with *f*1-macro (*log-loss*) scores of 0.82 (0.09) and 0.46 (0.20) for the 1 h and 12 h cases, respectively. The ranking was highly variable for the remaining ML techniques. According to the g-mean, LR models correctly forecast and show more stability at all states, while the MLP models perform better in the *Pre-alert* and *Alert* states. The proposed methodology for selecting the optimal ML technique for a FEWS can be extrapolated to other case studies. Future efforts are recommended to enhance the input data representation and develop communication applications to boost the awareness of society of floods.

Keywords: flood early warning; forecasting; hydrological extremes; machine learning; Andes

1. Introduction

Flooding is the most common natural hazard and results worldwide in the most damaging disasters [1–4]. Recent studies associate the increasing frequency and severity of flood events with a change in land use (e.g., deforestation and urbanization) and climate [2,5–7]. This particularly holds for the tropical Andes region, where complex hydro-meteorological conditions result in the occurrence of intense and patchy rainfall events [8–10].

According to the flood generation mechanism, floods can be classified into long- and short-rain floods [11,12]. A key for building resilience to short-rain floods is to anticipate in a timely way the event, in order to gain time for better preparedness. The response time between a rainfall event and its associated flood depends on the catchment properties and might vary from minutes to hours [13]. In this study special attention is given to flash-floods, which are floods that develop less than 6 h after a heavy rainfall with little or no forecast lead time [14].

Flood anticipation can be achieved through the development of a flood early warning system (FEWS). FEWSs have proved to be cost-efficient solutions for life preservation, damage mitigation, and resilience enhancement [15–18]. However, although crucial, flood forecasting remains a major challenge in mountainous regions due to the difficulty to



Citation: Muñoz, P.; Orellana-Alvear, J.; Bendix, J.; Feyen, J.; Célleri, R. Flood Early Warning Systems Using Machine Learning Techniques: The Case of the Tomebamba Catchment at the Southern Andes of Ecuador. *Hydrology* **2021**, *8*, 183. https:// doi.org/10.3390/hydrology8040183

Academic Editors: Marina Iosub and Andrei Enea

Received: 25 November 2021 Accepted: 13 December 2021 Published: 16 December 2021

Publisher's Note: MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Copyright: © 2021 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (https:// creativecommons.org/licenses/by/ 4.0/). effectively record the aerial distribution of precipitation due to the sparse density of the monitoring network and the absence of high-tech equipment by budget constraints [8,9].

To date, there has been no report of any operational FEWS in the Andean region for scales other than continental [17,19,20]. An alternative attempt in Peru aimed to derive daily maps of potential floods based on the spatial cumulated precipitation in past days [21]. Other endeavors in Ecuador and Bolivia focused on the monitoring of the runoff in the upper parts of the catchment to predict the likelihood of flood events in the downstream basin area [19,22]. However, such attempts are unsatisfactory as countermeasures against floods and especially flash-floods, where it is required to have reliable and accurate forecasts with lead times shorter than the response time between the farthest precipitation station and runoff control point.

There are two paradigms that drive the modeling of the precipitation-runoff response. First, the physically-based paradigm includes knowledge of the physical processes by using physical process equations [23]. This approach requires extensive ground data and, in consequence, intensive computation that hinders the temporal forecast window [24]. Moreover, it is argued that physically based models are inappropriate for real-time or short-term flood forecasting due to the inherent uncertainty of river-catchment dynamics and over-parametrization of this type of model [25]. The second data-driven paradigm assumes floods as stochastic processes with an occurrence distribution probability derived from historical data. Here, the idea is to exploit relevant input information (e.g., precipitation, past runoff) to find relations to the target variable (i.e., runoff) without requiring knowledge about the underlying physical processes. Among the traditional data-driven approaches, statistical modeling has proven to be unsuitable for short-term prediction due to lack of accuracy, complexity, model robustness, and even computational costs [24]. Previous encouraged the use of advanced data-driven models, e.g., machine learning (ML), to overcome the aforementioned shortcomings [7,24,26,27]. Particularly during the last decade, ML approaches have gained increasing popularity among hydrologists [24].

Different ML strategies for flood forecasting are implemented, generating either quantitative or qualitative runoff forecasts [18,28–38]. Qualitative forecasting consists of classifying floods into distinct categories or river states according to their severity (i.e., runoff magnitude), and use this as a base for flood class prediction [30,37,39]. The advantage of developing a FEWS is the possibility to generate a semaphore-like warning system that is easy to understand by decision-makers and the public (non-hydrologists). The challenge of FEWSs is the selection of the most optimal ML technique to obtain reliable and accurate forecasts with sufficient lead time for decision making. To date, the problem has received scant attention in the research literature, and as far as our knowledge extends no previous work examined and compared the potential and efficacy of different ML techniques for flood forecasting.

The present study compares the performance of five ML classification techniques for short-rain flood forecasting with special attention to flash floods. ML models were developed for a medium-size mountain catchment, the Tomebamba basin located in the tropical Andes of Ecuador. The ML models were tested with respect to their capacity to forecast three flood warning stages (*No-alert*, *Pre-alert* and *Alert*) for varying forecast lead times of 1, 4, and 6 h (flash-floods), but also 8 and 12 h to further test whether the lead time can be satisfactorily extended without losing the models' operational value.

This paper has been organized into four sections. The first section establishes the methodological framework for developing a FEWSs using ML techniques. It will then go on to describe the performance metrics used for a proper efficiency assessment. The second section presents the findings of the research following the same structure as the methodological section. Finally, the third and fourth sections presents the discussion and a summary of the main conclusions of the study, respectively.

2.1. Study Area and Dataset

The study area comprises the Tomebamba catchment delineated upstream of the Matadero-Sayausí hydrological station of the Tomebamba river (Figure 1), where the river enters the city. The Tomebamba is a tropical mountain catchment located in the southeastern flank of the Western Andean Cordillera, draining to the Amazon River. The drainage area of the catchment is approximately 300 km², spanning from 2800 to 4100 m above sea level (m a.s.l.). Like many other mountain catchments of the region, it is primarily covered by a páramo ecosystem, which is known for its important water regulation function [8].



Figure 1. The Tomebamba catchment located at the Tropical Andean Cordillera of Ecuador, South America (UTM coordinates).

The Tomebamba river plays a crucial role as a drinking water source for the city of Cuenca (between 25% to 30% of the demand). Other important water users are agricultural and industrial activities. Cuenca, which is the third-largest city of Ecuador (around 0.6 million inhabitants), is crossed by four rivers that annually flood parts of the city, causing human and significant economic losses.

The local water utility, the Municipal Public Company of Telecommunications, Water, Sewerage and Sanitation of Cuenca (ETAPA-EP), defined three flood alert levels associated with the Matadero-Sayausí station for floods originating in the Tomebamba catchment: (i) *No-alert* of flood occurs when the measured runoff is less than 30 m³/s, (ii) *Pre-alert* when runoff is between 30 and 50 m³/s, and (iii) the flood *Alert* is triggered when discharge exceeds 50 m³/s. With these definitions, and as shown in Figure 2, the discharge label for the *No-alert* class represents the majority of the data, whereas the *Pre-alert* and *Alert* classes comprise the minority yet the most dangerous classes.

To develop and operate forecasting models, we use data of two variables: precipitation in the catchment area and river discharge at a river gauge. For both variables, the available dataset comprises 4 years of concurrent hourly time series, from Jan/2015 to Jan/2019 (Figure 2). Precipitation information was derived from three tipping-bucket rain gauges: Toreadora (3955 m a.s.l.), Virgen (3626 m a.s.l.), and Chirimachay (3298 m a.s.l.) installed within the catchment and along its altitudinal gradient. Whereas for discharge, we used data of the Matadero-Sayausí station (2693 m a.s.l., Figure 1). To develop the ML modes, we split the dataset into training and test subsets. The training period ran from 2015 to 2017, whereas 2018 was used as the model testing phase.



Figure 2. Time series of precipitation (Toreadora) and discharge (Matadero-Sayausí). Horizontal dashed lines indicate the mean runoff and the currently employed flood alert levels for labeling the *Pre-alert* and *Alert* flood warnings classes.

2.2. Machine Learning (ML) Methods for Classification of Flood Alert Levels

ML classification algorithms can be grouped in terms of their functionality. According to Mosavi et al. (2018), five of the worldwide most-popular statistical method groups are commonly used for short-term flood prediction (extreme runoff), and include:

- i. Regression algorithms modeling the relationships between variables (e.g., logistic regression, linear regression, multivariate adaptive regression splines, etc.) [18,40].
- ii. Instance-based algorithms that rely on memory-based learning, representing a decision problem fed with data for training (e.g., K-nearest neighbor, learning vector quantification, locally weighted learning, etc.) [30].
- Decision tree algorithms, which progressively divide the whole data set into subsets based on certain feature values, until all target variables are grouped into one category (e.g., classification and regression tree, M5, random forest, etc.) [18,28,30,31,37].
- iv. Bayesian algorithms based on Bayes' theorem on conditional probability (e.g., naive Bayes, Bayesian network, Gaussian naïve Bayes, etc.) [18,31,35].
- v. Neural Network algorithms inspired by biological neural networks convert input(s) to output(s) through specified transient states that enable the model to learn in a sophisticated way (e.g., perceptron, multi-layer perceptron, radial basis function network, etc.) [18,31,36].

For this study, we selected five ML algorithms, one from each group, respectively, a logistic regression, K-nearest neighbor, random forest, naive Bayes, and a multi-layer perceptron.

2.2.1. Logistic Regression

Logistic Regression (LR) is a discriminative model, modeling the decision boundary between classes. In a first instance, linear regressions are applied to find existent relationships between model features. Thereafter, the probability (conditional) of belonging to a class is identified using a logistic (sigmoid) function that effectively deals with outliers (binary classification). From these probabilities, the LR classifies, with regularization, the dependent variables into any of the created classes. However, for multiclass classification problems are all binary classification possibilities considered, it is *No-alert* vs. *Pre-alert*, *No-alert* vs. *Alert*, and *Pre-alert* vs. *Alert*. Finally, the solution is the classification with the maximum probability (multinomial LR) using the *softmax* function Equation (1). With this function is the predicted probability of each class defined [41]. The calculated probability for each class is positive with the logistic function and normalized across all classes.

$$softmax(z)_{i} = \frac{e^{z_{i}}}{\sum_{l=1}^{k} e^{z_{l}}}$$
(1)

where z_i is the *i*th input of the *softmax* function, corresponding to class *i* from the *k* number of classes.

2.2.2. K-Nearest Neighbors

K-Nearest Neighbors (KNN) is a non-parametric statistical pattern recognition algorithm, for which no theoretical or analytical background exist but an intuitive statistical procedure (memory-based learning) for the classification. KNN classifies unseen data based on a similarity measure such as a distance function (e.g., Euclidean, Manhattan, Chebyshev, Hamming, etc.). The use of multiple neighbors instead of only one is recommended to avoid the wrong delineation of class boundaries caused by noisy features. In the end, the majority vote of the nearest neighbors (see the formulation in [41]) determines the classification decision. The number of nearest neighbors can be optimized to reach a global minimum avoiding longer computation times, and the influence of class size. The major advantage of the KNN is its simplicity. However, the drawback is that KNN is memory intensive, all training data must be stored and compared when added information is to be evaluated.

2.2.3. Random Forest

Random Forest (RF) is a supervised ML algorithm that ensembles a multitude of decorrelated decision trees (DTs) voting for the most popular class (classification). In practice, a DT (particular model) is a hierarchical analysis based on a set of conditions consecutively applied to a dataset. To assure decorrelation, the RF algorithm applies a bagging technique for a growing DT from different randomly resampled training subsets obtained from the original dataset. Each DT provides an independent output (class) of the phenomenon of interest (i.e., runoff), contrary to numerical labels for regression applications. The popularity of RF is due to the possibility to perform random subsampling and bootstrapping which minimizes biased classification [42]. An extended description of the RF functioning is available in [43,44].

The predicted class probabilities of an input sample are calculated as the mean predicted class probabilities of the trees in the forest. For a single tree, the class probability is computed as the fraction of samples of the same class in a leaf. However, it is well-known that the calculated training frequencies are not accurate conditional probability estimates due to the high bias and variance of the frequencies [45]. This deficiency can be resolved by controlling the minimum number of samples required at a leaf node, with the objective to induce a smoothing effect, and to obtain statistically reliable probability estimates.

2.2.4. Naïve Bayes

Naïve Bayes (NB) is a classification method based on Bayes' theorem with the "naive" assumption that there is no dependence between features in a class, even if there is dependence [46]. Bayes' theorem can be expressed as:

$$P(y|X) = \frac{P(X|y) P(y)}{P(X)}$$
⁽²⁾

where P(A|B) is the probability of y (hypothesis) happening, given the occurrence of X (features), and X can be defined as $X = x_1, x_2, ..., x_n$. Bayes' theorem can be written as:

$$P(y|x_1, x_2, \dots, x_n) = \frac{P(x_1|y) P(x_2|y) \dots P(x_n|y) P(y)}{P(x_1) P(x_2) \dots P(x_n)}$$
(3)

There are different NB classifiers depending on the assumption of the distribution of $P(x_i | y)$. In this matter, the study of Zhang [46] proved the optimality of NB under the Gaussian distribution even when the assumption of conditional independence is violated (real application cases). Additionally, for multiclass problems, the outcome of the algorithm is the class with the maximum probability. For the Gaussian NB algorithm no parameters have to be tuned.

2.2.5. Multi-Layer Perceptron

The Multi-Layer Perceptron (MLP) is a class of feedforward artificial neural networks (ANN). A perceptron is a linear classifier that separates an input into two categories with a straight line and produces a single outcome. Input is a feature vector multiplied by specific weights and added to a bias. Contrary to a single-layer case, the MLP can approximate non-linear functions using additional so-called hidden layers. Prediction of probabilities of belonging to any class is calculated through the *softmax* function. The MLP consists of multiple neurons in fully connected multiple layers. Determination of the number of neurons in the layers with a trial-and-error approach remains widely used [47]. Neurons in the first layer correspond to the input data, whereas all other nodes relate inputs to outputs by using linear combinations with certain weights and biases together with an activation function. To measure the performance of the MLP, the logistic loss function is defined with the limited memory Broyden–Fletcher–Goldfarb–Shanno (L-BFGS) method as the optimizer for training the network. A detailed and comprehensive description of ANN can be found in [48].

2.3. Methodology

Figure 3 depicts schematic the methodology followed in this study. The complete dataset for the study consists, as mentioned before, of precipitation and labeled discharge time-series (see Figure 2). The dataset was split in two groups, respectively, for training and testing purposes, and training and test feature spaces were composed for each lead time for the tasks of model hyperparameterization and model assessment. This procedure is repeated for each of the ML techniques studied. Finally, the ranking of the performance quality of all ML methods for every lead time, based on performance metrics and a statistical significance test, were determined.

2.3.1. Feature Space Composition

For each lead time, we used single training and testing feature spaces for all ML techniques. A feature space is composed by features (predictors) coming from two variables: precipitation and discharge. The process of feature space composition starts by defining a specific number of precipitation and discharge features (present time and past hourly lags) according to statistical analyses relying on Pearson's cross-, auto and partial-auto-correlation functions [49]. The number of lags from each station was selected by setting up a correlation threshold of 0.2 [28].

Similarly, for discharge, we used several features coming from past time slots of discharge selected for the analysis. It is worth noting that the number of discharge features triples since we replace each discharge feature with three features (one per flood warning class) in a process known as one-hot-encoding or binary encoding. Therefore, each created feature denotes 0 or 1 when the correspondent alarm stage is false or true, respectively. Finally, we performed a feature standardization process before the computation stage of the KNN, LR, NB, and NN algorithms. Standardization was achieved by subtracting the mean and scaling it to unit variance, resulting in a distribution with a standard deviation equal to 1 and a mean equal to 0.





Figure 3. Work schedule for the development and evaluation of the machine learning (ML) flood forecasting models.

2.3.2. Model Hyperparameterization

After the composition of the feature space the optimal architectures for each ML forecasting model, and for each lead time was set up. The optimal architectures were defined by the combination of hyperparameters under the concept of balance between accuracy, and computational cost, and speed. However, finding optimal architectures requires an exhaustive search of all combinations of hyperparameters. To overcome this issue, we relied on the randomized grid search (RGS) with a 10-fold cross-validation scheme. The RGS procedure randomly explores the search space for discretized continuous hyperparameters based on a cross-validation evaluation. Moreover, we selected the f1-macro score (see Section 2.3.4) as objective function.

2.3.3. Principal Component Analysis

ML applications require in general the analysis of high-dimension and complex data, involving substantial amounts of memory and computational costs. Reduction of the dimensionality was realized through the application of principal component analysis (PCA) enabling exclusion of correlating features that do not add information to the model. PCA was applied after feature scaling and normalization.

This method enables finding the dimension of maximum variance and the reduction of the feature space to that dimension so that the model performance remains as intact as possible when compared to performance with the full feature space. But considering that each ML technique assimilates data differently, we did not define the number of principal components to keep a fixed threshold of variance explanation (e.g., 80–90%), but performed an exploratory analysis to evaluate its influence on each model. As such, the number of PCAs was treated as an additional hyperparameter, and we optimized the number of principal components for each specific model (lead time and ML technique) with the objective to find the best possible model for each case.

All ML techniques and the RGS procedure were implemented through the scikit-learn package for ML in Python[®] [50]. Table 1 presents the relevant hyperparameters for each ML technique and their search space for tuning [38]. We employed default values for the hyperparameters which are depicted in Table 1.

ML Technique	Hyperparameters					
LR	C 0.001–1000	penalty {'11', '12'}				
	neighbor's	weights	metric	algorithm		
KNN	3–75	{'uniform', 'distance'}	{'euclidean', 'manhattan', 'minkowski'}	{'auto','ball_tree', 'kd_tree','brute'}		
	estimator's	max_features	hadeeth	min_samples_leaf	min_samples_split	
RF	50-1000	{'auto', 'sqrt', 'log2'}	50-1000	1–500	1–500	
MLP	solver {'lbfgs'}	max_iter 10–5000	alpha $1 imes 10^{-9}$ –0.1	hidden_layers 1–16		

Table 1. Model hyperparameters and their ranges/possibilities for tuning.

2.3.4. Model Performance Evaluation

Forecasting hydrological extremes such as floods turns into an imbalanced classification problem, and becomes even more complex when the interest lies in the minority class of the data (flood alert). This is because most ML classification algorithms focus on the minimization of the overall error rate, it is the incorrect classification of the majority class [51]. Resampling the class distribution of the data for obtaining an equal number of samples per class is one solution. In this study, we used another approach that relies on training ML models with the assumption of imbalanced data. The approach we used penalizes mistakes in samples belonging to the minority classes rather than under-sampling or over-sampling data. In practice, this implies that for a given metric efficiency, the overall score is the result of averaging each performance metric (for each class) multiplied by its corresponding weight factor. According to the class frequencies the weight factors for each class were calculated (inversely proportional), using Equation (4).

$$w_i = \frac{N}{C n_i} \tag{4}$$

where w_i is the weight of class *i*, *N* is the total number of observations, *C* is the number of classes, and n_j the number of observations in class *i*. This implies that higher weights will be obtained for minority classes.

Performance Metrics

The metrics for the performance assessment were derived from the well-known confusion matrix, especially suitable for imbalanced datasets and multiclass problems, and are respectively the f1 score, the geometric mean, and the logistic regression loss score [51–56]. Since neither of the metrics is adequate it is suggested to use a compendium

of metrics to properly explain the performance of the model. In addition, those metrics complement each other.

f1 Score

The f score is a metric that relies on precision and recall, which is an effective metric for imbalanced problems. When the f score as a weighted harmonic mean, we name this score f1. The latter score can be calculated with Equation (5).

$$f1 \text{ score} = \frac{2 \times \text{Precision } \times \text{Recall}}{(\text{Precision } + \text{Recall})}$$
(5)

where precision and recall are defined with the following equations:

$$Precision = \frac{TP}{TP + FP}$$
(6)

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}} \tag{7}$$

where *TP* stands for true positives, *FP* for false positives, and *FN* for false negatives.

The f1 score ranges from 0 to 1, indicating perfect precision and recall. The advantage of using the f1 score compared to the arithmetic or geometric mean is that it penalizes models most when either the precision or recall is low. However, classifying a *No-Alert* flood warning as *Alert* might have a different impact on the decision-making than when the opposite occurs. This limitation scales up when there is an additional state, e.g., *Pre-alert*. Thus, the interpretation of the f1 score must be taken with care. For multiclass problems, the f1 score is commonly averaged across all classes, and is called the f1-macro score to indicate the overall model performance.

Geometric Mean

The geometric-mean (g-mean) measures simultaneously the balanced performance of TP and TN rates. This metric gives equal importance to the classification task of both the majority (*No-alert*) and minority (*Pre-alert* and *Alert*) classes. The g-mean is an evaluation measure that can be used to maximize accuracy to balance TP and TN examples at the same time with a good trade-off [53]. It can be calculated using Equation (8)

$$G - mean = \sqrt{(TP_{rate} \times TN_{rate})}$$
(8)

where TP_{rate} and TN_{rate} are defined by:

$$TP_{rate} = Recall$$
 (9)

$$\Gamma N_{rate} = \frac{TN}{TN + FP}$$
(10)

The value of the g-mean metrics ranges from 0 to 1, where low values indicate deficient performance in the classification of the majority class even if the minority classes are correctly classified.

Logistic Regression Loss

The metric logistic regression loss (log-loss) measures the performance of a classification model when the input is a probability value between 0 and 1. It accounts for the uncertainty of the forecast based on how much it varies from the actual label. For multiclass classification, a separate log-loss is calculated for each class label (per observation), and the results are summed up. The log-loss score for multi-class problems is defined as:

$$Log \ loss = -\frac{1}{N} \sum_{i=1}^{N} \sum_{j=1}^{M} y_{ij} \log(p_{ij})$$
(11)

where *N* is the number of samples, *M* the number of classes, y_{ij} equal to 1 when the observation belongs to class *j*; else 0, and p_{ij} is the predicted probability that the observation belongs to class *j*. Starting from 0 (best score), the log-loss magnitudes increase as the probability diverges from the actual label. It punishes worse errors more harshly to promote conservative predictions. For probabilities close to 1, the log-loss slowly decreases. However, as the predicted probability decreases, the log-loss increases rapidly.

Statistical Significance Test for Comparing Machine-Learning (ML) Algorithms

Although we can directly compare performance metrics of ML alternatives and claim to have found the best one based on the score, it is not certain whether the difference in metrics is real or the result of statistical chance. Different statistical frameworks are available allowing us to compare the performance of classification models (e.g., a difference of proportions, paired comparison, binomial test, etc.).

Among them, Raschka [57] recommends using the chi-squared test to quantify the likelihood of the samples of skill scores, being observed under the assumption that they have the same distributions. The assumption is known as the null hypothesis, and aims to prove whether there is a statistically significant difference between two models (error rates). If rejected, it can be concluded that any observed difference in performance metrics is due to a difference in the models and not due to statistical chance. In our study we used the chi-squared test to assess whether the difference in the observed proportions of the contingency tables of a pair of ML algorithms (for a given lead time) is significant.

For the model comparison, we defined the statistical significance of improvements/ degradations for all lead times (training and test subsets) under a value of 0.05 (chi-squared test). In all cases, the MLP model was used as the base model to which the other models were compared.

3. Results

This section presents the results of the flood forecasting models developed with the LR, KNN, RF, NB, and MLP techniques, and for lead times of 1, 4, 6, 8, and 12 h. For each model, we addressed the forecast of three flood warnings, *No-alert, Pre-alert* and *Alert*. First, we present the results of the feature space composition process, taking the 1 h lead time case as an example. Then, we show the results of the hyperparameterization for all models, followed by an evaluation and ranking of the performance of the ML techniques.

3.1. Feature Space Composition

Figures 4 and 5 show the results of the discharge and precipitation lag analyses for the flood forecasting model 1-h before the flood would occur. Figure 4a depicts the discharge autocorrelation function (ACF) and the corresponding 95% confidence interval from lag 1 up to 600 (h). We found a significant correlation up to a lag of 280 h (maximum correlation at the first lag) and, thereafter, the correlation fell within the confidence band. On the other hand, Figure 4b presents the discharge partial-autocorrelation function (PACF) and its 95% confidence band from lag 1 to 30 h. We found a significant correlation up to lag 8 h (first lags outside the confidence band). As a result, based on the interpretation of the ACF and PACF analyses, and according to Muñoz et al. [28] we decided to include 8 discharge lags (hours) for the case of 1 h flood forecasting in the Tomebamba catchment.



Figure 4. (a) Autocorrelation function (ACF) and (b) partial-autocorrelation function (PACF) of the Matadero-Sayausí (Tomebamba catchment) discharge series. The blue hatch indicates in each case the correspondent 95% confidence interval.



Figure 5. Pearson's cross-correlation comparison between the Toreadora (3955 m a.s.l.), Virgen (3626 m a.s.l.), and Chirimachay (3298 m a.s.l.) precipitation stations and the Matadero-Sayausí discharge series. Note the blue horizontal line at a fixed correlation of 0.2 for determining past lags.

Figure 5 plots the Pearson's cross-correlation between the precipitation at each rainfall station and the discharge at the Matadero-Sayausí stream gauging station. For all stations, we found a maximum correlation at lag 4 (maximum 0.32 for Chirimachay). With the fixed correlation threshold of 0.2, we included 11, 14, and 15 lags for Virgen, Chirimachay, and Toreadora stations, respectively.

Similarly, the same procedure was applied for the remaining lead times (i.e., 4, 6, 8, and 12 h). In Table 2, we present the input data composition and the resulting total number of features obtained from the lag analyses for each forecasting model. For instance, for the 1 h case, the total number of features in the feature space equals 67, from which 43 are derived from precipitation (past lags and one feature from present time for each station), and 24 from discharge (one-hot-encoding).

	Discharge Lags * (h)		Precipitation Lags (h)		
Lead Time (h)	Matadero-Sayausí	Toreadora	Chirimachay	Virgen	Number of Features
1	8	15	14	11	67
4	12	18	17	14	88
6	14	20	19	16	100
8	16	22	21	18	112
12	20	26	25	22	136

Table 2. Input data composition (number of features) for all ML models of the Tomebamba catchment.

* Note that each discharge feature triples (three flood warning classes) after a one-hot-encoding process.

3.2. Model Hyperparameterization

The results of the hyperparameterization including the number of PCA components employed for achieving the best model efficiencies are presented in Table 3. No evident relation between the number of principal components and the ML technique nor the lead time was found. In fact, for some models we found differences in the f1-macro score lower than 0.01 for a low and high number of principal components. See for instance the case of the KNN models where the optimal number of components significantly decayed for lead times greater than 4 h. For the 1 h lead time, 96% of the components were used, whereas for the rest of the lead times only less than 8%.

Table 3. Model hyperparameters and number of principal components used for each specific model (ML technique and lead time).

ML Technique	Use on or other	Lead Time					
	nyperparameter	1 h	4 h	6 h	8 h	12 h	
LR	С	0.01	0.00001	0.0001	0.0001	0.001	
	penalty	'12'	'12'	'12'	'12'	'12'	
	PCA_components *	58	62	78	75	51	
	n_neighbors	15	15	23	33	55	
	weights	'uniform'	'uniform'	'uniform'	'uniform'	'uniform'	
KNN	metric	'minkowski'	'minkowski'	'minkowski'	'minkowski'	'minkowski'	
	Algorithm	'auto'	'auto'	'auto'	'auto'	'auto'	
	PCA_components *	64	6	6	6	4	
	n_estimators	700	700	700	700	800	
	max_features	'sqrt'	'auto'	auto	ʻlog2′	'auto'	
DE	max_depth	350	350	350	350	300	
KF	min_samples_leaf	450	450	480	480	450	
	min_samples_split	10	5	5	2	4	
	PCA_components *	66	79	90	45	78	
NB	PCA_components *	63	64	87	89	15	
MLP	solver	'lbfgs'	'lbfgs'	'lbfgs'	'lbfgs'	'lbfgs'	
	max_iter	2000	2000	2000	2000	2000	
	alpha	0.0001	0.0001	0.0001	0.0001	0.0001	
	hidden_layers	2	3	2	2	4	
	PCA_components *	63	51	64	76	4	

* From the total number of features: 1 h = 67, 4 h = 88, 6 h = 100, 8 h = 112, 12 h = 136 features.

If we turn to the evolution of models' complexity with lead time (Table 3) more complex ML architectures are needed to forecast greater lead times. This is underpinned by the fact that the corresponding optimal models require for greater lead times a stronger regularization (lower values of *C*) for LR, a greater number of neighbors (n_neighbors) for KNN, more specific trees (lower values of min_samples_split) for RF and more hidden layers (hidden_layers) for MLP.

3.3. Model Performance Evaluation

As mentioned before, model performances calculated with the f1-score, g-mean, and log-loss score were weighted according to class frequencies. Table 4 presents the frequency distribution for the complete dataset, respectively, for the training and test subsets. Here, the dominance of the *No-alert* flood class is evident, with more than 95% of the samples in both subsets. With this information, the class weights for the training period were calculated as $w_{No-alert} = 0.01$, $w_{Pre-alert} = 0.55$ and $w_{Alert} = 0.51$.

Table 4. The number of samples and relative percentage for the entire dataset and the training and test subsets.

Warning	Complete	Training	Test
No-alert	32,596 (96.1%)	24,890 (96.2%)	7706 (95.7%)
Pre-alert	720 (2.1%)	473 (1.8%)	247 (3.1%)
Alert	609 (1.8%)	(2.0%)	100 (1.2%)

The results of the model performance evaluation for all ML models and lead times (test subset) are summarized in Table 5. We proved for all models that the differences in performance metrics for a given lead time were due to the difference in the ML techniques rather than to the statistical chance. As expected, ML models' ability to forecast floods decreased for a longer lead time. For instance, for the case of 1 h forecasting, we found a maximum f1-macro score of 0.88 (MLP) for the training and 0.82 (LR) for the test subset. Whereas, for the 12 h case, the maximum f1-macro score was 0.71 (MLP) for the training and 0.46 (MLP) for the test subset.

Table 5. Models' performance evaluation on the test subset. Bold fonts indicate the best performance for a given lead time.

Lead Time (h)	RF	KNN	LR	NB	MLP			
F1-macro score								
1	0.59	0.73	0.82	0.57	0.78			
4	0.47	0.57	0.59	0.46	0.62			
6	0.47	0.45	0.50	0.41	0.51			
8	0.44	0.41	0.44	0.45	0.51			
12	0.42	0.36	0.44	0.43	0.46			
	G-mean							
1	0.86	0.77	0.88	0.81	0.83			
4	0.75	0.63	0.76	0.73	0.71			
6	0.70	0.56	0.72	0.68	0.62			
8	0.73	0.53	0.67	0.62	0.62			
12	0.69	0.50	0.69	0.64	0.56			
Log-loss score								
1	0.28	0.38	1.09	3.14	0.09			
4	0.38	0.46	0.74	4.10	0.11			
6	0.45	0.58	0.47	4.71	0.14			
8	0.50	0.65	0.53	0.59	0.16			
12	0.59	0.70	0.57	2.17	0.20			

Note: All improvements and degradations are statistically significant.

The extensive hyperparameterization (RGS scheme) powered by 10-fold cross-validation served to assure robustness in all ML models and reduced overfitting. We found only a small difference between the performance values by using the training and the test subsets. For all models, maximum differences in performances were lower than 0.27 for the f1-macro score and 0.19 for the g-mean.

In general, for all lead times, the MLP technique obtained the highest f1-macro score, followed by the LR algorithm. This performance dominance was confirmed by the ranking of the models according to the log-loss score. The ranking of the remaining models was highly variable and, therefore, not conclusive. For instance, the results of the KNN models obtained the second-highest score for the training subset, but the lowest for the test subset, especially for longer lead times. This is because the KNN is a memory-based algorithm and therefore more sensitive to the inclusion of information different to the training subset in comparison to the remaining ML techniques. This can be noted in Table 4, where the training and test frequency distributions are different for the *Pre-alert* and *Alert* classes.

On the other hand, for the g-mean score, we obtained a different ranking of the methods. We found the highest scores for the LR algorithm, followed by the RF and the MLP models. Despite this behavior, the values of the g-mean were superior to the f1-macro scores for all lead times and subsets. This is because the f1 score relies on the harmonic mean. Therefore, the f1 score penalizes more a low precision or recall in comparison with a metric based on a geometric or arithmetic mean. Results of the g-mean served to identify that the LR is the most stable method in terms of correctly classifying both the majority (*No-alert*) and the minority (*Pre-alert* and *Alert*) flood warning classes, while the MLP technique could be used to focus on the minority (flood alert) classes.

To extend the last idea, we analyzed the individual f1 scores of each flood warning class. This unveils the ability of the model to forecast the main classes of interest, i.e., Pre-alert and Alert. Figure 6 presents the evolution of the f1-score of each ML algorithm at the corresponding lead time. We found that for all ML techniques, the Alert class is clearly the most difficult to forecast when the f1-macro score was selected as the metric for the hyperparameterization task. An additional exercise consisted in choosing the individual f1-score for the *Alert* class as the target for hyperparameterization of all models. However, although we obtained comparable results for the Alert class, the scores of the Pre-alert class had significantly deteriorated, even reaching scores near zero. The most interesting aspect in Figure 6 is that the most efficient and stable models across lead times (test subset) were the models based on MLP and LR techniques. It is also evident that for all forecasting models, a lack of robustness for the Pre-alert warning class was found, and there were major differences between the f1-scores for the training and test subsets. An explanation for this might be that the *Alert* class implies a *Pre-alert* warning class, but not the opposite. Consequently, this might mislead the learning process causing overfitting during training leading to poor performances when assessing unseen data during the test phase.



Figure 6. Cont.



Figure 6. f1 scores per flood warning state (*No-alert*, *Pre-alert* and *Alert*) for all combinations of ML techniques across lead times. (**a**), Logistic Regression (**b**), K-Nearest Neighbors, (**c**) Random Forest, (**d**) Naïve Bayes, and (**e**) Multi-layer Perceptron. The brightest and dashed lines in each subfigure (color coding) represent the scores for the test subset.

Moreover, although we added a notion of class frequency distribution (weights) to the performance evaluation task, it can be noted that for all models, the majority class is most perfectly classified. This is because the *No-alert* class arises from low-to-medium discharge magnitudes. This eases and simplifies the learning process of the ML techniques since these magnitudes can be related to normal conditions (present time and past lags) of precipitation and discharge.

4. Discussion

In this study, we developed and evaluated five different FEWSs relying on the most common ML techniques for flood forecasting, and for short-term lead times of 1, 4, and 6 h for flash-floods, and 8 and 12 h to assess models' operational value for longer lead times. Historical runoff data were used to define and label the three flood warning scenarios to be forecasted (*No-alert*, *Pre-alert* and *Alert*). We constructed the feature space for the models according to the statistical analyses of precipitation and discharge data followed by a PCA analysis embedded in the hyperparameterization.

This was aimed to better exploit the learning algorithm of each ML technique. In terms of model assessment, we proposed an integral scheme based on the f1-score, the geometric mean, and the log-loss score to deal with data imbalance and multiclass characteristics. Finally, the assessment was complemented with a statistical analysis to provide

a performance ranking between ML techniques. For all lead times, we obtained the best forecasts for both, the majority and minority classes from the models based on the LR, RF, and MLP techniques (g-mean). The two most suitable models for the dangerous warning classes (*Pre-Alert* and *Alert*) were the MLP and LR (f1 and log-loss scores). This finding has important implications for developing FEWSs since real-time applications must be capable of dealing with both the majority and minority classes. Therefore, it can be suggested that the most appropriate forecasting models are based on the MLP technique.

The results on the evolution of model performances across lead times suggest that the models are acceptable for lead times up to 6 h, i.e., the models are suitable for flash-flood applications in the Tomebamba catchment. For lead times greater than 6 h, we found a strong decay in model performance. In other words, the utility of the 8 and 12 h forecasting models is limited by the models' operational value. This is because, in the absence of rainfall forecasts, the assumption of future rain is solely based on runoff measurements at past and present times. This generates forecasts that are not accurate enough for horizons greater than the concentration-time of the catchment. The concentration-time of the Tomebamba catchment was estimated between 2 and 6 h according to the equations of Kirpich, Giandotti, Ven Te Chow, and Temez, respectively. A summary of the equations can be found in Almeida et al. [58]. This results in an additional performance decay for the 8 and 12 h cases in addition to the error in modeling.

The study of Furquim et al. [31] is comparable. These authors analyzed the performance of different ML classification algorithms for flash-flood nowcasting (3 h) in a river located in an urban area of Brazil. They found that models based on neural networks and decision trees outperformed those based on the NB technique. In addition, the study of Razali et al. [30] proved that decision tree-based algorithms perform better than KNN models, which agrees with our findings. However, such studies only evaluated the percentage of correctly classified instances which is a simplistic evaluation. Thus, we recommend a more integral assessment of model performances, like the one in the current study, which allows for better support in decision making.

Other studies related to quantitative forecasting revealed that neural network-based models usually outperform the remaining techniques proposed in our study [32–34]. Similarly, the study of Khalaf et al. [37] proved the superiority of the RF algorithm when compared to the bagging decision trees and HyperPipes classification algorithms. Thus, in certain cases, the use of less expensive techniques regarding the computational costs produces comparable results as in [36]; this is also the case in our short-rain and flash-flood flood classification problem.

As a further step, we propose the development of ensemble models for improving the performance results of individual models. This can be accomplished by combining the outcomes of the ML models with weights obtained, for instance, from the log-log scores. Another alternative that is becoming popular is the construction of hybrid models as a combination of ML algorithms for more accurate and efficient models [24,35,36]. Moreover, as stated by Solomatine and Xue [36], inaccuracies in forecasting floods are mainly due to data-related problems. In this regard, Muñoz et al. [9] reported a deficiency in precipitation-driven models due to rainfall heterogeneity in mountainous areas, where orographic rainfall formation occurs. In most cases, rainfall events are only partially captured by punctual measurement, and even the entire storm coverage can be missing.

In general precipitation-runoff models will reach at a certain point an effectiveness threshold that cannot be exceeded without incorporating new types of data such as soil moisture [59,60]. In humid areas, the rainfall–runoff relationship also depends on other variables such as evapotranspiration, soil moisture, and land use, which leads to significant spatial variations of water storage. However, these variables are difficult to measure or estimate.

5. Conclusions

- The current study set out to propose a methodology and integral evaluation framework for developing optimal short-rain flood warning forecasting models using ML classification techniques. The proposed analyses were applied to forecast three flood warnings, *No-alert, Pre-alert* and *Alert* for the Tomebamba catchment in the tropical Andes of Ecuador. For this, the five most common ML classification techniques for short-term flood forecasting were used. From the results, the following conclusion can be drawn: results related to model comparison are statistically significant. This is important because this is not usually performed in other studies and it validates the performance comparison and ranking hereby presented.
- For all lead times, the most suitable models for flood forecasting are based on the MLP followed by the LR techniques. From the integral evaluation (i.e., several performance metrics), we suggest LR models as the most efficient and stable option for classifying both the majority (*No-alert*) and the minority (*Pre-alert* and *Alert*) classes whereas we recommend MLP when the interest lies in the minority classes.
- The forecasting models we developed are robust. Differences in the averaged f1, gmean and log-loss scores between training and test are consistent to all models. However, we limit the utility of the models for flash-flood applications (lead times up to 6 h). For longer lead times, we encourage improvement in precipitation representation, and even forecasting this variable for lead times longer than the concentration-time of the catchment.

A more detailed model assessment (individual f1 scores) demonstrated the difficulties of forecasting the *Pre-alert* and *Alert* flood warnings. This was evidenced when the hyperparameterization was driven for the optimization of the forecast for the alert class and this, however, did not improve the model performance of this specific class. This study can be extended with a deep exploration of the effect of input data composition, precipitation forecasting, and the feature engineering strategies for both the MLP and LR techniques. Feature engineering pursues the use of data representation strategies that could, for example, provide spatial and temporal information of the precipitation in the study area. This can be done by spatially discretizing precipitation in the catchments with the use of remotely sensed imagery. With this additional knowledge, it would be possible to improve the performance of the models hereby developed at longer lead times.

We recommend that future efforts should be put into applying the methodology and assessment framework proposed here in other tropical Andean catchments, and/or benchmarking the results obtained in this study with the outputs of physically based forecasting models. This was not possible for this study due to lack of data.

Finally, for FEWSs, the effectiveness of the models is strongly linked to the speed of communication to the public after a flood warning is triggered. Therefore, future efforts should focus on the development of a web portal and/or mobile application as a tool to boost the preparedness of society against floods that currently threaten people's lives, possessions, and environment in Cuenca and other comparable tropical Andean cities.

Author Contributions: Conceptualization, P.M. and R.C.; formal analysis, P.M.; funding acquisition, R.C.; investigation, P.M.; methodology, P.M. and J.O.-A.; project administration, R.C.; supervision, J.O.-A., J.B., J.F. and R.C.; writing—original draft, P.M.; writing—review and editing, J.O.-A., J.B., J.F. and R.C. All authors have read and agreed to the published version of the manuscript.

Funding: This research was funded by the project: "Desarrollo de modelos para pronóstico hidrológico a partir de datos de radar meteorológico en cuencas de montaña", funded by the Research Office of the University of Cuenca (DIUC), and the Empresa Pública Municipal de Telecomunicaciones, Agua Potable, Alcantarillado y Saneamiento de Cuenca (ETAPA-EP). Our thanks go to these institutions for their generous funding.

Institutional Review Board Statement: Not applicable.

Informed Consent Statement: Not applicable.

Data Availability Statement: All data, models, and code that support the findings of this study are available from the corresponding author upon reasonable request.

Acknowledgments: We acknowledge the Ministry of Environment of Ecuador (MAAE) for providing research permissions, and are grateful to the staff and students that contributed to the hydrometeorological monitoring. for experiments).

Conflicts of Interest: The authors declare no conflict of interest.

References

- 1. Stefanidis, S.; Stathis, D. Assessment of flood hazard based on natural and anthropogenic factors using analytic hierarchy process (AHP). *Nat. Hazards* **2013**, *68*, 569–585. [CrossRef]
- Paprotny, D.; Sebastian, A.; Morales-Nápoles, O.; Jonkman, S. Trends in flood losses in Europe over the past 150 years. *Nat. Commun.* 2018, 9, 1985. [CrossRef] [PubMed]
- 3. Ávila, Á.; Guerrero, F.C.; Escobar, Y.C.; Justino, F. Recent Precipitation Trends and Floods in the Colombian Andes. *Water* **2019**, *11*, 379. [CrossRef]
- 4. Mirza, M.M.Q. Climate change, flooding in South Asia and implications. Reg. Environ. Chang. 2011, 11, 95–107. [CrossRef]
- Sofia, G.; Roder, G.; Fontana, G.D.; Tarolli, P. Flood dynamics in urbanised landscapes: 100 years of climate and humans' interaction. *Sci. Rep.* 2017, 7, 40527. [CrossRef] [PubMed]
- 6. Min, S.-K.; Zhang, X.; Zwiers, F.W.; Hegerl, G.C. Human contribution to more-intense precipitation extremes. *Nature* **2011**, 470, 378–381. [CrossRef] [PubMed]
- Chang, L.-C.; Chang, F.-J.; Yang, S.-N.; Kao, I.-F.; Ku, Y.-Y.; Kuo, C.-L.; Amin, I.R.; bin Mat, M.Z. Building an Intelligent Hydroinformatics Integration Platform for Regional Flood Inundation Warning Systems. *Water* 2019, 11, 9. [CrossRef]
- Célleri, R.; Feyen, J. The Hydrology of Tropical Andean Ecosystems: Importance, Knowledge Status, and Perspectives. *Mt. Res. Dev.* 2009, 29, 350–355. [CrossRef]
- 9. Muñoz, P.; Célleri, R.; Feyen, J. Effect of the Resolution of Tipping-Bucket Rain Gauge and Calculation Method on Rainfall Intensities in an Andean Mountain Gradient. *Water* **2016**, *8*, 534. [CrossRef]
- 10. Arias, P.A.; Garreaud, R.; Poveda, G.; Espinoza, J.C.; Molina-Carpio, J.; Masiokas, M.; Viale, M.; Scaff, L.; van Oevelen, P.J. Hydroclimate of the Andes Part II: Hydroclimate Variability and Sub-Continental Patterns. *Front. Earth Sci.* **2021**, *8*, 666. [CrossRef]
- 11. Hundecha, Y.; Parajka, J.; Viglione, A. Flood type classification and assessment of their past changes across Europe. *Hydrol. Earth Syst. Sci. Discuss.* **2017**, 1–29.
- 12. Turkington, T.; Breinl, K.; Ettema, J.; Alkema, D.; Jetten, V. A new flood type classification method for use in climate change impact studies. *Weather. Clim. Extrem.* **2016**, *14*, 1–16. [CrossRef]
- 13. Borga, M.; Gaume, E.; Creutin, J.D.; Marchi, L. Surveying flash floods: Gauging the ungauged extremes. *Hydrol. Process.* **2008**, 22, 3883–3885. [CrossRef]
- 14. Knocke, E.T.; Kolivras, K.N. Flash Flood Awareness in Southwest Virginia. Risk Anal. Int. J. 2007, 27, 155–169. [CrossRef] [PubMed]
- 15. Sottolichio, A.; Hurther, D.; Gratiot, N.; Bretel, P. Acoustic turbulence measurements of near-bed suspended sediment dynamics in highly turbid waters of a macrotidal estuary. *Cont. Shelf Res.* **2011**, *31*, S36–S49. [CrossRef]
- 16. Borga, M.; Anagnostou, E.; Blöschl, G.; Creutin, J.-D. Flash flood forecasting, warning and risk management: The HYDRATE project. *Environ. Sci. Policy* **2011**, *14*, 834–844. [CrossRef]
- 17. del Granado, S.; Stewart, A.; Borbor, M.; Franco, C.; Tauzer, E.; Romero, M. Flood Early Warning Systems. Comparative Analysis in Three Andean Countries (Sistemas de Alerta Temprana para Inundaciones: Análisis Comparativo de Tres Países Latinoamericanos); Institute for Advanced Development Studies (INESAD): La Paz, Bolivia, 2016. (In Spanish)
- Noymanee, J.; Theeramunkong, T. Flood forecasting with machine learning technique on hydrological modeling. *Procedia Comput. Sci.* 2019, 156, 377–386. [CrossRef]
- 19. Dávila, D. Flood Warning Systems in Latin America (21 Experiencias de Sistemas de Alerta Temprana en América Latina); Soluciones Prácticas: Lima, Peru, 2016.
- 20. Boers, N.; Bookhagen, B.; Barbosa, H.D.M.J.; Marwan, N.; Kurths, J.; Marengo, J.A. Prediction of extreme floods in the eastern Central Andes based on a complex networks approach. *Nat. Commun.* **2014**, *5*, 5199. [CrossRef]
- 21. Aybar Camacho, C.L.; Lavado-Casimiro, W.; Huerta, A.; Fernández Palomino, C.; Vega-Jácome, F.; Sabino Rojas, E.; Felipe-Obando, O. Use of the gridded product 'PISCO' for precipitation studies, investigations and operationl systems of monitoring and hydrometeorological forecasting (Uso del Producto Grillado 'PISCO' de precipitación en Estudios, Investigaciones y Sistemas Operacionales de Monitoreo y Pronóstico Hidrometeorológico). *Nota Técnica* 2017. No. 001 SENAMHI-DHI.
- Fernández de Córdova Webster, C.J.; Rodríguez López, Y. First results of the current hydrometeorological network of Cuenca, Ecuador(Primeros resultados de la red actual de monitoreo hidrometeorológico de Cuenca, Ecuador). *Ing. Hidráulica Ambient*. 2016, 37, 44–56.
- Clark, M.P.; Bierkens, M.F.P.; Samaniego, L.; Woods, R.A.; Uijlenhoet, R.; Bennett, K.E.; Pauwels, V.R.N.; Cai, X.; Wood, A.W.; Peters-Lidard, C.D. The evolution of process-based hydrologic models: Historical challenges and the collective quest for physical realism. *Hydrol. Earth Syst. Sci.* 2017, 21, 3427–3440. [CrossRef]

- Mosavi, A.; Ozturk, P.; Chau, K.-W. Flood Prediction Using Machine Learning Models: Literature Review. Water 2018, 10, 1536. [CrossRef]
- Young, P.C. Advances in real-time flood forecasting. *Philos. Trans. R. Soc. Lond. Ser. A Math. Phys. Eng. Sci.* 2002, 360, 1433–1450. [CrossRef] [PubMed]
- Bontempi, G.; Ben Taieb, S.; Le Borgne, Y.-A. Machine Learning Strategies for Time Series Forecasting. In *Lecture Notes in Business Information Processing*; Springer: Berlin/Heidelberg, Germany, 2013; pp. 62–77.
- 27. Galelli, S.; Castelletti, A. Assessing the predictive capability of randomized tree-based ensembles in streamflow modelling. *Hydrol. Earth Syst. Sci.* 2013, 17, 2669–2684. [CrossRef]
- 28. Muñoz, P.; Orellana-Alvear, J.; Willems, P.; Célleri, R. Flash-Flood Forecasting in an Andean Mountain Catchment—Development of a Step-Wise Methodology Based on the Random Forest Algorithm. *Water* **2018**, *10*, 1519. [CrossRef]
- Furquim, G.; Neto, F.; Pessin, G.; Ueyama, J.; Joao, P.; Clara, M.; Mendiondo, E.M.; de Souza, V.C.; de Souza, P.; Dimitrova, D.; et al. Combining wireless sensor networks and machine learning for flash flood nowcasting. In Proceedings of the 2014 28th International Conference on Advanced Information Networking and Applications Workshops, Victoria, BC, Canada, 13–16 May 2014; pp. 67–72.
- 30. Toukourou, M.; Johannet, A.; Dreyfus, G.; Ayral, P.-A. Rainfall-runoff modeling of flash floods in the absence of rainfall forecasts: The case of 'Cévenol flash floods. *Appl. Intell.* **2011**, *35*, 178–189. [CrossRef]
- 31. Adamowski, J.F. Development of a short-term river flood forecasting method for snowmelt driven floods based on wavelet and cross-wavelet analysis. *J. Hydrol.* **2008**, 353, 247–266. [CrossRef]
- 32. Aichouri, I.; Hani, A.; Bougherira, N.; Djabri, L.; Chaffai, H.; Lallahem, S. River Flow Model Using Artificial Neural Networks. *Energy Procedia* **2015**, *74*, 1007–1014. [CrossRef]
- Khosravi, K.; Shahabi, H.; Pham, B.T.; Adamowski, J.; Shirzadi, A.; Pradhan, B.; Dou, J.; Ly, H.-B.; Gróf, G.; Ho, H.L.; et al. A comparative assessment of flood susceptibility modeling using Multi-Criteria Decision-Making Analysis and Machine Learning Methods. J. Hydrol. 2019, 573, 311–323. [CrossRef]
- 34. Solomatine, D.P.; Xue, Y. M5 Model Trees and Neural Networks: Application to Flood Forecasting in the Upper Reach of the Huai River in China. *J. Hydrol. Eng.* **2004**, *9*, 491–501. [CrossRef]
- Khalaf, M.; Hussain, A.J.; Al-Jumeily, D.; Fergus, P.; Idowu, I.O. Advance flood detection and notification system based on sensor technology and machine learning algorithm. In Proceedings of the 2015 International Conference on Systems, Signals and Image Processing (IWSSIP), London, UK, 10–12 September 2015; pp. 105–108.
- 36. Contreras, P.; Orellana-Alvear, J.; Muñoz, P.; Bendix, J.; Célleri, R. Influence of Random Forest Hyperparameterization on Short-Term Runoff Forecasting in an Andean Mountain Catchment. *Atmosphere* **2021**, *12*, 238. [CrossRef]
- Orellana-Alvear, J.; Célleri, R.; Rollenbeck, R.; Muñoz, P.; Contreras, P.; Bendix, J. Assessment of Native Radar Reflectivity and Radar Rainfall Estimates for Discharge Forecasting in Mountain Catchments with a Random Forest Model. *Remote Sens.* 2020, 12, 1986. [CrossRef]
- 38. Razali, N.; Ismail, S.; Mustapha, A. Machine learning approach for flood risks prediction. IAES Int. J. Artif Intell. 2020, 9, 73. [CrossRef]
- 39. Chen, S.; Xue, Z.; Li, M. Variable Sets principle and method for flood classification. *Sci. China Ser. E Technol. Sci.* 2013, 56, 2343–2348. [CrossRef]
- 40. Jati, M.I.H.; Santoso, P.B. Prediction of flood areas using the logistic regression method (case study of the provinces Banten, DKI Jakarta, and West Java). J. Phys. Conf. Ser. 2019, 1367, 12087. [CrossRef]
- 41. Bishop, C.M. Pattern Recognition and Machine Learning; Springer: New York, NY, USA, 2006.
- Dodangeh, E.; Choubin, B.; Eigdir, A.N.; Nabipour, N.; Panahi, M.; Shamshirband, S.; Mosavi, A. Integrated machine learning methods with resampling algorithms for flood susceptibility prediction. *Sci. Total. Environ.* 2020, 705, 135983. [CrossRef] [PubMed]
 Define the Definition of the Machine Learning and the Computer of the Computer of
- 43. Breiman, L. Random forests. *Mach. Learn.* 2001, 45, 5–32. [CrossRef]
- 44. Breiman, L. Classification and Regression Trees; Routledge: Oxfordshire, UK, 2017.
- 45. Zadrozny, B.; Elkan, C. Obtaining calibrated probability estimates from decision trees and naive bayesian classifiers. In Proceedings of the Eighteenth International Conference on Machine Learning, Williamstown, MA, USA, 28 June–1 July 2001; Morgan Kaufmann Publishers, Inc.: San Francisco, CA, USA; pp. 609–616.
- 46. Zhang, H. The optimality of naive Bayes. AA 2004, 1, 3.
- 47. Maier, H.R.; Jain, A.; Dandy, G.C.; Sudheer, K. Methods used for the development of neural networks for the prediction of water resource variables in river systems: Current status and future directions. *Environ. Model. Softw.* **2010**, 25, 891–909. [CrossRef]
- 48. Maier, H.R.; Dandy, G.C. Neural networks for the prediction and forecasting of water resources variables: A review of modelling issues and applications. *Environ. Model. Softw.* **2000**, *15*, 101–124. [CrossRef]
- Sudheer, K.P.; Gosain, A.K.; Ramasastri, K.S. A data-driven algorithm for constructing artificial neural network rainfall-runoff models. *Hydrol. Process.* 2002, 16, 1325–1330. [CrossRef]
- 50. Pedregosa, F.; Varoquaux, G.; Gramfort, A.; Michel, V.; Thirion, B.; Grisel, O.; Blondel, M.; Prettenhofer, P.; Weiss, R.; Dubourg, V.; et al. Scikit-learn: Machine Learning in {P}ython. *J. Mach. Learn. Res.* **2011**, *12*, 2825–2830.
- 51. Chen, C.; Liaw, A.; Breiman, L. *Using Random Forest to Learn Imbalanced Data*; University of California: Berkeley, CA, USA, 2004; Volume 110, pp. 1–12.
- 52. Read, J.; Pfahringer, B.; Holmes, G.; Frank, E. Classifier chains for multi-label classification. Mach. Learn. 2011, 85, 333–359. [CrossRef]

- 53. Gu, Q.; Zhu, L.; Cai, Z. Evaluation Measures of the Classification Performance of Imbalanced Data Sets. In *Computational Intelligence and Intelligent Systems*; Springer: Berlin/Heidelberg, Germany, 2009; Volume 51, pp. 461–471. [CrossRef]
- 54. Sun, Y.; Wong, A.K.C.; Kamel, M.S. Classification of imbalanced data: A review. Int. J. Pattern Recognit. Artif. Intell. 2009, 23, 687–719. [CrossRef]
- 55. Hossin, M.; Sulaiman, M.N. A review on evaluation metrics for data classification evaluations. *Int. J. Data Min. Knowl. Manag. Process.* **2015**, *5*, 1.
- 56. Akosa, J. Predictive accuracy: A misleading performance measure for highly imbalanced data. In Proceedings of the SAS Global Forum, Orlando, FL, USA, 2–5 April 2017; pp. 2–5.
- 57. Raschka, S. Model evaluation, model selection, and algorithm selection in machine learning. arXiv 2018, arXiv:1811.12808.
- 58. de Almeida, I.K.; Almeida, A.K.; Anache, J.A.A.; Steffen, J.L.; Sobrinho, T.A. Estimation on time of concentration of overland flow in watersheds: A review. *Geociencias* 2014, 33, 661–671.
- Loumagne, C.; Normand, M.; Riffard, M.; Weisse, A.; Quesney, A.; Le Hégarat-Mascle, S.; Alem, F. Integration of remote sensing data into hydrological models for reservoir management. *Hydrol. Sci. J.* 2001, 46, 89–102. [CrossRef]
- 60. Li, Y.; Grimaldi, S.; Walker, J.P.; Pauwels, V.R.N. Application of Remote Sensing Data to Constrain Operational Rainfall-Driven Flood Forecasting: A Review. *Remote Sens.* **2016**, *8*, 456. [CrossRef]