








Forecasting Solar Flares Using Magnetogram-based Predictors and Machine Learning

Kostas Florios^{1,2}  · Ioannis Kontogiannis¹  · Sung-Hong Park³  ·
Jordan A. Guerra³  · Federico Benvenuto⁴  · D. Shaun Bloomfield⁵  ·
Manolis K. Georgoulis¹ 

Received: 13 August 2017 / Accepted: 20 January 2018 / Published online: 30 January 2018
© Springer Science+Business Media B.V., part of Springer Nature 2018

Abstract We propose a forecasting approach for solar flares based on data from Solar Cycle 24, taken by the *Helioseismic and Magnetic Imager* (HMI) on board the *Solar Dynamics Observatory* (SDO) mission. In particular, we use the Space-weather HMI Active Region Patches (SHARP) product that facilitates cut-out magnetograms of solar active regions (AR) in the Sun in near-realtime (NRT), taken over a five-year interval (2012–2016). Our approach utilizes a set of thirteen predictors, which are not included in the SHARP metadata, extracted from line-of-sight and vector photospheric magnetograms. We exploit several machine learning (ML) and conventional statistics techniques to predict flares of peak magnitude $> M1$ and $> C1$ within a 24 h forecast window. The ML methods used are multi-layer perceptrons (MLP), support vector machines (SVM), and random

✉ K. Florios
cflorios@aueb.gr

I. Kontogiannis
jkonto@noa.gr

S.-H. Park
sunpark@tcd.ie

J.A. Guerra
guerraaj@tcd.ie

F. Benvenuto
benvenuto@dim.unige.it

D.S. Bloomfield
shaun.bloomfield@northumbria.ac.uk

M.K. Georgoulis
manolis.georgoulis@academyofathens.gr

¹ Research Center for Astronomy and Applied Mathematics, Academy of Athens, Athens, Greece

² Department of Statistics, Athens University of Economics and Business, Athens, Greece

³ School of Physics, Trinity College Dublin, Dublin, Ireland

⁴ Dipartimento di Matematica, Università di Genova, Genoa, Italy

⁵ Northumbria University, Newcastle upon Tyne, NE1 8ST, UK

forests (RF). We conclude that random forests could be the prediction technique of choice for our sample, with the second-best method being multi-layer perceptrons, subject to an entropy objective function. A Monte Carlo simulation showed that the best-performing method gives accuracy $ACC = 0.93(0.00)$, true skill statistic $TSS = 0.74(0.02)$, and Heideke skill score $HSS = 0.49(0.01)$ for $> M1$ flare prediction with probability threshold 15% and $ACC = 0.84(0.00)$, $TSS = 0.60(0.01)$, and $HSS = 0.59(0.01)$ for $> C1$ flare prediction with probability threshold 35%.

Keywords Flares, forecasting · Flares, relation to magnetic field · Active regions, magnetic fields

1. Introduction

Solar flares are sudden brightenings that occur in the solar atmosphere and release enormous amounts of energy over the entire electromagnetic spectrum. Flares are quite prominent in X-rays, UV, and optical lines (Fletcher *et al.*, 2011) and they are often (but not always) accompanied by eruptions that eject solar coronal plasma into the interplanetary space (coronal mass ejections, CMEs). These very intense phenomena – the largest explosions in the solar system – are associated with regions of enhanced magnetic field, called active regions (AR), and are associated, in white light, with sunspot groups. Depending on their peak X-ray intensity, as recorded by the National Oceanic and Atmospheric Administration’s (NOAA) *Geostationary Operational Environmental Satellite* (GOES) system, flares are categorized in classes, the strongest and most important being X, M, and C (in decreasing order). Flare classification is logarithmic, with a base of 10, and is complemented by decimal sub-classes (M5.0, C3.2 *etc.*).

The solar flare radiation may be detrimental to infrastructures, instruments and personnel in space, therefore flare forecasting is an integral part of contemporary space-weather forecasting. Forecast mainly employs measurements of the AR magnetic field in the solar photosphere. Magnetic-field-based predictors represent AR magnetic complexity or the energy budget available to power flares. Recent developments in instrumentation have led to a regular production of such measurements, offering the opportunity to produce extensive databases with properties suitable for solar flare prediction.

On the other hand, machine learning (ML) in recent years has become an increasingly popular approach for performing computer cognition tasks that were inherently possible only using human intelligence. Thus, ML is a subfield of artificial intelligence (AI), and it aims at using past data in order to train computers so that they can apply the accumulated knowledge to new, previously unseen, data. The acquisition of knowledge is the training phase, and the application of what was learned to future scenarios is the prediction phase. Typically, ML is more interested in prediction than conventional statistics. ML can also interface with conventional statistics in a field called statistical learning (Hastie, Tibshirani, and Friedman, 2009). Learning is either called supervised or unsupervised, depending on whether it is done with a teacher or not. Supervised learning comprises regression and classification, while unsupervised learning is also called clustering. In our study, we focus on classification, where a set of input variables or predictors belongs to one of two classes (binary classification). ML is more powerful than traditional statistical techniques such as, say, generalized linear models that include probit, logit, *etc.* for binary classification, because it can help model more complex nonlinear relationships. An introduction to ML research can be found in several textbooks (MacKay, 2003; Hastie, Tibshirani, and Friedman, 2009).

Several researchers have recently used ML techniques to effectively forecast solar flares. More often, the techniques used by researchers were neural networks (Wang *et al.*, 2008; Yu *et al.*, 2009; Colak and Qahwaji, 2009; Ahmed *et al.*, 2013), support vector machines (Li *et al.*, 2008; Yuan *et al.*, 2010; Bobra and Couvidat, 2015; Boucheron, Al-Ghraibah, and McAteer, 2015), ordinal logistic regression (Song *et al.*, 2009), decision trees (Yu *et al.*, 2009), and relevance vector machines (Al-Ghraibah, Boucheron, and McAteer, 2015). Very recently, random forests have also been used (Barnes *et al.*, 2016; Liu *et al.*, 2017).

We use predictors calculated from near-realtime (NRT) Space-weather HMI Active Region Patches (SHARP) data combined with state-of-the-art ML and statistical algorithms in order to effectively forecast flare events for an arbitrarily chosen 24-hour forecast window. Flare magnitudes of interest are $> M1$ and $> C1$. Prediction is binary, meaning that a given flare class is considered to either happen or not within the next 24 hours after prediction. Our predictions are effective immediately, therefore with zero latency. Analysis involves a comprehensive NRT SHARP sample including all calendar days between 2012 and 2016, at a cadence of 3 hours. Results in this work summarize the findings of the first 18 months of the “Flare Likelihood And Region Eruption foreCASTing” (FLARECAST) project and, while based on ongoing work, we took every effort to present robust and unbiased results.

The contribution of the present work is twofold:

- The use of novel magnetogram-based predictors in a multi-parameter solar flare prediction model.
- The use of classic and novel ML techniques, such as multi-layer perceptrons (MLP), support vector machines (SVM), and especially, for one of the first times,¹ random forests (RF), for the forecasting of $> M1$ and $> C1$ flares.

The application code is available at <https://doi.org/10.17632/4f6z2gf5d6.1>, along with the benchmark dataset used in this work. The run time for all methods is on the order of few minutes.

The analysis presented here is part of the EU Horizon 2020 FLARECAST project, aiming to develop an NRT online forecasting system for solar flares. The study is organized as follows: Section 2 describes the data selected to train and test the algorithms and presents the predictors we used, together with background information on the solar physics aspects of magnetogram-based calculations. Section 3 describes the ML algorithms in terms of their core principles, along with some additional remarks and comments. Section 4 is devoted to the forecast experiments and a comparison with similar published results and statistics. Section 5 presents the main conclusions and future integration of the present work in the FLARECAST operational system. Four appendices that describe multiple complementary aspects of this work are also included.

2. Data and Classification Predictors

2.1. Data

The *Helioseismic and Magnetic Imager* (HMI; Scherrer *et al.*, 2012) on board the *Solar Dynamics Observatory* (SDO; Pesnell, Thompson, and Chamberlin, 2012), provides regular

¹Liu *et al.* (2017) also used the random forest algorithm for solar flare prediction based on SDO/HMI data. Nevertheless, the specific details in that paper regarding the sampling strategy and the feature extraction are very different from our choices. For example, Liu *et al.* (2017) only considered flaring ARs (at the level $> B1$ class), and the sample size was $N = 845$, while we consider both flaring and non-flaring ARs with $N = 23,134$.

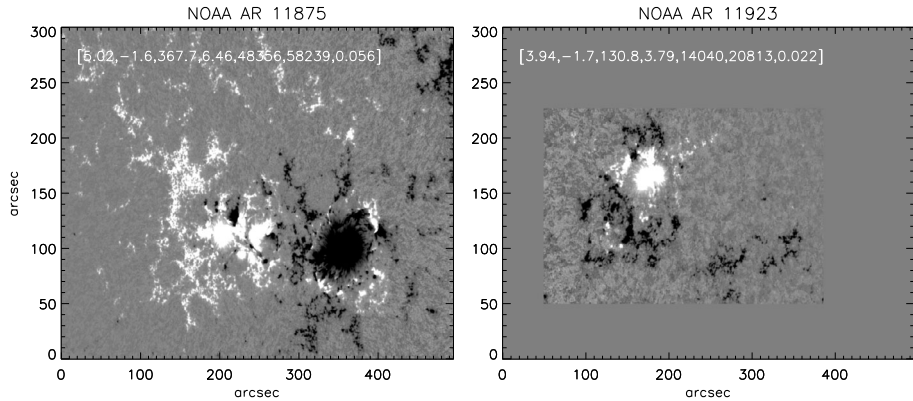


Figure 1 Two SHARP frames depicting an AR with very different levels of flaring activity. NOAA AR 11875 (*left*) produced 7 C-, 0 M-, and 0 X-class flares within 24 h, while NOAA AR 11923 (*right*) produced no flares. The two AR are scaled so as to retain their original relative size and, for comparison, vectors of the seven predictors used are included in the frames. The names of all $K = 7$ predictors [$\log R$, FSPI, TLMPIL, DI, WL_{SG} , $IsinEn1$, $IsinEn2$] are defined in Section 2.2. High values of the predictors statistically indicate a powerful AR (*left*), with low values indicating a quiescent, flare-quiet AR (*right*).

full-disk solar observations of the three components of the photospheric magnetic field. The HMI team has created the Space Weather HMI Active Region Patches (SHARPs), which are cut-outs of solar regions-of-interest along with a set of parameters that might be useful for solar flare prediction (Bobra *et al.*, 2014). For our analysis, we use the near-realtime (NRT), cylindrical equal area (CEA) SHARP data to calculate a set of predictors.

To associate SHARPs with flare occurrence, we use the *Geostationary Operational Environmental Satellite* (GOES) soft X-ray measurements. For each SHARP we search for flares within the next 24 hours by either matching the NOAA AR numbers with those of the recorded flares or by comparing the corresponding longitude and latitude ranges, considering also the differential solar rotation.

The algorithms of Section 3 were tested on a sample of the 2012–2016 SHARP dataset. We considered all days in the period October 1, 2012, to January 13, 2016, and for every given day, we computed the set of predictors (see Section 2.2) at a cadence of 3 hours, starting at 00:00 UT. For our analysis, only SHARP cut-outs that correspond to NOAA ARs were considered. In this way, we obtain a fairly representative sample of the solar activity, including several flares of interest, with a sufficiently high sampling frequency.

2.2. Predictors

The set of 13 predictors consists of both predictors that have previously been proposed in the literature and new ones, and comprises a subset of the parameter set developed for the FLARECAST project. In Figure 1 we show two sample magnetograms to demonstrate how the predictors reflect the complexity and size of the corresponding active region. The predictors used for this study are described below.

2.2.1. Magnetic Polarity Inversion Line (TLMPIL)

A magnetic polarity inversion line (MPIL) in the photosphere of an AR separates distinct patches of positive- and negative-polarity magnetic flux. Several studies have been carried

out to investigate the relationship between flare occurrence and MPIL characteristics (Schrijver, 2007; Falconer *et al.*, 2012). We determined a specific subset of an MPIL that has been also identified as an MPIL*, with i) a strong gradient in the vertical component of the field across the MPIL, and ii) a strong horizontal component of the field around the MPIL. MPIL* has been considered as the single most likely place in an AR where potential magnetic instabilities, such as, say, magnetic flux cancellation and/or magnetic flux rope formation (Fang *et al.*, 2012) can take place. Such processes seem intimately related to flares. We used the total length L_{tot} of MPIL* segments in active regions as an MPIL quantification parameter.

2.2.2. Decay Index (DI)

The decay index is a quantitative measure for the torus magnetic instability in a current-carrying magnetic flux rope (Kliem and Török, 2006). It has been found that the higher the value of the decay index in AR magnetic fields, the more likely a solar eruption involving a major solar flare (Zuccarello, Aulanier, and Gilchrist, 2015). We developed a decay index parameter derived by the ratio $L_{\text{hs}}/h_{\text{min}}$, where L_{hs} is the length of a highly sheared portion of an MPIL and h_{min} is the minimum height at which the decay index achieves a purported critical value of 1.5. This ratio can be used to measure the degree of instability in a flux rope. Note that if there was more than one MPIL in an AR, then we calculated the ratio $L_{\text{hs}}/h_{\text{min}}$ for every MPIL and took the peak value for a given time that represents the highest eruptive potential of the AR.

2.2.3. Gradient-Weighted Integral Length of the Neutral Line (WL_{SG})

The gradient-weighted integral length of the neutral line, WL_{SG} , is defined in Falconer, Moore, and Gary (2008) as

$$WL_{\text{SG}} = \int (\nabla B_z) dl, \quad (1)$$

and corresponds to the line integral of the vertical-field (B_z) horizontal gradient over all neutral line (or MPIL) segments on which the potential horizontal field is greater than 150 G. This MPIL-related property has been reported to show a useful empirical association with the occurrence of solar eruptions (flares, CMEs, SPEs; Falconer *et al.*, 2011, 2014) and is the main predictor used in the Magnetic Forecast (MAG4) forecasting service, developed in the University of Alabama (<http://www.uah.edu/cspar/research/mag4-page>).

For these calculations of WL_{SG} , two approximations of the vertical field B_z are used: B_{los} (line of sight; uncorrected) and B_r , keeping in mind that in the former case, only values for regions located within 30° from the central meridian are considered accurate. For each magnetogram, an MPIL mask is determined as in the calculation of MPIL characteristics, described previously. In order to select the strong-horizontal field segments of MPILs, the potential field extrapolation method developed by Alissandrakis (1981) is used. Finally, the horizontal gradient of B_z is calculated numerically and integrated over all MPIL segments. The accuracy of the calculated values was estimated by comparing flare rates derived from our calculations of WL_{SG} (using Equation 4 along with the values in Table 1 in Falconer *et al.*, 2011) with the flare rates from the text output of MAG4.

2.2.4. Ising Energy ($IsinEn1$, $IsinEn2$)

The Ising energy is a quantity that parameterizes the magnetic complexity of an AR (Ahmed *et al.*, 2010). For a two-dimensional distribution of positive and negative interacting mag-

netic elements, the Ising energy is defined as

$$E_{\text{Ising}} = - \sum_{ij} \frac{S_i S_j}{d^2}, \quad (2)$$

where S_i (S_j) equals +1 (−1) for positive (negative) pixels and d is the distance between opposite-polarity pairs. The interacting magnetic elements can be either the individual pixels with a minimum flux density value as in Ahmed *et al.* (2010) or the opposite-polarity partitions, produced using a flux-partitioning scheme (Barnes, Longcope, and Leka, 2005). The latter variation has been introduced for the first time in the FLARECAST project with promising results, and an assessment of its merit as a predictor is underway (Kontogiannis *et al.*, in preparation). The Ising energy calculation produces four predictors, two for the line-of-sight magnetic field, and two for the radial magnetic field component.

2.2.5. Fourier Spectral Power Index (FSPI)

The spectral power index, α , corresponds to the power-law exponent in fitting the one-dimensional power spectral density $E(k)$ extracted from magnetograms by the relation

$$E(k) \sim k^{-\alpha}. \quad (3)$$

This index parameterizes the power contained in magnetic structures of spatial scales l ($= k^{-1}$) belonging to the inertial range of magnetohydrodynamic (MHD) turbulence. Empirically, AR with spectral power index higher than 5/3 (Kolmogorov's exponent for turbulence) are thought to display an overall high productivity of flares (*e.g.* see Guerra *et al.*, 2015).

The spectral power index has been historically calculated from the vertical component of the photospheric magnetic field, as inferred from the line-of-sight component assuming perfectly radial magnetic fields. First, the magnetogram is processed using the fast Fourier transform (FFT). A two-dimensional power spectral density (PSD) is then obtained as

$$E(k_x, k_y) = |\text{FFT}[B(x, y)]|^2. \quad (4)$$

In order to express $E(k_x, k_y)$ from the Fourier k_x and k_y to the isotropic wavenumber $k = (k_x^2 + k_y^2)^{1/2}$, it is necessary to calculate $E(k)'$ – the integrated PSD over angular direction in Fourier space. From this last step, the one-dimensional PSD is obtained as $E(k) = 2\pi k E(k)'$. Finally, the power-law fit is performed as a linear fit in a logarithmic representation of $E(k)$ versus k and α is measured for the assumed turbulent inertial range of 2–20 Mm (*i.e.* 0.05–0.5 Mm^{−1}).

2.2.6. Schrijver's R Value (logR)

The R -value property quantifies the unsigned photospheric magnetic flux near strong MPILs. The presence of such MPILs indicates that twisted magnetic structures carrying electrical currents have emerged into the AR through the solar surface. Therefore, R represents a proxy for the maximum free magnetic energy that is available for release in a flare. This property and its usefulness in forecasting was first investigated by Schrijver (2007).

The algorithm for calculating R is relatively simple, computationally inexpensive, and was originally developed to use line-of-sight magnetograms from the *Michelson Doppler*

Imager (MDI) (Scherrer *et al.*, 1995) on board the *Solar and Heliospheric Observatory* (SoHO). First, a bitmap is constructed for each polarity in a magnetogram, indicating where the magnitude of positive and negative magnetic flux densities exceeds the threshold value of $\pm 150 \text{ Mx cm}^{-2}$. These bitmaps are then dilated by a square kernel of 3×3 pixels, and the areas where the bitmaps overlap are defined as strong-field MPILs. This combined bitmap is then convolved with a Gaussian filter with a full-width at half-maximum (FWHM) of $\approx 15 \text{ Mm}$. This particular value is constrained by how far from MPILs flares are observed to occur in extreme ultraviolet images of the solar corona. Finally, the convolved bitmap is multiplied by the absolute flux value of the line-of-sight magnetogram, and R is calculated as the sum over all pixels. Note that since the R value was implemented by Schrijver (2007) for MDI magnetograms, the SHARP magnetograms were resampled to the spatial scale of MDI before the kernel application and subsequent calculations.

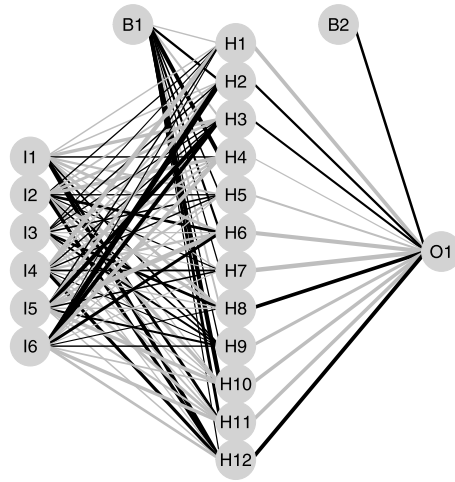
3. Machine-learning Algorithms and Conventional Statistics Models

The ML algorithms used in this study are MLPs, SVMs, and RFs. Among the hundreds of ML algorithms proposed for binary classification (*e.g.* Fernández-Delgado *et al.*, 2014), these three categories of algorithms are representative of three important approaches in ML: i) artificial neural networks (ANN), ii) kernel-based methods, and iii) classification and regression trees. This is the reason why they were used in the present study, in order to furthermore investigate whether the usage of RFs could bring any improvements in flare prediction in comparison to SVMs and MLPs. The RFs belong to the category of ensemble methods, while the MLPs use unconstrained optimization and SVMs use constrained optimization techniques (*e.g.* quadratic programming). In general, the working principle of ML comprises the following steps: i) train the model using a training set, ii) predict using the trained model and a testing set, and iii) check whether the algorithm predicted well, in what is called the validation of the overall ML procedure. For further study, we refer to Vapnik (1998), MacKay (2003), and Hastie, Tibshirani, and Friedman (2009).

3.1. Multi-layer Perceptrons

The MLP is a feed-forward network, thus it is described by the planar graph shown in Figure 2. It contains an input layer, a hidden layer, and an output layer of neurons. By the term neuron we denote a basic processing unit where inputs are summed using specific weights and the result is squashed *via* an activation function. The hidden layer might expand in a series of hidden layers. Nevertheless, the simplest MLP networks have just one hidden layer. In principle, the term hidden describes every layer that is neither the input nor the output layer, but resides in between, as presented in Figure 2. A sufficient number of hidden nodes allows the MLP to approximate any continuous nonlinear function of several inputs with a desired degree of accuracy (Hornik, Stinchcombe, and White, 1989), which is what characterizes the MLPs as universal approximators. It also holds that the greater the number of hidden nodes, the more complex the nonlinear function that can be approximated by the neural network with a desired degree of accuracy. The number of hidden nodes typically does not have to be more than twice the number of input nodes (or predictors). If too many hidden nodes are used, then the overfitting problem arises, which means that the MLP memorizes the sample observations and generalizes badly in the prediction phase. Usually, and in this study, the optimal number of hidden neurons (called size of the MLP) is determined with a fine-tuning procedure (*e.g.* cross-validation approach, see Section 4.2) before the training

Figure 2 Example MLP neural network with 6 inputs, 12 hidden nodes, 1 output, and 2 biases. **Bold, darker lines** indicate large positive weights ω .



phase starts. The tuning phase is relatively time consuming, so it need not be executed every time the training starts. It can be conducted for a single realization of the training set.

An MLP network is a kind of a nonlinear regression (classification) technique, equivalent to a nonlinear mapping from input I to an output $O = O(I; \omega, A)$. The output is a continuous function of the input and of the weights ω . The network is described by a given architecture A , which typically defines the number of nodes in every layer (*e.g.* input, hidden, and output). In general, MLP networks can be used to solve regression and classification problems. The statistical model of an MLP neural network for binary outcome, as described in the following, is based on MacKay (2003). For a recent survey on neural networks, we refer to Prieto *et al.* (2016).

3.1.1. Classification Networks

We consider an MLP with l inputs called I_l and bias B_1 . The network also contains a single hidden layer with j hidden nodes H_j and bias B_2 . We have in general i outputs O_i , while typically a single output is all that is needed ($i = 1$).

In the case of a classification problem, the propagation of the information from the inputs I to the output O is described by

$$\alpha_j^{(1)} = \sum_{l=1}^L \omega_{jl}^{(1)} I_l + B_j^{(1)}; \quad H_j = f(\alpha_j^{(1)}),$$

$$\alpha_i^{(2)} = \sum_{j=1}^J \omega_{ij}^{(2)} H_j + B_i^{(2)}; \quad O_i = g(\alpha_i^{(2)}),$$
(5)

where, for example, $f(\alpha) = \frac{1}{1+\exp(-\alpha)}$ and $g(\alpha) = \frac{1}{1+\exp(-\alpha)}$.

The index l is used for the inputs I_1, \dots, I_L , the index j is used for the hidden units, and the index i is used for the outputs ($i = 1$). The weights $\omega_{jl}^{(1)}$, $\omega_{ij}^{(2)}$, and biases $B_j^{(1)}$ and $B_i^{(2)}$ define the parameter vector ω to be estimated. The nonlinear logistic function f at the hidden layer (also known as activation function) helps the neural network approximate any generic continuous nonlinear function with a desirable degree of accuracy (Hornik,

Stinchcombe, and White, 1989). Visually, a neural network can be represented as a series of layers consisting of nodes, where every node is connected to nodes of the subsequent layer only (feed-forward networks).

In the case of binary classification, the MLP is trained using a dataset of examples $D = \{\mathbf{I}^{(n)}, \mathbf{T}^{(n)}\}$ by adjusting ω in order to minimize $G(\omega)$, the negative log-likelihood function,

$$G(\omega) = -\left(\sum_{n=1}^N \mathbf{T}^{(n)} \ln(\mathbf{O}(\mathbf{I}^{(n)}; \omega)) + (1 - \mathbf{T}^{(n)}) \ln(1 - \mathbf{O}(\mathbf{I}^{(n)}; \omega))\right). \tag{6}$$

Note that $\mathbf{I}^{(n)}$ is the matrix of the predictors and $\mathbf{T}^{(n)}$ is the vector of the targets for observation $n = 1, \dots, N$. In Equation 6, $\mathbf{T}^{(n)}$ is 0 (1) for the negative (positive) class, respectively, and $\mathbf{O}(\mathbf{I}^{(n)}; \omega)$ is strictly between 0 and 1 (a probability); this is ensured by Equations 5.

3.2. Support Vector Machines

The SVM variant we use is the *C*-Support Vector Classification (*C*-SVC) according to the widely used library LIBSVM (Chang and Lin, 2011; Meyer, Leisch, and Hornik, 2003).

Let us assume a vector of K predictor values at observation i , $\mathbf{x}_i \in R^K$, $i = 1, \dots, N$, which belongs in one of two classes, and an indicator vector $\mathbf{y} \in R^N$ such that $y_i \in \{1, -1\}$. Note that the positive class has the label $+1$ and the negative class has the label -1 . Then the *C*-SVC solves the optimization problem

$$\begin{aligned} &\text{minimize } \frac{1}{2} \omega^T \omega + C \sum_{i=1}^N \xi_i, \\ &\text{subject to} \\ & y_i (\omega^T \phi(\mathbf{x}_i) + b) \geq 1 - \xi_i, \quad i = 1, 2, \dots, N, \\ & \xi_i \geq 0, \quad i = 1, 2, \dots, N, \end{aligned} \tag{7}$$

where $\phi(\mathbf{x}_i)$ is an arbitrary unknown function that maps \mathbf{x}_i into a higher dimensional space, and $C > 0$ is the regularization parameter. The optimization in *C*-SVC model is performed by changing the decision variables ω , b , and ξ . LIBSVM solves the dual of *C*-SVC, which depends on a quantity $K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$, which is called the “kernel” function. While $\phi(\mathbf{x}_i)$ is unknown, the kernel function is known and is equal to the inner product of $\phi(\mathbf{x}_i)$ with itself, but for different pairs of observations i and j . This is the so-called kernel trick of the SVMs. As we show below, the kernel is a similarity measure and takes the maximum value of 1 when $\text{dist}(\mathbf{x}_i, \mathbf{x}_j) = 0$.

We used the radial basis function (RBF; or Gaussian) kernel, which is defined as $K(\mathbf{x}, \mathbf{x}') = \exp(-\gamma \|\mathbf{x} - \mathbf{x}'\|^2)$. A variant of the *C*-SVC model has been used for flare prediction in Bobra and Couvidat (2015).

For imbalanced datasets that account for rare events (e.g. in our case the $> M1$ flares), some researchers, e.g. Bobra and Couvidat (2015), have used two different values for the regularization parameter C in Equation 7, thereby penalizing the constraint violations for the minority class more strongly. These authors have used C_1 and C_2 with a ratio $C_2/C_1 \in [2, 15]$, where C_1 is the coefficient for the majority class (no events) and C_2 is the coefficient for the minority class (events). While we generally use the SVM in the original unweighted version in Equation 7, in auxiliary runs we also experimented with using different values for C_1 and C_2 with a ratio $C_2/C_1 \in [2, 15, 20]$ to account for the imbalanced nature of the $> M1$ flares dataset.

3.3. Random Forests

The RF is a relatively recent ML method and was introduced by Breiman (2001). The RF approach is an ensemble of tree predictors, where we let each tree vote for the most popular class. It has been reported (Fernández-Delgado *et al.*, 2014) that RF offers significant performance improvement over other classification algorithms. The RF approach relies on randomness and involves the concept of split purity and the Gini index for variable selection (Breiman *et al.*, 1984).

According to Hastie, Tibshirani, and Friedman (2009), the goal of the RF algorithm is to randomly build a set (or ensemble) of trees by repeating the tree-formation process B times to create B trees. In particular, the algorithm i) chooses a bootstrap sample from the training data, ii) grows a tree T_b to the bootstrapped sample by consequently applying the following two substeps: Substep 1 selects m variables randomly out of the M variables, and Substep 2 splits the current node into two children nodes after selecting the best variable (node) from the m chosen ones. By repeating steps i) and ii) (where ii) consists of Substeps 1–2), the algorithm creates a set (called ensemble) of trees $\{T_b\}_1^B$. Then, in the classification case studied in the present paper, a voting procedure for every tree T_b is followed in order to obtain the class prediction of the random forest.

This is one of the first times that RF is used for flare forecasting. Other related works are Liu *et al.* (2017) and Barnes *et al.* (2016). Furthermore, three recent applications of RF in astrophysics have been reported by (Vilalta, Gupta, and Macri, 2013; Schuh, Angrzyk, and Martens, 2015; Granett, 2017).

3.4. Implementation of ML Algorithms

3.4.1. Multi-layer Perceptrons

The MLPs were implemented using the R programming language and the `nnet` package (Venables and Ripley, 2002). The options used were `linout = FALSE`, to ensure that sigmoid activation functions are used at the output node; `entropy = TRUE`, to ensure that the negative log-likelihood objective function is minimized during the training phase (and not the default sum of squares error (SSE) criterion); and `size = iNode`, where `iNode` for both $> M1$ flares and for $> C1$ flares was chosen with a tuning procedure.

3.4.2. Support Vector Machines

Support vector machines were implemented using the R programming language and the `e1071` package (Meyer *et al.*, 2015). The option used was `probability = TRUE`, in order to obtain probability estimates for every element of the training set as well as probability estimates for every element of the testing set.

3.4.3. Random Forests

Random forests were implemented using `randomForest` package (Liaw and Wiener, 2002) in the R programming language. The options used were `importance = TRUE`, to create importance information for every predictor; `na.action = na.omit`, to exclude records of predictors with missing values appearing in preliminary versions of the dataset (but lacking from the final version of the dataset).

3.5. Conventional Statistics Models

Non-ML (or statistical) methods also considered are i) linear regression (LM), ii) probit regression (PR), and iii) logit regression (LG). Although multiple linear regression is known to be redundant for binary outcomes because it can yield probabilistic predictions outside the interval $[0, 1]$, we still included it in the array of tested methods. The reason is that some practitioners still use it for binary outcomes (calling it linear probability model (LPM), see Greene, 2002) and there is always interest to consider ordinary least squares (OLS) as an entry-level method for any regression analysis. An interesting article about the lack of use of probit and logit in astrophysics modeling is de Souza *et al.* (2015). The statistical algorithms were implemented in the statistical programming language R using the `lm` and `glm` functions.

For a description of these well-known methods we refer to (Greene, 2002; Winkelmann and Boes, 2006).

4. Data Preparation, Results, and Discussion

First, we implement ML predictions on $> M1$ flares. Second, we use statistical methods for the prediction of $> M1$ flares. Third, we predict $> C1$ flares with ML algorithms. Finally, we predict $> C1$ flares with the statistical algorithms. The following subsections describe these four experiments, presenting at first a single combination of a training/testing set for every flare class and category of techniques.

Results are presented for the prediction step in terms of i) skill score profiles (SSP) of ACC, TSS, and HSS as functions of the probability threshold, ii) ROC curves, and iii) RD plots for all methods: (for the explanation of metrics ACC, TSS, HSS, and ROC curves and RD diagrams, see Section 4.3). Skill score profiles were created by a code we developed in R, ROC curves were created using the `ROCR` package (Sing *et al.*, 2005), and reliability diagrams were created using the `verification` package (Laboratory, 2015).

All algorithms were implemented and run using the R programming language 3.3.2 (R Core Team, 2016) and the RStudio 0.99 IDE.

4.1. Data Pre-processing

The data comprise the $K = 7$ predictors [$\log R$, FSPI, TLMPI, DI, WL_{SG} , $IsinEn1$, and $IsinEn2$] described in Section 2.2 and computed using either the line-of-sight magnetograms, B_{los} , of SHARP data or the respective radial component, B_r (Bobra *et al.*, 2014). Hence, we tested $K = 2 \times 6 + 1 = 13$ predictors.² The sample comprised $N = 23,134$ observations, randomly split in half into $N_1 = 11,567$ observations for the training and $N_2 = 11,567$ observations for the testing set. The random split was performed for 200 replications, and all six prediction algorithms (*i.e.* MLP, SVM, RF, LM, probit, and logit) of Section 3 were trained and performed on identical training and test sets. The metrics ACC, TSS, and HSS of Section 4.3 were always computed for the testing (out-of-sample) set. We standardized all predictor variables to have a mean equal to 0 and a standard deviation equal to 1 because several ML algorithms involve non-linear optimization (*e.g.* MLPs). This helps to better train the ML algorithms and also explains the effect of every predictor variable on the studied outcome in the case of the statistical models LM, probit, and logit.

²This is because we considered only the B_r version for predictor WL_{SG} .

4.2. Tuning of ML Algorithms

As with any parameterized algorithm (*e.g.* simulated annealing, evolutionary algorithms, and other metaheuristics), the performance of ML algorithms depends on a number of crucial parameters that need to be fine-tuned before the application of the ML procedure (*e.g.* training, testing, and validation steps). The optimal tuning of ML algorithms is more or less still an open question in the ML community and always poses a great challenge for any practitioner. This choice of optimal options for the ML algorithms themselves is similar to the choice of optimal parameters for other numerical models (*e.g.* MHD models), where the analyst also has to explore the optimal parameter space in several crucial parameters before conducting numerical MHD simulations. The algorithms MLP, SVM, and RF have their critical hyperparameters (*e.g.* parameters that are critical for the forecasting performance of every algorithm) tuned *via* a 10-fold cross-validation study exploiting only the training set at one of its realizations. The set of plausible values for every ML algorithm is as follows: i) MLP: size (number of hidden neurons) $\in \{4, 13, 26\}$ and decay (weight decay parameter) $\in \{10^{-3}, 10^{-2}, 10^{-1}\}$, ii) SVM: γ (parameter in the RBF (or Gaussian) kernel) $\in \{10^{-6}, 10^{-5}, 10^{-4}, \dots, 10^{-1}\}$ and cost (regularization parameter) $\in \{10, 100\}$, and iii) RF: mtry (number of variables randomly sampled as candidates at each split) $\in \{\lfloor \sqrt{K} \rfloor = 3\}$ and ntree (number of trees to grow) $\in \{500\}$.

We tuned only the MLP and SVM classifiers because the default RF values mtry = 3 and ntree = 500 immediately provided satisfactory results. Tuning of the MLP and SVM was mostly needed in the > M1 flares case, which was found harder to predict than > C1 flares, but was also performed in the > C1 flares case. Thus, the hyperparameters for MLP and SVM needed tuning because, for example, the default values $\gamma = 1$ and cost = 1 for SVM provided unsatisfactory results. We used the `tune.nnet` and `tune.svm` functions of the R package `e1071` to tune the MLP and SVM, respectively. After the tuning, both MLP and SVM improved their performance significantly.

For the > M1 flares, the selected values are size = 26 and decay = 0.1 for the MLP and $\gamma = 0.1$ and cost = 10 for the SVM. These values are used throughout the remainder of this work. For the > C1 flares case, the selected values are size = 4 and decay = 0.1 for the MLP and $\gamma = 0.001$ and cost = 100 for the SVM.

4.3. Comparison Metrics

A wide variety of metrics exist in order to characterize the quality of binary classification. Of these, no single one is fit for all purposes. There exist two types of metrics, suitable for either categorical or probabilistic classification. In the former case, a strict class membership is returned from the model, and in the latter case, a probability of membership is returned. In this section we concentrate on categorical forecast metrics for binary classification. In what follows, let ACC denote accuracy, TSS denote true skill statistic, and HSS denote Heidke skill score. The performance of algorithms is measured using a number of metrics. These are derived from the so-called contingency table or confusion matrix, a representation of which is provided in Table 1.

Table 1 includes true positives (TP; events predicted and observed), true negatives (TN; events not predicted and not observed), false positives (FP; events predicted but not observed), and false negatives (FN; events not predicted but observed), where $N = TP + FP + FN + TN$ is the sample size. From these elements, the meaning of ACC is the proportion correct, namely the number of correct forecasts of both event and non-event, normalized by

Table 1 2×2 contingency table for binary forecasting.

		ACTUAL	
PREDICT	NO	NO	YES
	YES	TN	FN
	YES	FP	TP

the total sample size,

$$ACC = \frac{TP + TN}{N}. \tag{8}$$

The TSS (Hanssen and Kuipers, 1965) compares the probability of detection (POD) to the probability of false detection (POFD),

$$TSS = POD - POFD = \frac{TP}{TP + FN} - \frac{FP}{FP + TN}. \tag{9}$$

Moreover, the TSS is the maximum vertical distance from the diagonal in the ROC curve, which relates the POD and POFD for different probability thresholds; see Section 4. The TSS covers the range from -1 up to $+1$, while the value of zero indicates lack of skill. Values below zero are linked to forecasts behaving in a contrary way, namely mixing the role of the positive class with the role of the negative class. In any negative TSS value, by exchanging the roles of YES and NO events, we can obtain the corresponding positive TSS value that would be identical in absolute value terms with the negative TSS value.

The HSS (Heidke, 1926) measures the fractional improvement of the forecast over the random forecast,

$$HSS = \frac{2(TP \times TN - FP \times FN)}{(TP + FN)(FN + TN) + (TP + FP)(FP + TN)}, \tag{10}$$

which ranges from $-\infty$ to 1. Any negative value means that the random forecast is better, a zero value means that the method has no skill over the random forecast, and an ideal forecast method provides an HSS value equal to 1.

The TSS and HSS metrics are among the most popular metrics for comparison purposes in Meteorology and Space Weather and were conceptually compared in Bloomfield *et al.* (2012). In a probabilistic forecasting, such as the one for solar flares, they must be assigned a probability threshold, thus appearing as functions of this threshold.

To summarize, ACC is the most popular classification metric, but in rare events such as flares $> M1$, the ACC can be artificially high for the naive model, which will always predict the majority class (“no event”). Thus, TSS and HSS are more suitable for flare prediction. Moreover, TSS has the advantage of being invariant to the frequency of events in a sample (*e.g.* see Bloomfield *et al.*, 2012). Typically, both TSS and HSS need to be evaluated for a given probability threshold in order to assess the merit of a given probabilistic forecasting model, such as those we develop in this study.

Regarding the probabilistic assessment of classifiers, we used the visual approaches of receiver operating characteristic (ROC) curves and reliability diagrams (RD) (*e.g.* see Section 4). The ROC describes the relationship between the POD and the POFD for different probability thresholds (*e.g.* see Figure 3b). The area under the curve (AUC) in the ROC has an ideal value of one. The RD describes the relationship between the returned probabilities by the model and the actual observed frequencies of the data. A binning approach is used

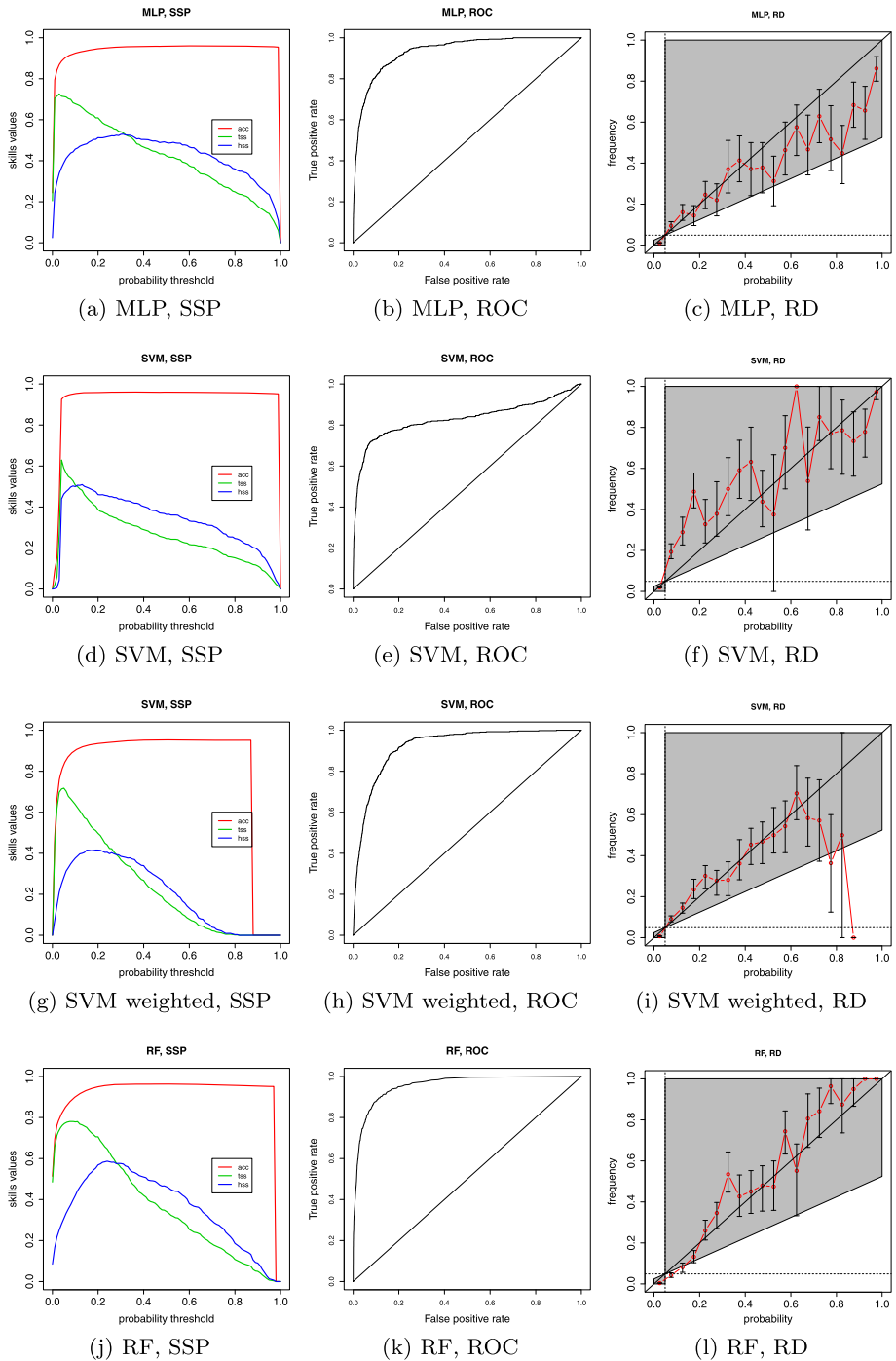


Figure 3 ML method comparison for $> M1$ GOES flare prediction for (from *top to bottom*) MLP, SVM, weighted SVM, and RF. From *left to right*, we present the corresponding SSP, ROC, and RD.

to construct the RD, in which probabilities are assigned to intervals of arbitrary length (for example, we used 20 bins of length 0.05 each). For an example of RD, see Figure 3c. To algebraically assess the probabilistic performance of classifiers, we also used the Brier score (BS) (Brier, 1950) and Brier skill score (BSS) (Wilks, 2011), as well as the AUC (Marzban, 2004).

4.4. Results on > M1 Flare Prediction

4.4.1. Prediction of > M1 Flare Events Using Machine Learning

Figure 3 shows the forecast performances of the three tested ML methods, using both binary scores (SSP [left]; ROC [middle]) and probabilistic ones (RD [right]).

- i) Regarding the MLPs, we note a wide plateau with a more or less flat profile for HSS and less so for TSS. This occurs because the number of hidden neurons (size = 26) is twice the number of input neurons, causing the MLP to provide probability estimates clustered around 0 and 1. The ROC curve is reasonably good, with maximum TSS = 0.726. Moreover, the RD shows a systematic overprediction above a forecast probability of 0.4.
- ii) For the SVMs, the SSP plateau noted in case of the MLPs is not present here, with nearly monotonically decreasing values of TSS and HSS appearing. The ROC curve shows a maximum TSS = 0.629, while the RD seems slightly better than for MLP, with some underprediction below a forecast probability of 0.4 and generally large uncertainties. When we used the weighted version of the SVM, with a ratio of $C_2/C_1 = 20$, the ROC curve improved, providing a maximum TSS = 0.718, but the overall forecasting ability as measured by the SSP and RD remained worse than the MLP.
- iii) With respect to the RFs, the SSP behavior is such that HSS shows a plateau around its peak value, although smaller than in case of MLPs, while TSS monotonically decreases. This said, we note that the peak HSS and TSS values are higher in this case (e.g. TSS = 0.780 and HSS = 0.587). The ROC curve is better than that of MLPs and SVMs, with a maximum TSS = 0.780. The RD, finally, appears clearly better than those of MLPs and SVMs, presenting some mild underprediction, mainly within error bars, above a forecast probability of 0.2.

4.4.2. Prediction of > M1 Flare Events Using Statistical Models

Figure 4 shows the forecast performances of the three tested statistical methods for > M1 flare prediction.

For the LM, the SSP is different between TSS and HSS, with TSS peaking more impulsively and for lower probabilities and then decreasing nearly monotonically. The ROC curve also shows a significant performance with maximum TSS = 0.744 that can also be seen in the RD, which shows a very good behavior, although with error bars, for the entire range of forecast probabilities.

The PR shows a slightly improved behavior in comparison with LM for the SSP, the ROC curves, and the RD. The RD seems also more reliable in this case compared to LM, although differences are mostly within the error bars.

We note a similar behavior for the LG as in the LM and especially PR method, and the RD in this case appears as good as the PR RD.

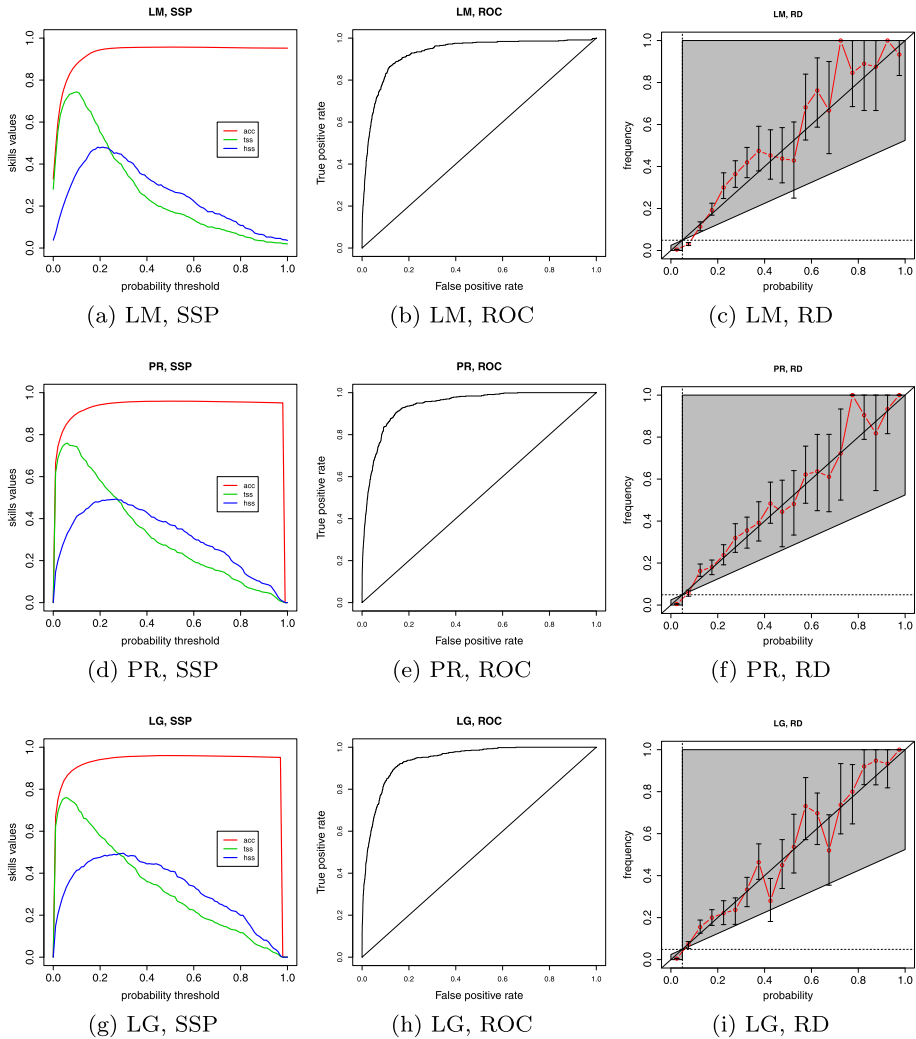


Figure 4 Same as Figure 3, but for statistical methods: linear regression (LM; *top*), probit regression (PR; *middle*), and logit regression (LG; *bottom*).

4.4.3. Monte Carlo Simulation for > M1 Flares

In Table 2 we provide the average values of the skill scores ACC, TSS, and HSS for all prediction methods after the 200 replications of the Monte Carlo experiment regarding > M1 flare prediction. Table 2 shows that the maximum HSS = 0.57 is obtained with the RF method for a probability threshold of 25%. The corresponding RF score values are ACC = 0.96 ± 0.00 , TSS = 0.63 ± 0.02 , and HSS = 0.57 ± 0.02 . The second-best method in Table 2 for the same probability threshold is MLP, with ACC = 0.95 ± 0.00 , TSS = 0.56 ± 0.02 , and HSS = 0.50 ± 0.02 . For the threshold where the maximum TSS is observed, we obtain the best results for the RF method and a threshold of 10%, with values ACC = 0.90 ± 0.00 , TSS = 0.77 ± 0.01 , and HSS = 0.42 ± 0.01 . The second-best method

Table 2 Monte Carlo scenario 1, based on 200 SHARP datasets, on > M1 GOES flare prediction. Numbers in boldface correspond to the most significant results of a given method (MLP: multi-layer perceptron; LM: linear regression; PR: probit regression; LG: logit regression; RF: random forest; SVM: support vector machine).

Par	%	MLP			LM			PR			LG			RF			SVM		
		ACC	TSS	HSS	ACC	TSS	HSS	ACC	TSS	HSS	ACC	TSS	HSS	ACC	TSS	HSS	ACC	TSS	HSS
val ₀	0.00	0.15	0.11	0.01	0.32	0.27	0.04	0.00	0.00	0.00	0.00	0.00	0.00	0.51	0.48	0.08	0.00	0.00	0.00
val ₅	0.05	0.90	0.70	0.39	0.78	0.70	0.22	0.84	0.75	0.30	0.85	0.75	0.31	0.85	0.77	0.32	0.94	0.59	0.47
val ₁₀	0.10	0.93	0.66	0.45	0.88	0.73	0.35	0.90	0.71	0.39	0.90	0.69	0.40	0.90	0.77	0.42	0.95	0.51	0.49
val ₁₅	0.15	0.94	0.62	0.48	0.92	0.65	0.43	0.93	0.64	0.45	0.93	0.63	0.45	0.93	0.74	0.49	0.96	0.46	0.50
val ₂₀	0.20	0.95	0.59	0.50	0.94	0.54	0.46	0.94	0.58	0.48	0.94	0.58	0.48	0.95	0.69	0.54	0.96	0.42	0.48
val ₂₅	0.25	0.95	0.56	0.50	0.95	0.45	0.45	0.95	0.52	0.49	0.95	0.52	0.49	0.96	0.63	0.57	0.96	0.39	0.47
val ₃₀	0.30	0.95	0.53	0.50	0.95	0.38	0.44	0.96	0.47	0.49	0.96	0.48	0.50	0.96	0.57	0.57	0.96	0.37	0.46
val ₃₅	0.35	0.96	0.50	0.50	0.96	0.31	0.39	0.96	0.41	0.47	0.96	0.43	0.49	0.96	0.51	0.57	0.96	0.35	0.45
val ₄₀	0.40	0.96	0.47	0.50	0.96	0.26	0.35	0.96	0.36	0.45	0.96	0.39	0.47	0.97	0.46	0.55	0.96	0.33	0.44
val ₄₅	0.45	0.96	0.44	0.49	0.96	0.21	0.31	0.96	0.32	0.42	0.96	0.35	0.45	0.97	0.41	0.52	0.96	0.31	0.42
val ₅₀	0.50	0.96	0.42	0.48	0.96	0.18	0.28	0.96	0.28	0.39	0.96	0.31	0.42	0.97	0.37	0.49	0.96	0.29	0.41
val ₅₅	0.55	0.96	0.39	0.47	0.96	0.16	0.25	0.96	0.25	0.36	0.96	0.27	0.39	0.96	0.32	0.45	0.96	0.28	0.39
val ₆₀	0.60	0.96	0.37	0.45	0.96	0.14	0.23	0.96	0.21	0.33	0.96	0.24	0.36	0.96	0.28	0.41	0.96	0.26	0.38
val ₆₅	0.65	0.96	0.34	0.44	0.96	0.11	0.19	0.96	0.18	0.29	0.96	0.21	0.32	0.96	0.24	0.36	0.96	0.24	0.36
val ₇₀	0.70	0.96	0.32	0.42	0.96	0.09	0.16	0.96	0.16	0.25	0.96	0.18	0.28	0.96	0.19	0.31	0.96	0.22	0.34
val ₇₅	0.75	0.96	0.29	0.40	0.96	0.08	0.13	0.96	0.13	0.22	0.96	0.15	0.25	0.96	0.15	0.25	0.96	0.21	0.32
val ₈₀	0.80	0.96	0.27	0.37	0.95	0.06	0.12	0.96	0.11	0.18	0.96	0.12	0.21	0.96	0.11	0.19	0.96	0.18	0.29
val ₈₅	0.85	0.96	0.24	0.34	0.95	0.05	0.09	0.96	0.08	0.14	0.96	0.09	0.16	0.96	0.08	0.14	0.96	0.16	0.26
val ₉₀	0.90	0.96	0.20	0.31	0.95	0.03	0.06	0.95	0.05	0.10	0.95	0.06	0.10	0.95	0.04	0.08	0.96	0.13	0.22
val ₉₅	0.95	0.96	0.15	0.25	0.95	0.02	0.04	0.95	0.03	0.06	0.95	0.03	0.05	0.92	0.01	0.02	0.96	0.08	0.14
val ₁₀₀	1.00	0.00	0.00	0.00	0.95	0.02	0.03	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00

may be considered the LM at 10% threshold, with $ACC = 0.88 \pm 0.00$, $TSS = 0.73 \pm 0.01$, and $HSS = 0.35 \pm 0.01$. The difference between RF and LM is statistically significant at the 0.01% level, as shown in Table 4 in row 1. For the range of thresholds 10% to 25%, the RF method yields increasing values of HSS and decreasing values of TSS. For example, an appealing forecasting model could be RF with threshold 15% and metrics $ACC = 0.93 \pm 0.00$, $TSS = 0.74 \pm 0.02$, and $HSS = 0.49 \pm 0.01$ in Table 2, but this would depend on the needs and requirements of a given decision maker.

4.5. Results on > C1 Flare Prediction

4.5.1. Prediction of > C1 Flare Events Using Machine Learning

We continued our computational experiments by training and performing our algorithms to the prediction of GOES > C1 flares. Figure 5 shows the forecast performances of the three tested ML methods for > C1 flare prediction.

For the MLP, we note that since for the > C1 flares the number of hidden nodes selected is size = 4, plateaus in HSS and TSS are not so eminent, in contrast to the case of > M1 flare prediction. The ROC curve seems satisfactory with maximum $TSS = 0.574$, and the RD is quite significant, showing no systematic over- or underprediction.

A purely monotonic decrease of TSS can be seen in the SVM, following an instantaneous peak. Some plateau in HSS is also noted, followed by a monotonic decrease. The ROC curve appears less satisfactory than in case of MLPs with maximum $TSS = 0.566$, and the RD shows some systematic underprediction for most of the forecast probability range.

For the RFs, we note a relatively similar behavior with MLPs, although with a slightly more pronounced HSS peak. The ROC curve seems better behaved than in the previous two methods with maximum $TSS = 0.615$, and the RD is arguably the best achieved together with the MLP RD.

4.5.2. Prediction of > C1 Flare Events Using Statistical Models

Figure 6 shows the forecast performances of the three tested statistical methods for > C1 class flare prediction.

For the LM, we note a decrease in the ACC of the method and some more or less similar behavior of HSS and TSS. The ROC curve seems satisfactory with maximum $TSS = 0.562$, while the RD appears to show a systematic overprediction below a forecast probability of 0.4 and a systematic underprediction above a forecast probability of 0.4 (excluding probabilities > 0.9).

A similar behavior with LM appears for the SSPs in the PR, while the ROC curve seems slightly better with maximum $TSS = 0.566$. The RD curve shows some systematic underprediction, although generally within the error bars.

Finally, for the LG, we note a similar behavior in the SSP as in the case of LM and PR, but arguably a better-behaved ROC curve with maximum $TSS = 0.567$. The RD seems to be the best behaved, compared to those of LM and PR.

4.5.3. Monte Carlo Simulation for > C1 Flares

In Table 3 we provide the average values of the skill scores ACC, TSS, and HSS for all prediction methods after the 200 replications of the Monte Carlo experiment regarding > C1 flares prediction. Table 3 shows that the maximum $HSS = 0.60$ is obtained with

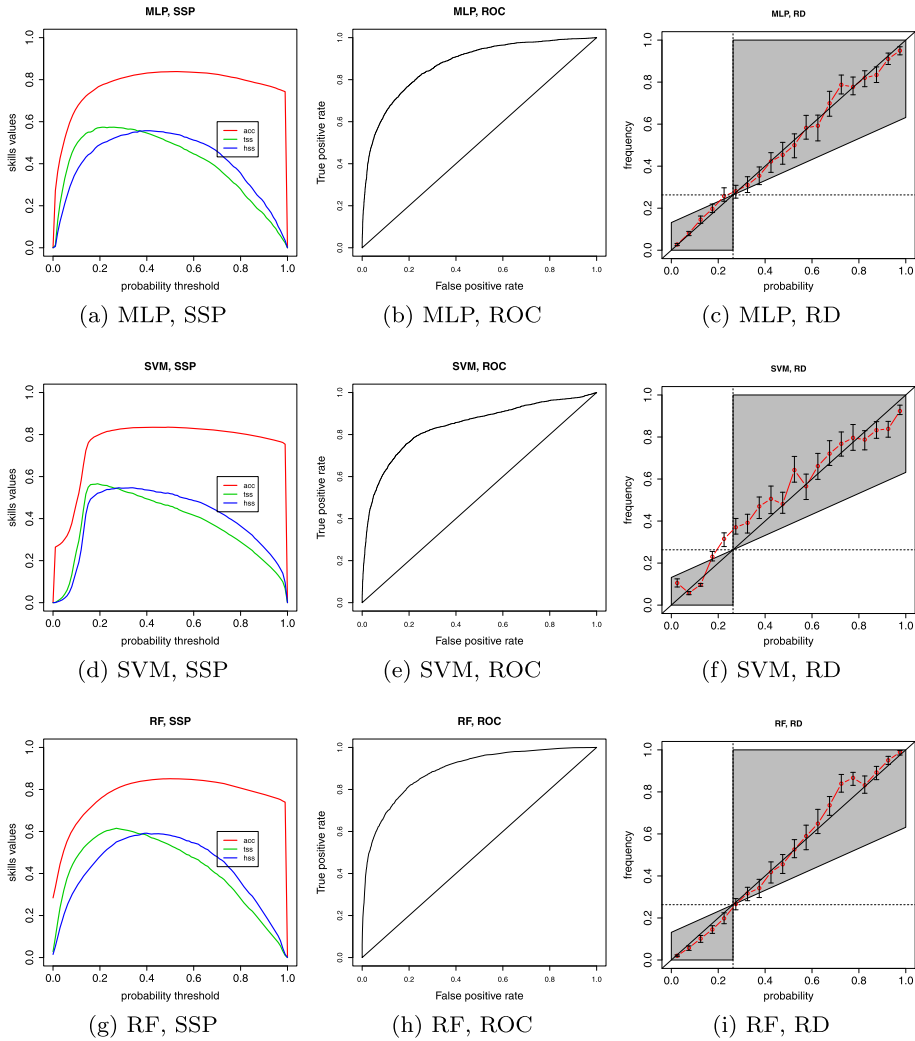


Figure 5 Same as Figure 3, but for >C1 flare prediction.

the RF method for a probability threshold of 40%. The corresponding skill score values are $ACC = 0.85 \pm 0.00$, $TSS = 0.59 \pm 0.01$, and $HSS = 0.60 \pm 0.01$. The second-best method in Table 3 for the same probability threshold is obtained with the LG method, with $ACC = 0.83 \pm 0.00$, $TSS = 0.54 \pm 0.01$, and $HSS = 0.56 \pm 0.01$. Considering again the probability threshold where the maximum TSS is observed, we obtain the optimal results for the RF method and threshold 30% with values $ACC = 0.82 \pm 0.00$, $TSS = 0.61 \pm 0.01$, and $HSS = 0.57 \pm 0.01$. The second-best method may be considered the MLP (or the LG in a tie) at a 30% threshold with $ACC = 0.81 \pm 0.00$, $TSS = 0.57 \pm 0.01$, and $HSS = 0.53 \pm 0.01$. For a range of probability thresholds (30%–40%), the method RF yields increasing values of HSS and decreasing values of TSS. As a result, it is again not clear which the best-fit value of the threshold probability is if we choose to simultaneously optimize both TSS and HSS. For example, an appealing RF forecasting model is with a threshold 35% and

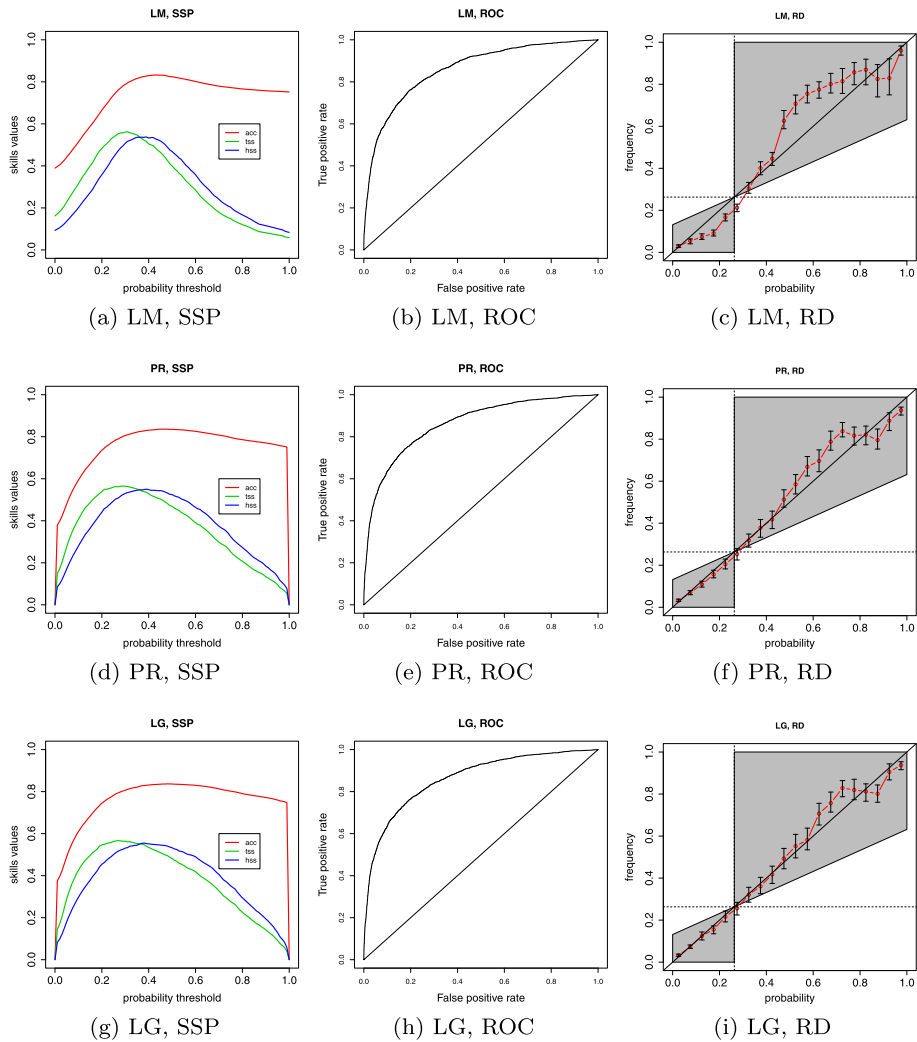


Figure 6 Same as Figure 4, but for > C1 flare prediction.

skill scores $ACC = 0.84 \pm 0.00$, $TSS = 0.60 \pm 0.01$ and $HSS = 0.59 \pm 0.01$ in Table 3. These results are generally above those reported for > C1 class flare predictability, namely $TSS \in [0.50, 0.55]$ and $HSS \in [0.40, 0.45]$ (Al-Ghraibah, Boucheron, and McAteer, 2015; Boucheron, Al-Ghraibah, and McAteer, 2015). In brief, our data samples, both training and testing, are comprehensive and generally unbiased.

4.6. Assessment of Prediction Methods and Predictor Strength

Following the presentation of results in Tables 2 and 3, we can see that both for > M1 and > C1 flare prediction, RF delivers the best skill score metrics for a wide range of probability thresholds. The second-best method is MLP, together with LG. In this setting, we performed some additional evaluation that confirms these results.

Table 3 Monte Carlo scenario 2, based on 200 SHARP datasets, on > C1 GOES flare prediction. Numbers in boldface correspond to the most significant results of a given method (MLP: multi-layer perceptron; LM: linear regression; PR: probit regression; LG: logit regression; RF: random forest; SVM: support vector machine).

Par	%	MLP			LM			PR			LG			RF			SVM		
		ACC	TSS	HSS	ACC	TSS	HSS	ACC	TSS	HSS	ACC	TSS	HSS	ACC	TSS	HSS	ACC	TSS	HSS
val ₀	0.00	0.00	0.00	0.00	0.39	0.16	0.09	0.00	0.00	0.00	0.00	0.00	0.00	0.28	0.03	0.01	0.00	0.00	0.00
val ₅	0.05	0.52	0.33	0.21	0.44	0.23	0.14	0.47	0.27	0.17	0.48	0.28	0.17	0.52	0.34	0.21	0.29	0.03	0.02
val ₁₀	0.10	0.66	0.49	0.35	0.51	0.32	0.20	0.58	0.40	0.27	0.60	0.42	0.29	0.64	0.48	0.34	0.44	0.22	0.13
val ₁₅	0.15	0.72	0.55	0.43	0.58	0.40	0.27	0.66	0.49	0.36	0.68	0.50	0.38	0.71	0.55	0.42	0.75	0.54	0.45
val ₂₀	0.20	0.76	0.57	0.48	0.66	0.48	0.35	0.73	0.55	0.44	0.74	0.55	0.45	0.76	0.59	0.48	0.80	0.57	0.53
val ₂₅	0.25	0.79	0.57	0.51	0.74	0.55	0.45	0.78	0.56	0.49	0.78	0.57	0.50	0.79	0.61	0.53	0.82	0.56	0.55
val ₃₀	0.30	0.81	0.57	0.53	0.79	0.57	0.51	0.80	0.57	0.53	0.81	0.57	0.53	0.82	0.61	0.57	0.83	0.54	0.55
val ₃₅	0.35	0.82	0.56	0.55	0.82	0.55	0.54	0.82	0.56	0.55	0.82	0.56	0.55	0.84	0.60	0.59	0.84	0.52	0.55
val ₄₀	0.40	0.83	0.55	0.55	0.83	0.52	0.54	0.83	0.53	0.55	0.83	0.54	0.56	0.85	0.59	0.60	0.84	0.50	0.54
val ₄₅	0.45	0.84	0.53	0.55	0.83	0.47	0.52	0.84	0.50	0.55	0.84	0.51	0.55	0.85	0.56	0.59	0.84	0.48	0.53
val ₅₀	0.50	0.84	0.50	0.55	0.83	0.40	0.47	0.84	0.47	0.53	0.84	0.48	0.53	0.85	0.54	0.59	0.84	0.46	0.52
val ₅₅	0.55	0.84	0.48	0.53	0.81	0.34	0.41	0.83	0.43	0.50	0.84	0.45	0.52	0.85	0.51	0.57	0.83	0.44	0.51
val ₆₀	0.60	0.84	0.45	0.51	0.80	0.28	0.35	0.82	0.38	0.46	0.83	0.41	0.48	0.85	0.48	0.55	0.83	0.42	0.49
val ₆₅	0.65	0.83	0.41	0.49	0.79	0.22	0.29	0.82	0.34	0.42	0.82	0.37	0.44	0.84	0.44	0.52	0.83	0.38	0.46
val ₇₀	0.70	0.82	0.37	0.45	0.78	0.18	0.24	0.81	0.29	0.37	0.81	0.32	0.40	0.83	0.40	0.48	0.82	0.35	0.43
val ₇₅	0.75	0.82	0.33	0.41	0.77	0.15	0.20	0.80	0.25	0.32	0.80	0.27	0.35	0.82	0.34	0.43	0.81	0.32	0.39
val ₈₀	0.80	0.81	0.28	0.36	0.77	0.12	0.17	0.79	0.20	0.27	0.79	0.22	0.29	0.81	0.28	0.36	0.81	0.28	0.36
val ₈₅	0.85	0.79	0.22	0.29	0.76	0.10	0.14	0.78	0.16	0.22	0.78	0.18	0.24	0.79	0.21	0.29	0.80	0.24	0.31
val ₉₀	0.90	0.78	0.16	0.21	0.76	0.08	0.11	0.77	0.13	0.18	0.77	0.14	0.19	0.78	0.15	0.21	0.79	0.19	0.25
val ₉₅	0.95	0.75	0.08	0.11	0.76	0.07	0.10	0.76	0.09	0.13	0.76	0.09	0.13	0.76	0.08	0.12	0.77	0.14	0.19
val ₁₀₀	1.00	0.00	0.00	0.00	0.75	0.06	0.08	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00

We present analytical results in Appendix A for the predictor strength. It seems that $\log R$ and WL_{SG} rank in first place for both $> C1$ and $> M1$ flare prediction, closely followed by the Ising energy and the TLMPIIL.

In order to investigate the robustness of our results, we present additional results in Appendix C, where we make predictions once a day (at 00:00 UT). The mean evolution (over 200 Monte Carlo iterations) of ACC, TSS, and HSS with respect to the probability threshold is presented. Likewise, the BS, AUC, and BSS are presented. The main finding is that issuing forecasts once a day keeps similar average skill scores as are achieved for issuing forecasts eight times a day, but the associated uncertainties (*e.g.* standard deviations) are higher in the case of daily predictions.

A final word for the comparison of ML algorithms *versus* conventional statistics models for this specific dataset and positive/negative class definitions is provided in Appendix D. There, we have included an auxiliary meta-analysis of the results in Tables 2 and 3 in order to clearly show whether the ML category of prediction algorithms performs better than the conventional statistics models in the $> M1$ and $> C1$ flare prediction cases. A multicriteria analysis using the weighted-sum (WS) method (Greco, Figueira, and Ehrgott, 2016) seems appropriate in order to aggregate the performance metrics ACC, TSS, and HSS of all classifiers as a function of the probability threshold (*e.g.* using equal weights for the aggregation). In this way, a composite index (CI) as a measure of overall utility is computed for every algorithm and probability threshold combination. There exist $21 \times 6 = 126$ such alternatives when we use a 5% probability threshold grid, such as the grid in Tables 2 and 3. The ranking, in non-increasing order, of the CI reveals the overall merit of every probabilistic classifier and also allows us to draw conclusions for groups of classifiers, such as the group of ML methods (comprising RF, SVM, and MLP) and the group of conventional statistics methods (comprising LM, PR, and LG). Appendix D presents this multicriteria WS analysis, revealing that overall, in $> C1$ flare prediction ML outperforms conventional statistics methods by 71% *versus* 29% in the synthesis of the top $100(1/6) = 16.6\%$ performing methods (top 21 methods out of a total of 126). Likewise, in the $> M1$ flare prediction case, ML outperforms conventional statistics methods by 62% *versus* 38% in the synthesis of the top $100(1/6) = 16.6\%$ performing methods. This shows that $> C1$ flare prediction is more advantageous for ML *versus* statistical methods, in comparison to the $> M1$ flare case. This is due to the low performance of the SVM in $> M1$ flare prediction, which in turn is due to the way in which we have implemented, for simplicity, the SVM for a highly unbalanced sample in $> M1$ flare prediction,³ using a single C constant and not two different C_1, C_2 constants during the SVM training with Equation 7.

In auxiliary runs (available upon request), we also noted that when the sample size was very low, using ML algorithms posed no advantage over conventional statistics models. In order to have proper training, the ML algorithms need $N > 2,000$ for $K = 13$, especially for the $> M1$ flare prediction.

4.7. Statistical Tests for Random Forest *Versus* MLP and Calculation of AUC and Brier Skill Scores

In Section 4.7.1 we present results of a t -test between the two best-performing methods according to maximizer thresholds for either TSS or HSS for $> M1$ class and $> C1$ class flare cases. Section 4.7.2 presents additional calculations reporting on BS, BSS, and AUC, which we used to assess the classification in the prediction.

³Even by using the SVM weighted variant and recomputing the WS ranking using this variant (*e.g.* see Figure 3 and Table 5), the qualitative results of the ranking still hold.

Table 4 Unpaired t -tests to compare the means of TSS and HSS metrics (out-of-sample) for the best and the second-best methods in $> M1$ and $> C1$ flare forecasting.

No.	Metric	Threshold (%)	Best	Second-best	p -value
$> M1$ class flare prediction					
1	TSS	10	RF	LM	$< 10^{-4}$
2	HSS	10	RF	LM	$< 10^{-4}$
3	TSS	25	RF	MLP	$< 10^{-4}$
4	HSS	25	RF	MLP	$< 10^{-4}$
$> C1$ class flare prediction					
5	TSS	30	RF	MLP	$< 10^{-4}$
6	HSS	30	RF	MLP	$< 10^{-4}$
7	TSS	40	RF	LG	$< 10^{-4}$
8	HSS	40	RF	LG	$< 10^{-4}$

4.7.1. Unpaired t -Tests to Compare Two Means for TSS and HSS of Random Forest versus MLP

A t -test compares the means of two groups. Here, we used the t -test to compare the mean TSS (HSS, respectively) of the RF method *versus* those of the MLP method (or in general the second-best performing method). Means were considered with respect to the Monte Carlo simulations performed on the 200 replications of the previous section. The TSS (HSS, respectively) values considered are those for specific probability thresholds maximizing either TSS or HSS. Table 4 presents the t -test results regarding the best and the second-best methods with respect to either TSS or HSS for these specific probability thresholds.

We find that RF is always (*i.e.* 8/8 of times) statistically better than the second-best method (which is the MLP 4/8 of times), with respect to both TSS and HSS.

4.7.2. Calculation of AUC and Brier Skill Scores

Tables 5 and 6 present the calculated mean values of BS, AUC, and BSS for the $> M1$ and $> C1$ flare prediction cases, respectively.

For the $> M1$ flare case (Table 5), results show that on average, the best BS and BSS results are achieved with the RF method (BS = 0.0266; BSS = 0.4163). The best AUC results are achieved with the RF method (AUC = 0.9556), but also for the PR (AUC = 0.9392) and LG (AUC = 0.9391) methods.

For the $> C1$ flare case (Table 6), results show that on average, the best BS and BSS results are achieved with the RF method (BS = 0.1074; BSS = 0.4426). The best AUC results are also achieved with the RF method (AUC = 0.8927), with other methods (except SVM) following closely. The SVM probably needs better fine-tuning, given its sensitivity on γ and cost (see Section 4.2).

4.8. Related Published Work and Comparison to Our Results

Ahmed *et al.* (2013) presented prediction results for $> C1$ class flares using cross-validation with 60% training and 40% testing subsets, with ten iterations in operational and segmented mode. Since our analysis focuses in operational mode, the gold standard for near-real-time operational systems such as FLARECAST, we present here their results on the operational mode for the period April 1996 – December 2010: POD = 0.455 & POFD = 0.010 thus

Table 5 Mean values for BS, BSS, and AUC for all tested models on the prediction of > M1 flares. Means are obtained after 200 Monte Carlo replications. Parentheses below the values denote standard deviations. Lower values indicate a better performance for BS, and higher values indicate a better performance for AUC and BSS.

MLP	LM	PR	LG	RF	SVM	SVM _{weighted}
BS						
0.0324 (0.0013)	0.0331 (0.0008)	0.0305 (0.0008)	0.0302 (0.0008)	0.0266 (0.0008)	0.0327 (0.0012)	0.0357 (0.0011)
AUC						
0.9301 (0.0067)	0.9278 (0.0043)	0.9392 (0.0033)	0.9391 (0.0033)	0.9556 (0.0035)	0.8320 (0.0168)	0.9175 (0.0059)
BSS						
0.2903 (0.0267)	0.2745 (0.0117)	0.3323 (0.0128)	0.3375 (0.0134)	0.4163 (0.0126)	0.2829 (0.0159)	0.2181 (0.0154)

Table 6 Same as Table 5, but for the prediction of > C1 flares.

MLP	LM	PR	LG	RF	SVM
BS					
0.1167 (0.0014)	0.1292 (0.0012)	0.1201 (0.0012)	0.1191 (0.0012)	0.1074 (0.0012)	0.1226 (0.0015)
AUC					
0.8731 (0.0029)	0.8638 (0.0029)	0.8665 (0.0027)	0.8669 (0.0027)	0.8927 (0.0026)	0.8466 (0.0033)
BSS					
0.3940 (0.0069)	0.3293 (0.0052)	0.3767 (0.0055)	0.3818 (0.0058)	0.4426 (0.0056)	0.3636 (0.0063)

TSS = POD – POFD = 0.445 and HSS = 0.539. Hence, Ahmed *et al.* reported (using a variant of a neural network, and threshold 50%) results for flares > C1: TSS = 0.445 and HSS = 0.539.

Li *et al.* (2008) presented results using an SVM coupled with k-nearest neighbor (KNN) for flare prediction > M1 in a way that, unfortunately, cannot be used to recover TSS and HSS values. Instead, they reported Equal = TN + TP, High = FP, Low = FN. The accuracy achieved is only ACC = 57.02% for SVM and ACC = 63.91% for SVM-KNN for the testing year 2002.

Song *et al.* (2009) presented results using an ordinal logistic regression model classifying the C-, M-, and X-class flares with response values 1, 2, and 3, respectively. The B-class flare (or no flares) category received class 0 (baseline). Their sample contained 34 X-class flares, 68 M-class flares, 65 C-class flares, and 63 B-class or no-flare cases. A clear drawback of this sample is that it was not taken using a random number generator but seems to be hand-picked aiming at studying the considered 230 events during the period 1998–2005. As a result, the sample is biased in that the occurrence rates of the various flare classes are not representative of an actual solar cycle. Perhaps not surprisingly, these authors presented high TSS and HSS values that, given the sample, might be taken with a conservative outlook.

From the results of Model 4 in that study (*i.e.* Table 8 of Song *et al.*, 2009), we are able to infer that for C-class flares, Song *et al.* computed values $TSS = 0.65$ and $HSS = 0.623$ (C1 – C9 flares). Moreover, we maintain an impression that these numbers are obtained in-sample for the dataset with 230 events in Song *et al.* (2009).

Yu *et al.* (2009) used a sliding window approach to account for the evolution of three magnetic flare predictors with importance index above 10 (for the definition of the flare importance index, see Yu *et al.*, 2009). The time period was 1996 to 2004, with a cadence of 96 minutes. The authors used the C4.5 decision tree algorithm and the learning vector quantization (LVQ) neural network, both implemented in WEKA (Witten *et al.*, 2016; Hall *et al.*, 2009). The authors used a 10-fold cross-validation approach with 90% training and 10% testing sets from the original sample. The sliding window size was 45 observations. Their results showed that the sliding window versions of the C4.5 and LVQ neural network algorithms improved the results obtained with the same algorithms for a sliding window size equal to 0. Since the authors presented only the TP rate and the TN rate results, we are not able to recover their HSS value. Their recovered TSS is $TSS = 0.651$ for the C4.5 algorithm with a sliding window of 45 observations and $TSS = 0.667$ for the LVQ, also with a sliding window of 45 observations.

Yuan *et al.* (2010) used the same dataset as in Song *et al.* (2009) and proposed a cascading approach, using first an ordinal logistic regression model to produce probabilities for GOES flare classes B, C, M, and X (associated with response levels 0, 1, 2, and 3, respectively), and second, feeding the probability values to an SVM in order to obtain the final class membership. Their results, according to Yuan *et al.* (2010), improved the prediction especially for X-class flares (response level = 3 in the ordinal logistic regression), but were still not exceptionally high. For example, for level = 1, therefore for C-class flares, we were able to recover the following TSS values for the used methods: logistic regression: $TSS = 0.22$, SVM: $TSS = 0.08$, logistic regression + SVM: $TSS = 0.09$. These rather fair results, as can be seen from the contingency tables presented in Yuan *et al.* (2010), may be due to the selection of a probability threshold value at 50% for levels 0, 1, and 3 in the ordinal logistic regression model and at 25% for the level 3 (X-class) flares in the same model. Choosing a threshold equal to 50% maximizes ACC but not TSS/HSS, as can be seen both here and in Bloomfield *et al.* (2012).

Colak and Qahwaji (2009) developed an online solar flare forecasting system called ASAP. Their prediction algorithm is a combination of two neural networks with the sum-of-squared error (SSE) objective function, where the first neural network predicts whether a flare of all types (C, M, or X) will occur, and if the prediction is yes, the second neural network predicts whether a C-, M-, or X-class flare will occur. The ASAP system was developed in C++ and has been validated with data from 1999 to 2002 (around the peak of Solar Cycle 23). The predictors were the sunspot area and characteristics from the McIntosh classification of sunspots (Zpc scheme). They obtained $HSS = 49.3\%$ (C-class flares) and $HSS = 47\%$ (M-class flares) for a forecast window of 24 h.

Wang *et al.* (2008) developed an MLP neural network using three input variables for the prediction of solar flares of class $> M1$. The predictors were the maximum horizontal gradient $|\text{grad}_h(B_z)|$, the length L of the neutral line, and the number of singular points η . A limitation of the study is that only flaring active regions (at GOES C1 and above) were sampled and considered. The forecast window was 48 h. The authors presented prediction results for the period 1996–2002 (training set: April 1996 to December 2001, testing set: January 2002 to December 2002). The results were presented as plots of the X-ray flux associated with the predicted/observed flares for the test year 2002, therefore a comparison with the authors' skill scores is not possible. This work reported $ACC = 69\%$ for the test year.

Bobra and Couvidat (2015) applied an SVM to a sample of 5,000 non-flaring and 303 flaring (at the GOES $> M1$ level) AR. Those $N = 5,303$ AR with $N = 5,000$ negative examples and $P = 303$ positive examples (ratio $N/P = 16.5$), were sampled from the ≈ 1.5 million patches of the SHARP product (Bobra *et al.*, 2014) between 2010 and 2014. The authors selected 285 M-class flares and 18 X-class flares observed between 2010 May and 2014 May. By comparison, our study relies on a representative sample of flaring/non-flaring AR in the period 2012–2016 and for flares $> M1$, with a ratio $N/P = 19.9$ ($P = 1,108$ and $N = 22,026$). By inspecting Table 3 of Bobra and Couvidat (2015), we see that the authors report the results as $ACC = 0.924 \pm 0.007$, $TSS = 0.761 \pm 0.039$, and $HSS_2 = 0.517 \pm 0.035$, while our results are $ACC = 0.93 \pm 0.00$, $TSS = 0.74 \pm 0.02$, and $HSS = 0.49 \pm 0.01$ (their definition of HSS_2 is the same as the HSS definition in Section 4.3). Thus, our results with random forests are competitive with those of Bobra and Couvidat (2015). We note that we used a 50/50 rule for splitting training/testing sets, while Bobra and Couvidat (2015) use a 70/30 rule. Moreover, N/P in Bobra and Couvidat (2015) is 16.5, while in our case, N/P is 19.9. Finally, we used no fine-tuning in the parameters of the random forest, while Bobra and Couvidat (2015) carefully tuned the C , γ and C_1/C_2 of their Equations 2, 5, and 6, respectively. Regardless, Bobra and Couvidat (2015) still represent the state-of-the art in solar flare forecasting so far.

Boucheron, Al-Ghraibah, and McAteer (2015) applied support vector regression (SVR) to 38 predictors characterizing the magnetic field of solar AR in order to predict i) the flare size and ii) the time-to-flare using SVR modeling. The forecast window they used varied between 2 and 24 hours with a step of 2 hours (12 cases of forecast windows). By using the size regression with appropriate thresholds (different to the usual probability thresholds, for example, in Bloomfield *et al.*, 2012), the authors achieved prediction results for $> C1$ flares with $TSS = 0.55$ and $HSS = 0.46$, while reporting that using the same data, Al-Ghraibah, Boucheron, and McAteer (2015) achieved $TSS \approx 0.50$ and $HSS \approx 0.40$, respectively, for the prediction of $> C1$ class flares.

Al-Ghraibah, Boucheron, and McAteer (2015) applied relevance vector machines (RVM), a technique that is a generalization of SVM, to a set of 38 magnetic properties characterizing 2124 AR in a total of 122,060 images across different time points for all AR. They predicted $> C1$ flares using either the full set of properties or suitable subsets thereof. The magnetic properties are of three types: i) snapshots in space and time, ii) evolution in time, and iii) structures of multiple size scales. Al-Ghraibah, Boucheron, and McAteer (2015) reported results (*e.g.* see their Table 5 and Figure 6) in the range $TSS \approx 0.51$ and $HSS \approx 0.39$, which is a baseline result for the literature when no temporal information is included in the predictor set (*i.e.* static images are used).

5. Conclusions

We presented a new approach for the efficient prediction of $> M1$ and $> C1$ solar flares: classic and modern machine-learning (ML) methods, such as multi-layer perceptrons (MLP), support vector machines (SVM), and random forests (RF) were used in order to build the prediction models. The predictor variables were based on the SDO/HMI SHARP data product, available since 2012.

The sample was representative of the solar activity during a five-year period of Solar Cycle 24 (2012–2016), with all calendar days within this period included in the sample. The cadence of properties, or predictors, within the chosen days was 3 hours.

We showed that the RF method could be our prediction method of choice, both for the prediction of $> M1$ flares (with a relative frequency of 4.8%, or 1,108 events) and for the prediction of $> C1$ flares (with a relative frequency of 26.1%, or 6,029 events). In terms of categorical skill scores, a probability threshold of 15% for $> M1$ flares gives rise to mean (after 200 replications) RF skill scores on the order $TSS = 0.74 \pm 0.02$ and $HSS = 0.49 \pm 0.01$, while a probability threshold of 35% for $> C1$ flares gives rise to mean $TSS = 0.60 \pm 0.01$ and $HSS = 0.59 \pm 0.01$. The respective accuracy values are $ACC = 0.93$ and $ACC = 0.84$. In terms of probabilistic skill scores, the ranking of the ML techniques with respect to their BSS against climatology is RF (0.42), MLP (0.29), and SVM (0.28) for $> M1$ flares and RF (0.44), MLP (0.39), and SVM (0.36) for $> C1$ flares.

We further indicate that for $> M1$ flare prediction, SVM and MLP need additional tuning of their hyperparameters (Section 4.2) in order to produce comparable results with RF. Moreover, several statistical methods (linear regression, probit, and logit) produced acceptable forecast results when compared with the ML methods. By increasing the number of hidden nodes, the MLP networks provide flatter skill score profiles (*i.e.* ACC, TSS, and HSS as a function of the threshold probability), but the peak values of the corresponding curves are lower than those achieved by MLP networks with fewer hidden nodes. Regarding the $> C1$ flares, all forecast methods work acceptably, although the best method is, again, RF. A Monte Carlo experiment showed that results are robust with respect to different realizations of the training/testing pair, with different random seeds. Monte Carlo modeling also decreases the amplitudes of the applicable standard deviations of the skill scores. Typically standard deviations are larger for the $> M1$ flare case than for that of $> C1$ flares. This is to be attributed to the different occurrence frequency of flares in the two cases.

The RF is a relatively new approach to solar flare prediction. Nonetheless, it may be preferable over other widely used ML algorithms, at least for the datasets exploited so far, giving competitive results without much tuning of the RF hyperparameters. This generates hope for future meaningful developments in the formidable solar flare prediction problem, at the same time aligning with excellent performance for RF reported in several classification benchmarks (Fernández-Delgado *et al.*, 2014). This important statement made, it appears that even with the application of RF, solar flare prediction in the foreseeable future will likely continue to be probabilistic (*i.e.* 0.0–1.0, continuous), rather than binary (*i.e.* 0 or 1).

In terms of the predictors importance, Schrijver's R is found to be among the most statistically significant predictors, together with WL_{SG} . The Ising energy and the TLMPIIL are also considered as important, ranking slightly below the previous two predictors. This stems from the importance calculations according to the Fisher score and RF importance for the $> C1$ and $> M1$ flare cases in Appendix A. This result is also in line with the common knowledge that flares occur mostly when strong and highly sheared MPILs are formed. Other MPIL-highlighting predictors, such as the effective connected magnetic field strength, B_{eff} (Georgoulis and Rust, 2007), remain to be tested, in conjunction with R and WL_{SG} , as their cadence was lower than 3 h at the time this study was performed.

An interesting finding for the RF technique (Appendix B) is obtained by the predictors' ranking information according to their importance, as measured by the Fisher score. Namely, when we create prediction models with a varying number of the most important predictors included, the RF prediction performance (in terms of TSS and HSS) continues to improve monotonically with the number of included parameters. Conversely, the MLP and SVM algorithms achieve only slight improvements in prediction results (again in terms of TSS and HSS) by adding more than, say, the six most important predictors. This interesting finding may further improve forecasting when more viable predictors become available.

For future FLARECAST-supported research, we plan to enlarge our analysis sample by reducing the property cadence from 3 h to 1 h or even less (the limit is the inherent cadence

of SDO/HMI SHARP data, namely 12 min). Another direction of future research is to investigate the robustness of our results for samples created with a higher cadence of 12 h (24 h) coupled with a forecast window of 12 h (24 h). Furthermore, we plan to exploit the substantial time-series aspect of our data using recurrent neural networks, possibly trained with evolutionary algorithms. The present work, along with a series of similar concluded or still ongoing studies, is considered for possible integration in the final FLARECAST online system and forecasting tool, to be deployed by early 2018.

Acknowledgements We would like to thank the anonymous referee for very helpful comments that greatly improved the initial manuscript. This research has been supported by the EU Horizon 2020 Research and Innovation Action under grant agreement No.640216 for the “Flare Likelihood And Region Eruption foreCASTing” (FLARECAST) project. Data were provided by the MEDOC data and operations centre (CNES/CNRS/Univ. Paris-Sud), <http://medoc.ias.u-psud.fr/> and the GOES team.

Disclosure of Potential Conflicts of Interest The authors declare that they have no conflicts of interest.

Appendix A: Importance of Predictors for Flare Prediction

We computed the Fisher score (Bobra and Couvidat, 2015; Chang and Lin, 2008; Chen and Lin, 2006) and the Gini importance (Breiman, 2001) for every predictor in the case of > M1 and > C1 flares. The obtained values for the importance of several predictors are presented in Figures 7 and 8 for > C1 and > M1 flare prediction, respectively. The Fisher score, F , is defined for the j th predictor as

$$F(j) = \frac{(\bar{x}_j^{(+)} - \bar{x}_j)^2 + (\bar{x}_j^{(-)} - \bar{x}_j)^2}{\frac{1}{n^+ - 1} \sum_{k=1}^{n^+} (x_{k,j}^{(+)} - \bar{x}_j^{(+)})^2 + \frac{1}{n^- - 1} \sum_{k=1}^{n^-} (x_{k,j}^{(-)} - \bar{x}_j^{(-)})^2}. \quad (\text{A.1})$$

In Equation A.1, \bar{x}_j , $\bar{x}_j^{(+)}$, and $\bar{x}_j^{(-)}$ are the mean values for the j th predictor over the entire sample, the positive class, and the negative class, respectively. Furthermore, n^+ (n^-) are the number of positive (negative) class observations. In addition, $x_{k,j}^{(+)}$ ($x_{k,j}^{(-)}$) are the values for the k th observation of the j th predictor belonging in the positive (negative) class. The higher the value of $F(j)$, the more important the j th predictor.

The Gini importance is returned with the randomForest function of the randomForest package in R. The higher the Gini importance of the j th predictor, the more important this predictor.

Table 7 Abbreviations for predictors used in the main text (Symbol1) and in Figures 7 and 8 (Symbol2).

Abbreviations for predictors		
Symbol1	Symbol2	Description
logR	r_value_logr	Schrijver’s R value
FSPI	alpha_exp_fft	Fourier spectral power index
TLMPIL	mpil	Magnetic polarity inversion line
DI	decay_index	Decay index
WLSG	wlsg	Gradient-weighted integral length of the neutral line
IsinEn1	ising_energy	Ising energy original
IsinEn2	ising_energy_part	Ising energy partitioned

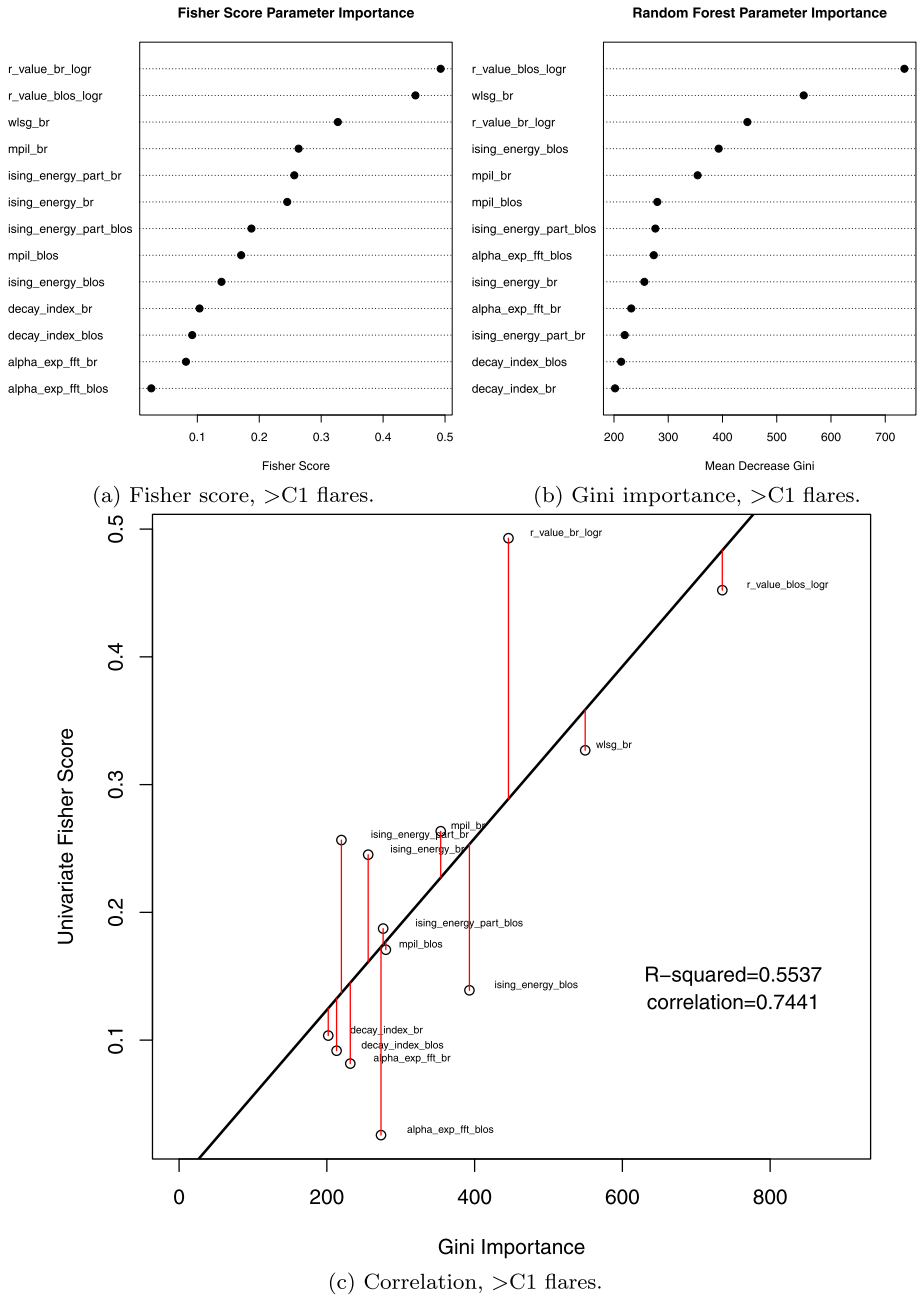
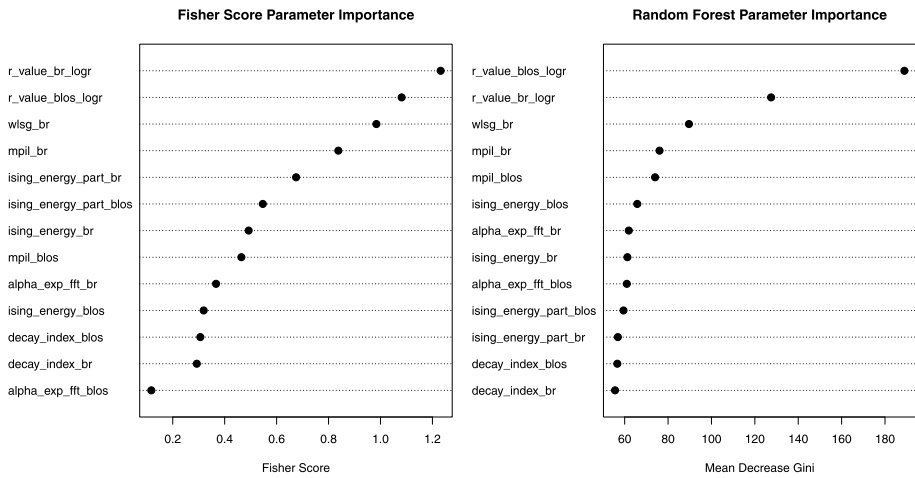


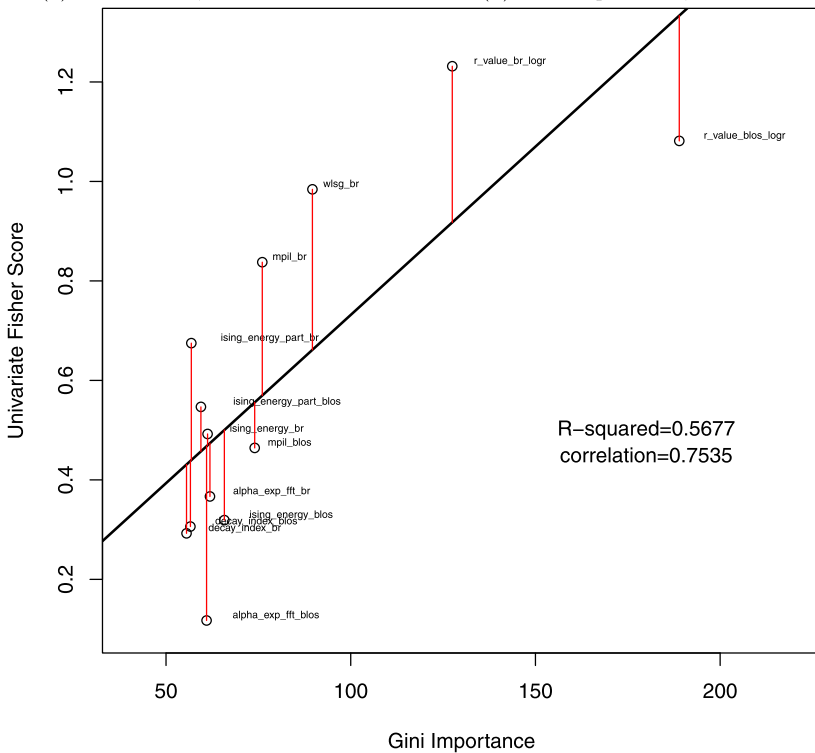
Figure 7 Importance of several predictors while predicting >C1 flares.

We note that the correlation between the two quantities (Fisher score and Gini importance) is $r = 0.7441$ for >C1 flares and $r = 0.7535$ for >M1 flares, respectively. This shows that the two methods qualitatively agree on describing which predictors are the most



(a) Fisher score, >M1 flares.

(b) Gini importance, >M1 flares.



(c) Correlation, >M1 flares.

Figure 8 Importance of several predictors while predicting > M1 flares.

important regarding flare prediction, in both classes of flare prediction. Figure 7 also shows that for > C1 flares, the top three ranked predictors for both Fisher score and Gini importance are the two versions of Schrijver’s R and WL_{SG} . For the > M1 flares, from Figure 8

the top four ranked predictors for either Fisher score or Gini importance are the two versions of Schrijver's R , WL_{SG} , and $TLMPIL_{Br}$. In Appendix A the terminology for every predictor is explained in Table 7.

Appendix B: Prediction Models Resulting from Ranking the Predictors

We employed a backward elimination procedure, eliminating gradually predictors according to their Fisher score rank, starting from the model with all $K = 13$ predictors included. In every step, we eliminated the least important predictor from the set of currently included predictors. In this way, we obtained prediction results for models with 2, 3, ..., 11, 12 predictors included for the ML methods, RF, SVM, and MLP and the conventional statistics methods LM, PR, and LG. The results of this iterative procedure for flares $> C1$ and $> M1$ are presented in Figures 9 and 10.

Figure 9 shows that there is a cut-off for the number of parameters included in the RF equal to the 6 most important ones (according to the Fisher score in Equation A.1) above which the RF is advantageous over the other two ML algorithms. For low-dimensional prediction models (e.g. fewer than 6 included parameters) there is no special advantage in using RF, and MLP or SVM seem a better choice then. This finding shows that of the highly correlated set of predictors, the MLP and SVM perform well using only a handful of them (fewer than 6), but the RF continues to improve its performance in higher-dimension settings, when the prediction model includes all 12 most important predictors. There is interest in investigating the performance of RF when the number of (correlated) predictors would be twice or three times that of the present study (24–36 predictors). Would the upward trend in Figure 9a continue to hold when the number of included parameters increased to 24 or 36? We note that RF is the only ML algorithm in the present study that belongs in the category of “ensemble” methods. Moreover, in Figure 9, the performance of the three conventional statistics methods LM, PR, and LG is presented. Clearly, the LM presents the worst forecasting ability, and we also note that the other two methods, PR and LG, score similar values for the TSS and HSS in general. It is also noteworthy that the profiles of PR and LG are very flat as a function of the number of included predictors, even flatter than the profiles from SVM and MLP.

Likewise, Figure 10 shows that for low-dimensional settings, RF is worse than MLP. The cut-off seems again to be six included parameters. Above this value, the RF provides better out-of-sample TSS and HSS than MLP. There seems to be a problematic region between three and six parameters included for the SVM, where adding more parameters to the SVM degrades its performance. With more than six parameters, the SVM performance improves again. Similarly to the $> C1$ -class flares case, we again note in Figure 10 rather flat profiles for the TSS and HSS for the conventional statistics methods, with PR and LG showing better behavior than LM.

One general conclusion is that for very few predictors $K < 6$, all methods work the same, so for parsimony, the conventional statistics methods could be preferred. This is also true for very small samples $N < 2,000$ (results available upon request). Conversely, when $K > 6$ and $N \geq 10,000$, the ML methods and especially the RF are better.

We note that in Appendix B, the MLP always has four hidden nodes and the SVM has γ and cost parameters analogously to the full $K = 13$ SVM model for $> C1$ and $> M1$ flares cases.

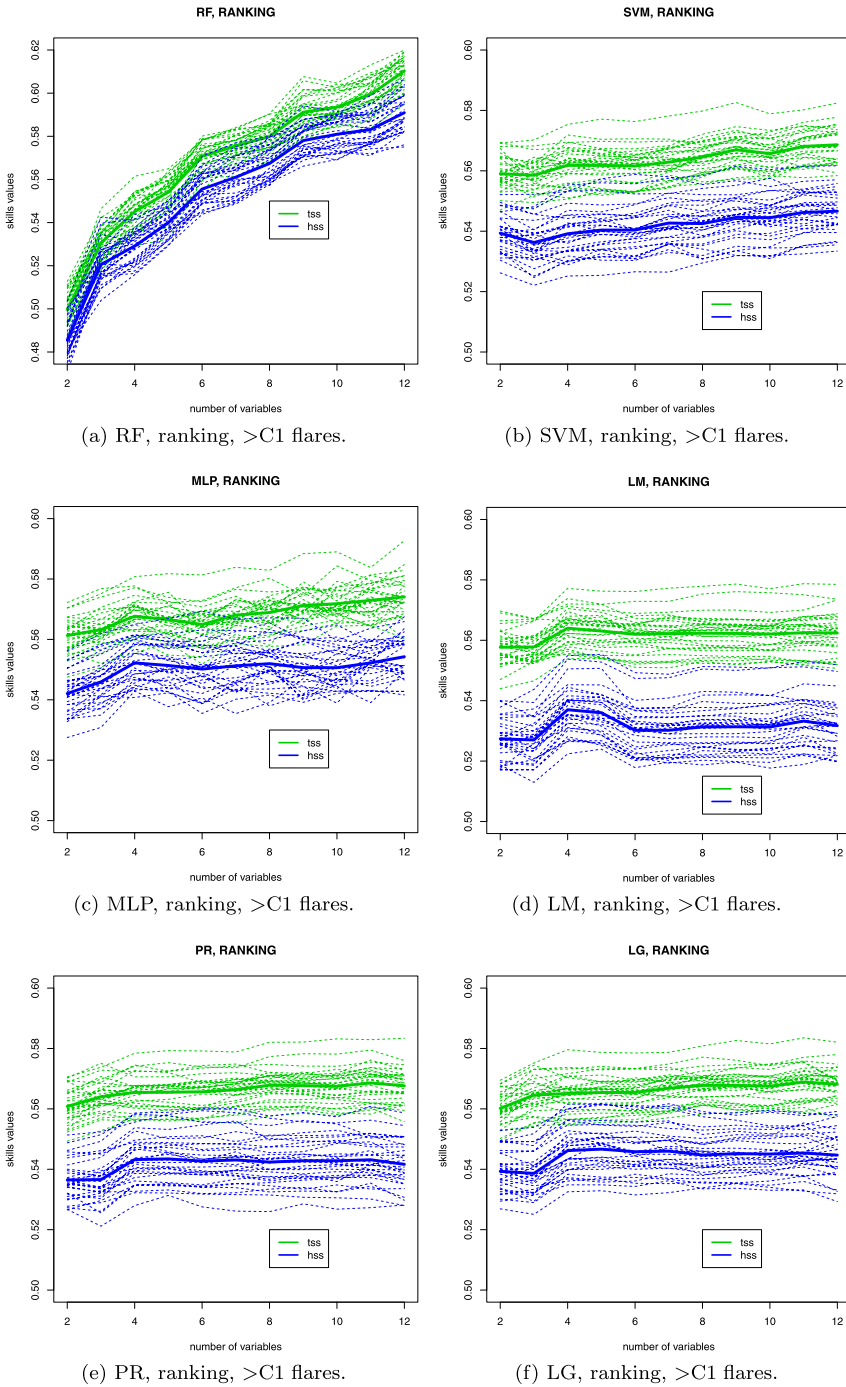


Figure 9 Out-of-sample skill scores (TSS and HSS) for the three ML prediction methods and the three statistics methods during the ranking procedure for > C1 flares. The *thick continuous lines* depict the averages of the skill scores over 30 randomized runs.

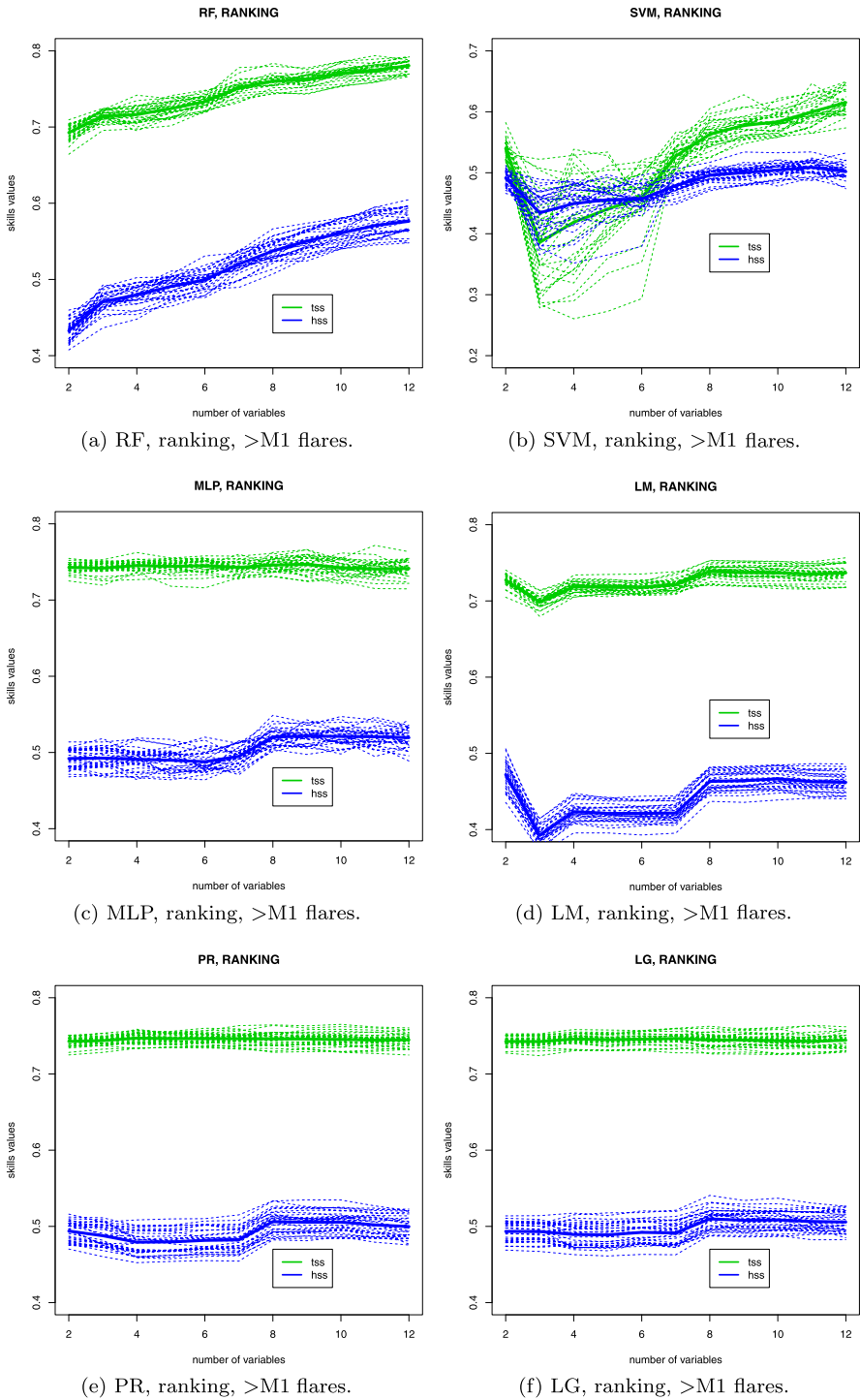


Figure 10 Same as in Figure 9, but for > M1 flares.

Appendix C: Validation Results when Predictions Are Issued Only Once a Day (at Midnight)

We present here forecasting results in the following scenario:

- i) The training is performed as in the main scenario.
- ii) The testing is performed only for the observations in the testing set of the main scenario, which correspond to a time of 00:00 UT. To achieve this, we filter for the observations in the previous testing set with `midnightStatus = TRUE`.

This method of training-testing is called the “hybrid method”, where training is done with a cadence of 3 h and a forecast window of 24 h, and testing is done with a cadence of 24 h and a forecast window of 24 h. The hybrid method is preferable over using a training phase with a cadence of 24 h, which would result in undertrained models because of the limited sample size during training.

Tables 8 and 9 are analogous to Tables 5 and 6 of the main scenario, but for midnight-only (so once a day) predictions. For completeness, we recall that Table 8 is for > M1 flare prediction and Table 9 is for > C1 flare prediction.

By comparing Table 8 to Table 5, we see that BS and AUC do not change much on average when we move from the baseline scenario to the midnight prediction scenario. Nevertheless, the associated uncertainty increases in the case of midnight-only predictions, since the size of the testing set is smaller (only one, rather than eight, predictions *per day*). More

Table 8 Same as Table 5, but for predictions issued only at midnight.

	MLP	LM	PR	LG	RF	SVM
BS						
	0.0320	0.0328	0.0305	0.0305	0.0262	0.0333
	(0.0031)	(0.0025)	(0.0025)	(0.0025)	(0.0022)	(0.0032)
AUC						
	0.9342	0.9245	0.9419	0.9412	0.9558	0.8361
	(0.0131)	(0.0111)	(0.0087)	(0.0089)	(0.0081)	(0.0339)
BSS						
	0.2311	0.2122	0.2681	0.2686	0.3722	0.2007
	(0.0671)	(0.0316)	(0.0391)	(0.0421)	(0.0397)	(0.0538)

Table 9 Same as Table 6, but for predictions issued only at midnight.

	MLP	LM	PR	LG	RF	SVM
BS						
	0.1142	0.1273	0.1181	0.1169	0.1023	0.1187
	(0.0041)	(0.0034)	(0.0037)	(0.0038)	(0.0037)	(0.0041)
AUC						
	0.8771	0.8673	0.8696	0.8696	0.9004	0.8620
	(0.0078)	(0.0079)	(0.0079)	(0.0079)	(0.0074)	(0.0086)
BSS						
	0.3970	0.3281	0.3767	0.3826	0.4597	0.3735
	(0.0190)	(0.0143)	(0.0163)	(0.0169)	(0.0169)	(0.0172)

significant differences are observed for BSS since the associated climatology is also different. Nevertheless, the finding that RF is the best overall method continues to hold.

Similar conclusions can be drawn for the $> C1$ flare prediction case, that is, through Tables 9 and 6. Here, noticeably, not even the BSS changes significantly, since the underlying climatology seems similar in both cases. This is because in contrast to $> M1$ class flares with a mean frequency of $\approx 5\%$, $> C1$ class flares show a mean frequency of $\approx 25\%$.

Finally, Tables 10 and 11 present the skill scores ACC, TSS, and HSS for the midnight prediction scenario analogously to Tables 2 and 3 for the baseline scenario. For completeness, we note that Table 10 pertains to $> M1$ flare prediction and Table 11 to $> C1$ flare prediction. We see that on average, the issuing of midnight-only predictions does not change the ACC, TSS, and HSS much with respect to the probability threshold. For example, on $> M1$ flare midnight-only predictions, the RF provides $ACC = 0.93 \pm 0.01$, $TSS = 0.73 \pm 0.04$, and $HSS = 0.47 \pm 0.03$ for a probability threshold of 15%. Furthermore, for $> C1$ flare midnight-only predictions, the RF yields $ACC = 0.85 \pm 0.01$, $TSS = 0.63 \pm 0.02$, and $HSS = 0.61 \pm 0.02$ for a probability threshold of 35%.

Appendix D: Concluding Remarks on ML *versus* Statistical Methods for Flare Forecasting

In order to assess the overall forecasting ability of ML *versus* statistical approaches in our dataset and problem definition, we employed the weighted-sum (WS) multicriteria ranking approach (Greco, Figueira, and Ehrgott, 2016), using a composite index (CI) defined in Equation D.1:

$$CI = \frac{1}{3} \left(\frac{ACC - ACC_{\min}}{ACC_{\max} - ACC_{\min}} + \frac{TSS - TSS_{\min}}{TSS_{\max} - TSS_{\min}} + \frac{HSS - HSS_{\min}}{HSS_{\max} - HSS_{\min}} \right). \quad (D.1)$$

The CI value was computed for $6 \times 21 = 126$ probabilistic classifiers using the set of methods {MLP, LM, PR, LG, RF, SVM} and a probability threshold grid of 5%. Then, the 126 probabilistic classifiers were ranked in non-increasing values of the CI index. A normalization was made for ACC, TSS, and HSS, so that each metric over the set of alternatives took values in the range [0, 1]. The normalization is useful because the range of values for ACC is different from the range of values for TSS and HSS. Moreover, ACC_{\min} is the minimum of ACC over all 126 alternative models. Likewise, ACC_{\max} is the maximum ACC obtained over all 126 alternative models. Similar facts hold for TSS_{\min} , TSS_{\max} , HSS_{\min} , and HSS_{\max} . Analytically, Table 12 presents the results of the multicriteria ranking approach for all methods we used with various probability thresholds, especially for the $> C1$ flare forecasting case. Table 13 conveys a similar ranking of all methods developed in this paper, but for the $> M1$ flare prediction.

Figure 11 summarizes the results shown in Tables 12 and 13, so that the differences between ML and statistical methods are highlighted (*e.g.* see Figures 11b and 11d). Similarly, conclusions for the merit of all methods developed in this paper can be drawn in Figures 11a and 11c. The top 100τ percentile methods are those ranked in the corresponding positions of Tables 12 and 13. For example, the top 16.6%(1/6) methods are those ranked in positions 1–21. For low values of τ , we obtain the best methods designated as the top $100\tau\%$ methods. From Figure 11 we see that both for $> C1$ and $> M1$ flares, the RF has the greatest frequency in the top 16.6% percentile of methods, with a frequency of 33.3%. This means that in Tables 12 and 13, in positions 1–21, the RF method appears seven times in each

Table 10 Same as Table 2, but for predictions issued only at midnight.

Par	%	MLP		LM		PR		LG		RF		SVM				
		ACC	TSS	HSS	ACC	TSS	HSS	ACC	TSS	HSS	ACC	TSS	HSS			
val0	0.00	0.18	0.14	0.01	0.33	0.29	0.04	0.00	0.00	0.00	0.50	0.47	0.07	0.00	0.00	
val5	0.05	0.90	0.70	0.37	0.76	0.69	0.19	0.84	0.85	0.77	0.30	0.85	0.78	0.31	0.58	0.42
val10	0.10	0.93	0.66	0.43	0.88	0.72	0.32	0.90	0.90	0.73	0.37	0.90	0.78	0.41	0.95	0.44
val15	0.15	0.94	0.62	0.46	0.92	0.59	0.38	0.93	0.93	0.65	0.42	0.93	0.73	0.47	0.96	0.42
val20	0.20	0.95	0.58	0.47	0.94	0.45	0.38	0.94	0.94	0.56	0.43	0.94	0.67	0.52	0.96	0.38
val25	0.25	0.95	0.55	0.48	0.95	0.36	0.37	0.95	0.95	0.49	0.44	0.95	0.61	0.53	0.96	0.42
val30	0.30	0.95	0.52	0.48	0.95	0.30	0.36	0.96	0.96	0.43	0.44	0.96	0.54	0.53	0.96	0.40
val35	0.35	0.96	0.49	0.47	0.96	0.22	0.29	0.96	0.96	0.35	0.40	0.96	0.48	0.52	0.96	0.39
val40	0.40	0.96	0.46	0.47	0.96	0.19	0.27	0.96	0.96	0.31	0.39	0.96	0.42	0.50	0.96	0.38
val45	0.45	0.96	0.43	0.46	0.96	0.16	0.23	0.96	0.96	0.26	0.35	0.96	0.38	0.48	0.96	0.36
val50	0.50	0.96	0.40	0.45	0.96	0.13	0.21	0.96	0.96	0.21	0.31	0.96	0.34	0.45	0.96	0.35
val55	0.55	0.96	0.37	0.43	0.96	0.09	0.16	0.96	0.96	0.18	0.27	0.96	0.29	0.41	0.96	0.34
val60	0.60	0.96	0.34	0.41	0.96	0.07	0.11	0.96	0.96	0.14	0.23	0.96	0.25	0.37	0.96	0.33
val65	0.65	0.96	0.32	0.40	0.96	0.06	0.10	0.96	0.96	0.11	0.18	0.96	0.21	0.32	0.96	0.31
val70	0.70	0.96	0.29	0.38	0.96	0.05	0.10	0.96	0.96	0.09	0.15	0.96	0.16	0.26	0.96	0.30
val75	0.75	0.96	0.26	0.36	0.94	0.04	0.08	0.96	0.96	0.07	0.13	0.96	0.12	0.20	0.96	0.28
val80	0.80	0.96	0.23	0.33	0.87	0.02	0.04	0.96	0.96	0.05	0.09	0.96	0.09	0.15	0.96	0.26
val85	0.85	0.96	0.20	0.30	0.66	0.01	0.02	0.95	0.95	0.04	0.06	0.95	0.07	0.12	0.96	0.23
val90	0.90	0.96	0.17	0.25	0.50	0.01	0.01	0.91	0.91	0.02	0.03	0.91	0.04	0.07	0.96	0.20
val95	0.95	0.96	0.12	0.19	0.48	0.01	0.01	0.54	0.65	0.00	0.00	0.53	0.01	0.01	0.96	0.14
val100	1.00	0.00	0.00	0.00	0.45	0.01	0.01	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00

Table 11 Same as Table 3, but for predictions issued only at midnight.

Par	%	MLP		LM		PR		LG		RF		SVM	
		ACC	TSS	HSS	ACC	TSS	HSS	ACC	TSS	HSS	ACC	TSS	HSS
val ₀	0.00	0.00	0.00	0.00	0.18	0.10	0.00	0.00	0.00	0.00	0.27	0.02	0.01
val ₅	0.05	0.53	0.35	0.22	0.44	0.14	0.48	0.29	0.17	0.50	0.19	0.31	0.19
val ₁₀	0.10	0.67	0.51	0.37	0.51	0.20	0.59	0.41	0.28	0.61	0.44	0.44	0.30
val ₁₅	0.15	0.73	0.56	0.44	0.58	0.40	0.66	0.49	0.35	0.68	0.51	0.38	0.72
val ₂₀	0.20	0.77	0.59	0.49	0.65	0.49	0.73	0.55	0.43	0.74	0.56	0.45	0.77
val ₂₅	0.25	0.80	0.59	0.52	0.73	0.55	0.78	0.57	0.49	0.78	0.58	0.50	0.81
val ₃₀	0.30	0.81	0.58	0.54	0.79	0.50	0.80	0.57	0.53	0.81	0.57	0.53	0.83
val ₃₅	0.35	0.83	0.57	0.56	0.82	0.53	0.82	0.56	0.54	0.82	0.56	0.54	0.85
val ₄₀	0.40	0.83	0.55	0.56	0.83	0.54	0.84	0.54	0.55	0.84	0.54	0.56	0.86
val ₄₅	0.45	0.84	0.53	0.55	0.84	0.53	0.84	0.51	0.55	0.84	0.52	0.56	0.86
val ₅₀	0.50	0.84	0.51	0.55	0.83	0.48	0.84	0.48	0.54	0.84	0.50	0.55	0.86
val ₅₅	0.55	0.84	0.48	0.54	0.82	0.33	0.84	0.44	0.51	0.84	0.46	0.46	0.86
val ₆₀	0.60	0.84	0.45	0.51	0.81	0.28	0.83	0.38	0.45	0.83	0.41	0.48	0.85
val ₆₅	0.65	0.83	0.42	0.49	0.79	0.20	0.82	0.33	0.41	0.82	0.36	0.36	0.85
val ₇₀	0.70	0.83	0.37	0.45	0.78	0.15	0.81	0.28	0.35	0.81	0.31	0.39	0.84
val ₇₅	0.75	0.82	0.32	0.40	0.77	0.12	0.80	0.22	0.29	0.80	0.26	0.26	0.82
val ₈₀	0.80	0.81	0.27	0.35	0.77	0.10	0.78	0.17	0.23	0.79	0.20	0.20	0.81
val ₈₅	0.85	0.80	0.22	0.29	0.77	0.09	0.78	0.14	0.19	0.78	0.15	0.15	0.80
val ₉₀	0.90	0.78	0.14	0.20	0.76	0.07	0.77	0.11	0.15	0.77	0.11	0.11	0.78
val ₉₅	0.95	0.73	0.06	0.08	0.76	0.06	0.76	0.07	0.10	0.76	0.07	0.07	0.76
val ₁₀₀	1.00	0.00	0.00	0.00	0.76	0.05	0.07	0.00	0.00	0.00	0.00	0.00	0.00

Table 12 Ranking of all models with ML and statistical methods with the multicriteria WS method with respect to the three criteria ACC, TSS, and HSS and using a weight vector equal to $w = [1/3, 1/3, 1/3]$ for $> C1$ flare forecasting. The methods are ranked in decreasing order of CI for varying levels of probability thresholds used, with a grid of 5% for the probability thresholds. The six best positions of the ranking are covered by the RF method for various probability thresholds.

rank	model	CI	rank	model	CI	rank	model	CI
PANEL A $> C1$ flares								
1	RF-val35	0.983	22	PR-val30	0.919	43	MLP-val60	0.856
2	RF-val40	0.983	23	LG-val45	0.914	44	SVM-val55	0.849
3	RF-val45	0.971	24	LM-val40	0.907	45	SVM-val15	0.843
4	RF-val30	0.970	25	MLP-val50	0.907	46	LG-val20	0.842
5	RF-val50	0.953	26	PR-val45	0.907	47	PR-val55	0.840
6	RF-val25	0.938	27	MLP-val25	0.906	48	LM-val25	0.835
7	LG-val35	0.935	28	SVM-val40	0.906	49	PR-val20	0.828
8	MLP-val35	0.934	29	LM-val30	0.898	50	SVM-val60	0.824
9	PR-val35	0.933	30	RF-val60	0.898	51	MLP-val65	0.821
10	MLP-val40	0.932	31	LG-val25	0.892	52	MLP-val15	0.821
11	SVM-val25	0.932	32	LG-val50	0.889	53	LG-val60	0.818
12	LG-val40	0.930	33	SVM-val45	0.887	54	RF-val15	0.810
13	RF-val55	0.929	34	PR-val25	0.886	55	RF-val70	0.809
14	SVM-val30	0.928	35	RF-val20	0.886	56	LM-val50	0.803
15	MLP-val30	0.925	36	MLP-val55	0.885	57	SVM-val65	0.789
16	PR-val40	0.925	37	PR-val50	0.877	58	PR-val60	0.787
17	LM-val35	0.923	38	MLP-val20	0.875	59	MLP-val70	0.778
18	MLP-val45	0.923	39	LM-val45	0.869	60	LG-val65	0.768
19	LG-val30	0.921	40	SVM-val50	0.867	61	SVM-val70	0.751
20	SVM-val35	0.920	41	LG-val55	0.861	62	LG-val15	0.748
21	SVM-val20	0.920	42	RF-val65	0.859	63	RF-val75	0.747
PANEL B $> C1$ flares								
64	PR-val65	0.738	85	LG-val80	0.593	106	PR-val95	0.420
65	LM-val55	0.735	86	LM-val15	0.592	107	LG-val95	0.417
66	MLP-val75	0.727	87	RF-val85	0.587	108	RF-val95	0.410
67	PR-val15	0.724	88	PR-val80	0.566	109	LM-val90	0.405
68	MLP-val10	0.720	89	SVM-val90	0.552	110	MLP-val95	0.400
69	LM-val20	0.718	90	LM-val70	0.537	111	LM-val95	0.387
70	LG-val70	0.715	91	LG-val85	0.533	112	LM-val5	0.372
71	SVM-val75	0.709	92	PR-val85	0.515	113	LM-val100	0.370
72	RF-val10	0.699	93	MLP-val90	0.510	114	SVM-val10	0.363
73	PR-val70	0.682	94	RF-val5	0.506	115	LM-val0	0.291
74	RF-val80	0.671	95	MLP-val5	0.501	116	SVM-val5	0.142
75	SVM-val80	0.667	96	RF-val90	0.500	117	RF-val0	0.131
76	MLP-val80	0.666	97	LM-val75	0.494	118	MLP-val0	0.000
77	LM-val60	0.664	98	LM-val10	0.488	119	MLP-val100	0.000
78	LG-val75	0.657	99	LG-val90	0.484	120	PR-val0	0.000
79	PR-val75	0.622	100	SVM-val95	0.484	121	PR-val100	0.000
80	LG-val10	0.621	101	PR-val90	0.476	122	LG-val0	0.000
81	SVM-val85	0.613	102	LM-val80	0.461	123	LG-val100	0.000
82	MLP-val85	0.596	103	LG-val5	0.441	124	RF-val100	0.000
83	PR-val10	0.594	104	LM-val85	0.431	125	SVM-val0	0.000
84	LM-val65	0.593	105	PR-val5	0.427	126	SVM-val100	0.000

Table 13 Same as Table 12, but for > M1 flare forecasting.

rank	model	CI	rank	model	CI	rank	model	CI
PANEL A > M1 flares								
1	RF-val20	0.938	22	PR-val25	0.839	43	SVM-val25	0.774
2	RF-val25	0.932	23	MLP-val5	0.836	44	MLP-val55	0.773
3	RF-val15	0.927	24	MLP-val35	0.836	45	LG-val40	0.770
4	RF-val30	0.912	25	SVM-val10	0.834	46	SVM-val30	0.756
5	RF-val10	0.890	26	LG-val30	0.827	47	MLP-val60	0.754
6	RF-val35	0.883	27	LM-val20	0.823	48	LM-val30	0.749
7	MLP-val15	0.870	28	LM-val10	0.822	49	PR-val40	0.749
8	MLP-val20	0.867	29	MLP-val40	0.820	50	LG-val45	0.743
9	MLP-val10	0.866	30	PR-val30	0.817	51	SVM-val35	0.741
10	MLP-val25	0.859	31	SVM-val15	0.816	52	RF-val55	0.736
11	PR-val15	0.856	32	RF-val45	0.815	53	MLP-val65	0.735
12	PR-val20	0.856	33	RF-val5	0.809	54	SVM-val40	0.725
13	LG-val15	0.854	34	MLP-val45	0.805	55	PR-val45	0.713
14	LG-val20	0.851	35	LG-val35	0.799	56	MLP-val70	0.713
15	RF-val40	0.851	36	LG-val5	0.796	57	LG-val50	0.709
16	MLP-val30	0.849	37	SVM-val20	0.792	58	SVM-val45	0.709
17	LM-val15	0.849	38	MLP-val50	0.790	59	LM-val5	0.699
18	SVM-val5	0.848	39	PR-val5	0.786	60	LM-val35	0.694
19	PR-val10	0.845	40	LM-val25	0.786	61	RF-val60	0.693
20	LG-val10	0.843	41	PR-val35	0.782	62	SVM-val50	0.692
21	LG-val25	0.840	42	RF-val50	0.778	63	MLP-val75	0.689
PANEL B > M1 flares								
64	PR-val50	0.679	85	LM-val50	0.571	106	LM-val80	0.424
65	SVM-val55	0.679	86	SVM-val85	0.550	107	LG-val90	0.413
66	LG-val55	0.676	87	PR-val70	0.546	108	PR-val90	0.410
67	MLP-val80	0.663	88	LM-val55	0.546	109	LM-val85	0.404
68	SVM-val60	0.662	89	RF-val75	0.543	110	RF-val90	0.392
69	PR-val55	0.648	90	LG-val75	0.539	111	LM-val90	0.381
70	LM-val40	0.646	91	MLP-val95	0.539	112	PR-val95	0.378
71	RF-val65	0.645	92	LM-val60	0.521	113	LG-val95	0.372
72	SVM-val65	0.644	93	PR-val75	0.516	114	LM-val95	0.364
73	LG-val60	0.644	94	SVM-val90	0.512	115	LM-val100	0.352
74	MLP-val85	0.633	95	LG-val80	0.506	116	RF-val95	0.334
75	SVM-val70	0.625	96	RF-val80	0.492	117	LM-val0	0.251
76	PR-val60	0.612	97	LM-val65	0.489	118	MLP-val0	0.108
77	LG-val65	0.607	98	PR-val80	0.483	119	MLP-val100	0.000
78	SVM-val75	0.604	99	LG-val85	0.466	120	PR-val0	0.000
79	LM-val45	0.603	100	LM-val70	0.461	121	PR-val100	0.000
80	MLP-val90	0.594	101	SVM-val95	0.450	122	LG-val0	0.000
81	RF-val70	0.593	102	PR-val85	0.444	123	LG-val100	0.000
82	SVM-val80	0.579	103	RF-val85	0.442	124	RF-val100	0.000
83	PR-val65	0.576	104	LM-val75	0.441	125	SVM-val0	0.000
84	LG-val70	0.572	105	RF-val0	0.433	126	SVM-val100	0.000

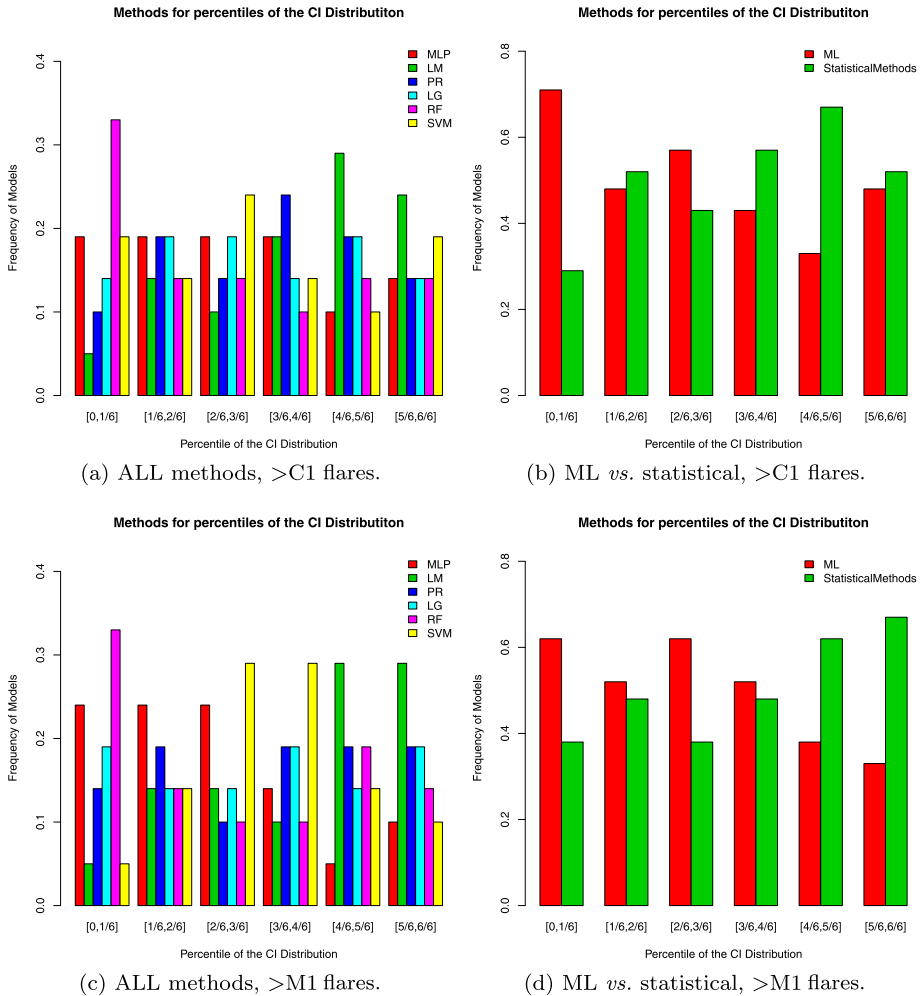


Figure 11 Descriptive statistics on the frequency with which every forecasting method for any probability threshold presents itself in the top 100τ% percentile of the CI distribution. Panels (a) and (c) describe frequencies for all methods, and panels (b) and (d) group the results by category of methods (e.g. ML vs. statistical methods). For example, for >C1 flares in panel (a), the top 16.6% methods are dominated by RF with a frequency of $7/21 = 33\%$. Likewise, for >M1 flares in panel (c), the top 16.6% methods are again dominated by RF with a frequency of $7/21 = 33\%$.

table. We also see in Figure 11b that for >C1 flares, the top 16.6% methods are of type ML with a frequency 71% (versus 29% for statistical methods). Similarly, in Figure 11d, ML dominates in the top 16.6% methods with a frequency of 62% (versus 38% for statistical methods).

References

Ahmed, O., Qahwaji, R., Colak, T., Dudok De Wit, T., Ipson, S.: 2010, A new technique for the calculation and 3D visualisation of magnetic complexities on solar satellite images. *Vis. Comput.* **26**, 385. DOI.

- Ahmed, O.W., Qahwaji, R., Colak, T., Higgins, P.A., Gallagher, P.T., Bloomfield, D.S.: 2013, Solar flare prediction using advanced feature extraction, machine learning, and feature selection. *Solar Phys.* **283**(1), 157. DOI.
- Al-Ghraibah, A., Boucheron, L.E., McAteer, R.T.J.: 2015, An automated classification approach to ranking photospheric proxies of magnetic energy build-up. *Astron. Astrophys.* **579**, A64. DOI.
- Alissandrakis, C.E.: 1981, On the computation of constant alpha force-free magnetic field. *Astron. Astrophys.* **100**, 197.
- Barnes, G., Longcope, D.W., Leka, K.D.: 2005, Implementing a magnetic charge topology model for solar active regions. *Astrophys. J.* **629**, 561. DOI.
- Barnes, G., Schanche, N., Leka, K., Aggarwal, A., Reeves, K.: 2016, A comparison of classifiers for solar energetic events. *Proc. Int. Astron. Union* **12**(S325), 201. DOI.
- Bloomfield, D.S., Higgins, P.A., McAteer, R.T.J., Gallagher, P.T.: 2012, Toward reliable benchmarking of solar flare forecasting methods. *Astrophys. J. Lett.* **747**(2). DOI.
- Bobra, M.G., Couvidat, S.: 2015, Solar flare prediction using SDO/HMI vector magnetic field data with a machine-learning algorithm. *Astrophys. J.* **798**(2), 135. DOI.
- Bobra, M.G., Sun, X., Hoeksema, J.T., Turmon, M., Liu, Y., Hayashi, K., Barnes, G., Leka, K.D.: 2014, The Helioseismic and Magnetic Imager (HMI) vector magnetic field pipeline: SHARPs – Space-weather HMI active region patches. *Solar Phys.* **289**(9), 3549. DOI.
- Boucheron, L.E., Al-Ghraibah, A., McAteer, R.T.J.: 2015, Prediction of solar flare size and time-to-flare using support vector machine regression. *Astrophys. J.* **812**(1), 51. DOI.
- Breiman, L.: 2001, Random forests. *Mach. Learn.* **45**(1), 5. DOI.
- Breiman, L., Friedman, J., Stone, C.J., Olshen, R.A.: 1984, *Classification and Regression Trees*.
- Brier, G.W.: 1950, Verification of forecasts expressed in terms of probability. *Mon. Weather Rev.* **78**, 1. DOI.
- Chang, Y.-W., Lin, C.-J.: 2008, Feature ranking using linear svm. In: Guyon, I., Aliferis, C., Cooper, G., Elisseeff, A., Pellet, J.-P., Spirtes, P., Statnikov, A. (eds.) *Proceedings of the Workshop on the Causation and Prediction Challenge at WCCI 2008, Proceedings of Machine Learning Research* **3**, 53.
- Chang, C.-C., Lin, C.-J.: 2011, LIBSVM: A library for support vector machines. *ACM Trans. Intell. Syst. Technol.* **2**, 27:1. Software available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm>. DOI.
- Chen, Y.-W., Lin, C.-J.: 2006, Combining SVMs with various feature selection strategies. In: *Feature Extraction*, Springer, Berlin, 315. DOI.
- Colak, T., Qahwaji, R.: 2009, Automated solar activity prediction: A hybrid computer platform using machine learning and solar imaging for automated prediction of solar flares. *Space Weather* **7**(6), S06001. DOI.
- de Souza, R.S., Cameron, E., Kiledar, M., Hilbe, J., Vilalta, R., Maio, U., Biffi, V., Ciardi, B., Riggs, J.D.: 2015, The overlooked potential of generalized linear models in astronomy, I: binomial regression. *Astron. Comput.* **12**, 21. DOI.
- Falconer, D.A., Moore, R.L., Gary, G.A.: 2008, Magnetogram measures of total nonpotentiality for prediction of solar coronal mass ejections from active regions of any degree of magnetic complexity. *Astrophys. J.* **689**, 1433. DOI.
- Falconer, D., Barghouty, A.F., Khazanov, I., Moore, R.: 2011, A tool for empirical forecasting of major flares, coronal mass ejections, and solar particle events from a proxy of active-region free magnetic energy. *Space Weather* **9**(4), n/a. S04003. DOI.
- Falconer, D.A., Moore, R.L., Barghouty, A.F., Khazanov, I.: 2012, Prior flaring as a complement to free magnetic energy for forecasting solar eruptions. *Astrophys. J.* **757**, 32. DOI.
- Falconer, D.A., Moore, R.L., Barghouty, A.F., Khazanov, I.: 2014, MAG4 versus alternative techniques for forecasting active region flare productivity. *Space Weather* **12**(5), 306. DOI.
- Fang, F., Manchester, W. IV, Abbott, W.P., van der Holst, B.: 2012, Buildup of magnetic shear and free energy during flux emergence and cancellation. *Astrophys. J.* **754**, 15. DOI.
- Fernández-Delgado, M., Cernadas, E., Barro, S., Amorim, D.: 2014, Do we need hundreds of classifiers to solve real world classification problems? *J. Mach. Learn. Res.* **15**, 3133.
- Fletcher, L., Dennis, B.R., Hudson, H.S., Krucker, S., Phillips, K., Veronig, A., Battaglia, M., Bone, L., Caspi, A., Chen, Q., Gallagher, P., Grigis, P.T., Ji, H., Liu, W., Milligan, R.O., Temmer, M.: 2011, An observational overview of solar flares. *Space Sci. Rev.* **159**, 19. DOI.
- Georgoulis, M.K., Rust, D.M.: 2007, Quantitative forecasting of major solar flares. *Astrophys. J. Lett.* **661**, L109. DOI.
- Granett, B.R.: 2017, Probing the sparse tails of redshift distributions with Voronoi tessellations. *Astron. Comput.* **18**, 18. DOI.
- Greco, S., Figueira, J., Ehrgott, M.: 2016, *Multiple Criteria Decision Analysis*, 2nd edn.
- Greene, W.H.: 2002, *Econometric Analysis*, 5th edn.
- Guerra, J.A., Pulkkinen, A., Uritsky, V.M., Yashiro, S.: 2015, Spatio-temporal scaling of turbulent photospheric line-of-sight magnetic field in active region NOAA 11158. *Solar Phys.* **290**, 335. DOI.

- Hall, M., Frank, E., Holmes, G., Pfahringer, B., Reutemann, P., Witten, I.H.: 2009, The WEKA data mining software: An update. *ACM SIGKDD Explor. Newsl.* **11**(1), 10. [DOI](#).
- Hanssen, A., Kuipers, W.: 1965, *On the Relationship Between the Frequency of Rain and Various Meteorological Parameters: (with Reference to the Problem of Objective Forecasting)*.
- Hastie, T., Tibshirani, R., Friedman, J.: 2009, *The Elements of Statistical Learning: Data Mining, Inference and Prediction*, 2nd edn.
- Heidke, P.: 1926, Berechnung des erfolges und der güte der windstärkevorhersagen im sturmwarnungsdienst. *Geogr. Ann.* **8**, 301. [DOI](#).
- Hornik, K., Stinchcombe, M., White, H.: 1989, Multilayer feedforward networks are universal approximators. *Neural Netw.* **2**(5), 359. [DOI](#).
- Kliem, B., Török, T.: 2006, Torus instability. *Phys. Rev. Lett.* **96**(25), 255002. [DOI](#).
- Laboratory, N.-R.A.: 2015, Verification: Weather forecast verification utilities. R package version 1.42.
- Li, R., Cui, Y., He, H., Wang, H.: 2008, Application of support vector machine combined with k-nearest neighbors in solar flare and solar proton events forecasting. *Adv. Space Res.* **42**(9), 1469. [DOI](#).
- Liaw, A., Wiener, M.: 2002, Classification and regression by randomforest. *R News* **2**(3), 18.
- Liu, C., Deng, N., Wang, J.T.L., Wang, H.: 2017, Predicting solar flares using SDO/HMI vector magnetic data products and the random forest algorithm. *Astrophys. J.* **843**(2), 104. [DOI](#).
- MacKay, D.J.C.: 2003, *Information Theory, Inference, and Learning Algorithms*.
- Marzban, C.: 2004, The ROC curve and the area under it as performance measures. *Weather Forecast.* **19**(6), 1106. [DOI](#).
- Meyer, D., Leisch, F., Hornik, K.: 2003, The support vector machine under test. *Neurocomputing* **55**(1), 169. [DOI](#).
- Meyer, D., Dimitriadou, E., Hornik, K., Weingessel, A., Leisch, F.: 2015, e1071: Misc functions of the department of statistics, probability theory group (formerly: E1071), TU Wien. R package version 1.6-7.
- Pesnell, W.D., Thompson, B.J., Chamberlin, P.C.: 2012, The Solar Dynamics Observatory (SDO). *Solar Phys.* **275**, 3. [DOI](#).
- Prieto, A., Prieto, B., Ortigosa, E.M., Ros, E., Pelayo, F., Ortega, J., Rojas, I.: 2016, Neural networks: An overview of early research, current frameworks and new challenges. *Neurocomputing* **214**, 242. [DOI](#).
- R Core Team: 2016, *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna.
- Scherrer, P.H., Bogart, R.S., Bush, R.I., Hoeksema, J.T., Kosovichev, A.G., Schou, J., *et al.*: 1995, The solar oscillations investigation – Michelson Doppler imager. *Solar Phys.* **162**, 129. [DOI](#).
- Scherrer, P.H., Schou, J., Bush, R.I., Kosovichev, A.G., Bogart, R.S., Hoeksema, J.T., Liu, Y., Duvall, T.L., Zhao, J., Title, A.M., Schrijver, C.J., Tarbell, T.D., Tomczyk, S.: 2012, The Helioseismic and Magnetic Imager (HMI) investigation for the Solar Dynamics Observatory (SDO). *Solar Phys.* **275**, 207. [DOI](#).
- Schrijver, C.J.: 2007, A characteristic magnetic field pattern associated with all major solar flares and its use in flare forecasting. *Astrophys. J. Lett.* **655**, L117. [DOI](#).
- Schuh, M.A., Angryk, R.A., Martens, P.C.: 2015, Solar image parameter data from the SDO: Long-term curation and data mining. *Astron. Comput.* **13**, 86. [DOI](#).
- Sing, T., Sander, O., Beerenwinkel, N., Lengauer, T.: 2005, ROCr: Visualizing classifier performance in R. *Bioinformatics* **21**(20), 7881. [DOI](#).
- Song, H., Tan, C., Jing, J., Wang, H., Yurchyshyn, V., Abramenko, V.: 2009, Statistical assessment of photospheric magnetic features in imminent solar flare predictions. *Solar Phys.* **254**(1), 101. [DOI](#).
- Vapnik, V.: 1998, *Statistical Learning Theory*.
- Venables, W.N., Ripley, B.D.: 2002, *Modern Applied Statistics with S*, 4th edn. Springer, New York. [DOI](#).
- Vilalta, R., Gupta, K.D., Macri, L.: 2013, A machine learning approach to Cepheid variable star classification using data alignment and maximum likelihood. *Astron. Comput.* **2**, 46. [DOI](#).
- Wang, H.N., Cui, Y.M., Li, R., Zhang, L.Y., Han, H.: 2008, Solar flare forecasting model supported with artificial neural network techniques. *Adv. Space Res.* **42**(9), 1464. [DOI](#).
- Wilks, D.S.: 2011, *Statistical Methods in the Atmospheric Sciences* **100**.
- Winkelmann, R., Boes, S.: 2006, *Analysis of Microdata*, Springer, Berlin.
- Witten, I.H., Frank, E., Hall, M.A., Pal, C.J.: 2016, *Data Mining: Practical Machine Learning Tools and Techniques*.
- Yu, D., Huang, X., Wang, H., Cui, Y.: 2009, Short-term solar flare prediction using a sequential supervised learning method. *Solar Phys.* **255**(1), 91. [DOI](#).
- Yuan, Y., Shih, F.Y., Jing, J., Wang, H.-M.: 2010, Automated flare forecasting using a statistical learning technique. *Res. Astron. Astrophys.* **10**(8), 785. [DOI](#).
- Zuccarello, F.P., Aulanier, G., Gilchrist, S.A.: 2015, Critical decay index at the onset of solar eruptions. *Astrophys. J.* **814**, 126. [DOI](#).