

Interpretable Machine Learning to Forecast SEP Events for Solar Cycle 23

Spiridon Kasapis¹, Lulu Zhao², Yang Chen³, Xiantong Wang², Monica Bobra⁴,
Tamas Gombosi²

¹Department of Naval Architecture and Marine Engineering, University of Michigan

²Department of Climate and Space Sciences and Engineering, University of Michigan

³Department of Statistics, University of Michigan

⁴Hansen Experimental Physics Laboratory, Stanford University

Key Points:

- In an experimental setting, SMARP data can correctly predict whether a solar flare will lead to a solar energetic particle (SEP) event 72% of the times.
- Flare peak intensity is the strongest SEP predictor and can be coupled with SMARP data to achieve accuracy $\leq 0.92 \pm 0.07$.
- The SMARP dataset provides a leading time of 55.3 ± 28.6 minutes for forecasting the SEP events.

Corresponding author: Spiridon Kasapis, skasapis@umich.edu

This is the author manuscript accepted for publication and has undergone full peer review but has not been through the copyediting, typesetting, pagination and proofreading process, which may lead to differences between this version and the [Version of Record](#). Please cite this article as [doi: 10.1029/2021SW002842](https://doi.org/10.1029/2021SW002842).

This article is protected by copyright. All rights reserved.

Abstract

We use machine learning methods to predict whether an active region (AR) which produces flares will lead to a solar energetic particle (SEP) event using Space-Weather Michelson Doppler Imager (MDI) Active Region Patches (SMARPs). This new data product is derived from maps of the solar surface magnetic field taken by the Michelson Doppler Imager (MDI) aboard the Solar and Heliospheric Observatory (SOHO). We survey the SMARP active regions associated with flares that appear on the solar disk between June 5, 1996 and August 14, 2010, label those that produced SEPs as positive and the rest as negative. The AR SMARP features that correspond to each flare are used to train two different types of machine learning methods, the support vector machines (SVMs) and the regression models. The results show that the SMARP data can predict whether a flare will lead to an SEP with accuracy (ACC) $\leq 0.72 \pm 0.12$ while allowing for a competitive leading time of 55.3 ± 28.6 minutes for forecasting the SEP events.

Plain Language Summary

We train machine learning algorithms to predict whether sunspots on the solar disk will produce solar phenomena such as Solar Energetic Particle (SEP) events. We use a newly published piece of solar data (SMARP) captured by the Space-Weather Michelson Doppler Imager (MDI), an instrument aboard the Solar and Heliospheric Observatory (SOHO) which is a spacecraft launched on 1995 to study the Sun. The results show that the SMARP data can correctly predict 72% of the times whether a sunspot will lead to an SEP while allowing for a potentially competitive forecast window.

1 Introduction

Large solar eruptions can potentially harm modern civilization in several different ways. Events such as large solar flares and coronal mass ejections (CMEs) that lead to solar particle acceleration, can adversely affect the near-earth environment, degrade high frequency (HF) radio communications, incapacitate satellites, expose airline passengers to elevated radiation levels and even endanger life in outer space. Therefore, predicting and monitoring such events is an important task for the community.

Solar Energetic Particles are rare events that involve protons, electrons and heavy ions accelerated to high energies (up to tens of GeV while the fastest ones can accelerate to speeds of up to 80% of the speed of light) by two solar processes (Reames, 2013), the energization at a solar flare site or the shock waves associated with Coronal Mass Ejections (CMEs). Solar charged particles are accelerated in flares or CME shock waves (Wild et al., 1963) and travel preferentially along the interplanetary magnetic field to their detection point in space (McCracken & Ness, 1966).

The study of solar energetic particle (SEP) events is a relatively recent science as the identification of the first event took place on 28 February, 1942 (Forbush, 1946). Observations of solar proton events (subset of solar energetic particle events) were made using ground-based instruments that detected ionization, neutrons, or radio disturbances caused by them. The largest solar proton event recorded using these modern techniques (particles exceeded 15 GeV at the top of the atmosphere) was on the 23rd of February, 1956. In the mid-1960s spacecraft were deployed that began directly measuring solar proton events. This was also the time when the first flare was associated with an SEP event (Shea & Smart, 1995).

During the so-called Halloween storms in late October 2003, SEP events caused 47 satellites to report malfunctions, more than 10 satellites to go out of action for days, the Mars Odyssey spacecraft went into deep safe mode (Lopez et al., 2004), a Japanese satellite costing 640m USD was completely lost, the US FAA issued their first-ever high ra-

64 diation dosage alert for high-altitude aircraft, and astronauts in the ISS had to seek safety
65 into their heavily shielded service module (Webb & Allen, 2004; Horne et al., 2013).

66 One of the sources of solar activity phenomena that cause SEPs are the magnetically
67 strong regions on the solar sphere that we refer to as active regions (van Driel-Gesztelyi
68 & Green, 2015). The most flare productive active regions (ARs) are the ones that under-
69 go large changes in sunspot area and show magnetic flux imbalance (Choudhary et
70 al., 2013). Large active regions are also generally strong, produce a number of flares, evolve
71 rapidly and their lifetime spans from days to months (Choudhary et al., 2013). Using
72 instruments carried onboard satellites such as the Michelson Doppler Imager (MDI) on
73 the Solar and Heliospheric Observatory (SOHO) or the Helioseismic and Magnetic Im-
74 ager (HMI) on the Solar Dynamics Observatory (SDO), we are able to retrieve compo-
75 nents of the magnetic field at the solar surface, allowing us to calculate physical char-
76 acteristics of the ARs (Scherrer et al., 1995; Schou et al., 2012).

77 Solar particle prediction studies mainly use the flare and nearEarth space environ-
78 ment data to forecast SEP events given the knowledge that large SEPs are almost al-
79 ways accompanied by a flare (Schrijver et al., 2012). Laurenza et al. (2009) used data
80 such as flare longitude, time-integrated soft X-ray intensity, and time-integrated inten-
81 sity of type III radio emission at around 1 MHz to provide short-term warnings for SEP
82 events. Similarly, Núñez (2011) used the soft X-ray, differential and integral proton fluxes
83 data to forecast the SEP events of Solar Cycle 23 recorded on the NOAA/SWPC list.
84 Although both flare and CME data are found to be useful inputs to predictive models,
85 García-Rigo et al. (2016) deemed it sufficient to only use flare properties as they noticed
86 that the CME information offers insignificant increase in SEP prediction accuracy, be-
87 cause of the difficulty to obtain real-time information on the true radial CME speeds due
88 to line-of-sight issues.

89 Recently, machine learning (ML) methods like neural networks (in the multi-layer
90 perceptron implementation), random forests, decision trees, extremely randomized trees
91 and other, have been used in predicting SEP events. The preliminary results obtained
92 by Bain et al. (2018) show that machine learning classification techniques such as the
93 logistic regression (LR), decision trees (DTs) and support vector machine (SVM) algo-
94 rithms give an improved forecasting skill over the current SWPC Proton Prediction Model
95 (Balch, 2008) based on physical parameters associated with solar flares and coronal mass
96 ejections. An even more comprehensive study that assesses the predictability of Solar
97 Energetic Particles using ML techniques was recently published by Lavasa et al. (2021).
98 In their work they conclude that random forests (RF) could be the prediction technique
99 of choice for an optimal sample comprised by both flares and CMEs while proving that
100 the most important features are the CME speed, width and flare soft X-ray (SXR) flu-
101 ence. Lastly, Sadykov et al. (2021) recently indicated the possibility of developing ro-
102 bust all-clear SPE forecasts by employing machine learning methods. Their approach
103 indicates that for AR-based predictions, it is necessary to take into account western limb
104 and far-side ARs, characteristics of the preceding proton flux represent the most valu-
105 able input for prediction, daily median characteristics of ARs and the counts of type II,
106 III, and IV radio bursts may be excluded from the forecast and that ML-based forecasts
107 outperform SWPC NOAA forecasts.

108 Different studies have used a variety of sources to obtain the data necessary for so-
109 lar particle event prediction. Richardson et al. (2018) predict the SEP events peak pro-
110 ton intensity at one energy interval (1424 MeV) using the CME data (speed, width, di-
111 rection and type II and type III radio emissions associated with the CME) in the Space
112 Weather Database of Notifications, Knowledge, Information (DONKI). Papaioannou et
113 al. (2016) have presented a catalogue which includes solar variables (such as logarithm
114 of the solar flare magnitude (\log SXRs), solar flare longitude, duration, and rise time)
115 for 314 SEP events obtained from the Energetic Particle Sensor (EPS) aboard the Geo-
116 stationary Operational Environmental Satellites (GOES; Rodriguez et al., 2014) and CME

117 data (width/size and velocity) obtained by the Large Angle and Spectrometric Coronagraph (LASCO; Brueckner et al., 1995) carried onboard the SOHO spacecraft. Using
118 this information, Papaioannou et al. (2018) classify the solar energetic particle (SEP)
119 event radiation impact with respect to the characteristics of their parent solar events while
120 attempting to infer the possible prediction of SEP events.
121

122 Similarly, Anastasiadis et al. (2017) provide full-disk Helioseismic and Magnetic
123 Imager (HMI) magnetograms to their novel integrated prediction system which nowcasts
124 SEP events. The HMI instrument aboard the Solar Dynamics Laboratory (SDO) mea-
125 sures the solar surface magnetic field from which the Space-Weather HMI Active Region
126 Patches (SHARPs) are derived. SHARPs have been used to identify flares or SEPs in
127 Chen et al. (2019) and Inceoglu et al. (2018) respectively.

128 A new data product recently published by Bobra et al. (2021) called Space-Weather
129 MDI Active Region Patches (SMARPs) will be used in this work to predict SEPs. SMARPs
130 are derived from the solar surface magnetic field taken by the Michelson Doppler Imager
131 (MDI) on the SOHO spacecraft and provide a continuous and seamless set of keywords
132 that describe every active region observed during Solar Cycle 23. The big difference be-
133 tween the HMI (Schou et al., 2012) and the MDI (Scherrer et al., 1995) is that the first
134 measures the vector magnetic field at the solar surface whereas the later only measures
135 the line-of-sight component of the solar magnetic field. The main aim of this study is
136 to evaluate the predictive power of MDI Active Region Patches (SMARPs) on SEP events
137 as it is desirable for the space weather community to explore new datasets that, when
138 used on machine learning algorithms in the future, will be able to predict when solar pro-
139 ton events will occur, along with how energetic and how intense they will be.

140 2 Database

141 In this work, we will evaluate the predictive power of the MDI solar magnetogram
142 on SEP events. In particular, we focus on whether an active region which is associated
143 with a solar flare will lead to an SEP event. Compared to other works that predict whether
144 an SEP is likely to occur in a defined future time window, our models simply forecast
145 whether flares have a resulting particle increase. To achieve this, five different predic-
146 tors obtained from the SMARP dataset (SMARP Predictors) are used, while two more
147 predictors from the NOAA solar X-ray flare dataset (Flare Predictors) are used for com-
148 parison. While we are specifically interested in the responses of the ML models when only
149 SMARP Predictors are used, the ability to forecast SEPs by using flare data will serve
150 as a baseline capability.

151 2.1 SMARP Predictors

152 The magnetogram is measured by the Michelson Doppler Imager (MDI Scherrer
153 et al., 1995) onboard SOHO between June 5, 1996 and August 14, 2010. Based on the
154 magnetogram, Bobra et al. (2021) derived a new database called Space-Weather MDI
155 Active Region Patches (SMARPs), which contains characteristics of the active regions
156 on the solar surface. A Tracked Active Region Patch (TARP) Number is assigned to each
157 active region as its identification number and a NOAA active region number, if avail-
158 able, is assigned to each active region patch. Three physical keywords, total unsigned
159 flux (USFLUXL), mean gradient of the vertical field (MEANGBL), and the logarithm
160 of the total unsigned flux near polarity inversion line (RVALUE) are calculated using the
161 pixels in the active region and stored in the SMARP header file. In addition, the SMARP
162 header file also contains four spacial features specifying the location of the correspond-
163 ing AR on the solar surface: the minimum and maximum latitude (LATDMIN, LAT-
164 DMAX) and the minimum and maximum longitude (LONDTMIN, LONDTMAX). The
165 SMARP data is available on the Joint Science Operations Center database (Mumford
166 et al., 2015; Barnes et al., 2020).

167 Besides the three physical keywords stored in the SMARP header file, we calcu-
 168 late the angular distance between the AR and the magnetic foot-point of the earth. The
 169 longitude and latitude location of the active region on the sun is approximated by the
 170 geometric center of the active region using the latitude and longitude keywords. The mag-
 171 netic foot-point of the earth on the sun is assumed to be at $W45^\circ$. Note that the mag-
 172 netic foot-point varies from event to event. One way of characterizing this variability is
 173 to calculate the magnetic foot-point location using the solar wind speed measured at 1
 174 AU assuming an ideal Parker spiral up to the solar source surface and reconstruct the
 175 coronal magnetic fields using potential field source surface model. However, interplan-
 176 etary magnetic field can also be disturbed by corotating interaction regions (CIRs), in-
 177 terplanetary coronal mass ejections (ICMEs) and other solar transient events, especially
 178 in solar maximum. In this work, for simplicity, we use $W45^\circ$ as an approximation. We
 179 also calculate the size of the active region by multiplying the difference of longitude by
 180 the difference of latitude.

181 2.2 Construction of SEP Event List

182 The SEP event list we use in this work is documented in the NOAA Space Envi-
 183 ronment Service Center website. This study is based strictly on the Solar Energetic Par-
 184 ticle categorization that NOAA SWPC follows, therefore the catalogue lists the SEP events
 185 observed by GOES that are accelerated to energies greater than 10 MeV, noting only
 186 the times when the proton flux exceeds this threshold. This means that NOAA some-
 187 times groups more than one SEP particle injection together as one event, which is an
 188 intrinsic problem of the dataset. Really strong ARs are also known to produce multi-
 189 ple SEP events, therefore there are a number of homologous events coming from the same
 190 AR. For each SEP event, a solar flare and the corresponding NOAA active region num-
 191 ber is assigned if exists. The solar flare list is obtained from the NOAA Solar Flare Data
 192 website. For each solar flare, the list contains the start, peak and end times, the peak
 193 intensity of the flare, the active region location and the corresponding NOAA active re-
 194 gion number.

195 We match the solar flare list with the SMARP database using the AR numbers.
 196 If a flare does not have a registered AR number, matching based on their occurrence time
 197 and spatial coordinates is performed. We also discard those solar flares whose flare class,
 198 coordinates or AR numbers are undefined or missing. Out of the $\sim 25,000$ flares (A,
 199 B, C, M, X) recorded during the 14 year span between 1996 and 2010, only 6,510 flares
 200 have sufficient information and therefore are matched with SMARP files.

201 During this 14 year span, 93 SEP events are detected by the GOES spacecraft. Miss-
 202 ing information about the SEP's associated flare or AR such as the Location and its Im-
 203 portance (Xray/Opt), leave only 70 SEPs with information adequate to label the 6,510
 204 flares. We assign a label to each flare: Positive if it led to an SEP and Negative if it did
 205 not. An additional 5 SEP-flare couples were discarded due to missing physical feature
 206 data about their corresponding SMARP Active Region. Therefore, the dataset used for
 207 training has a Positive and a Negative component comprised of 65 and 6,510 flares re-
 208 spectively, making it vastly unbalanced.

209 The SMARP header files contains rows with the physical and spatial features of
 210 each active region at a 96-minute cadence throughout its entire lifetime, starting two days
 211 before it emerges or rotates onto the solar disk until two days after it submerges or dis-
 212 appears from view behind the limb (Bobra et al., 2021). We select the SMARP header
 213 file row at the time right before the flare peak time. In Figure 1, a histogram of the time
 214 difference between the selected SMARP file row and the flare peak time is plotted with
 215 the left panel corresponding to the positive dataset and right panel corresponding to the
 216 negative dataset. The importance of the SMARP time characterization varies between
 217 different SMARP active regions. Picking earlier SMARP points can potentially yield slightly

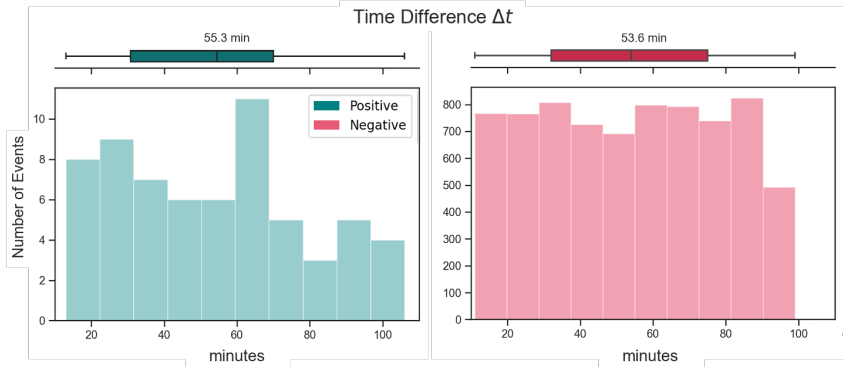


Figure 1. Histograms for the time difference between the flare peak time and the selected SMARP data (selected row in the TARP file). The distributions range between 10 and 100 minutes. The mean time differences shown in the error bars above the graphs are 55.3 and 53.6 minutes with a standard deviation of 28.6 and 24.8 minutes for the Positive (green) and Negative (red) datasets respectively.

different results than the ones presented below but the scope of the paper is to give an evaluation of the SMARP dataset rather than examine the optimal prediction window. The distributions range between 10 and 100 minutes. The mean time difference is 55.3 and 53.6 minutes with a standard deviation of 28.6 and 24.8 minutes for the Positive and Negative dataset respectively. The time difference Δt should not be confused with forecast windows or a lead times presented in similar works as the non-operational nature of our study does not allow for foreknowledge of whether a flare will occur.

2.3 Flare Predictors

We will evaluate the prediction power of SMARP dataset on SEP events by comparing the prediction results with those obtained by only using the flare information, i.e. flare peak intensity and flare location. We use the solar long wavelength X-ray flare data that NOAA's Geostationary Operational Environmental Satellites (GOES) continuously provides since 1975. Similarly to the SMARP Predictors, we calculate the flare angular distance from the earth's magnetic foot-point location, $W45^\circ$, on the sun.

3 Preliminary Data Analysis

We conduct preliminary analysis/assessment of different predictors, i.e. the SMARP and the Flare Predictors, via comparing the histogram of each predictor for the positive with that of negative samples. Figure 2 shows the density histograms of each predictor from the SMARP dataset on the top and from the GOES flare information on the bottom. The positive data is shown in green and negative data in red.

As shown in Figure 2, the flare peak intensity is a powerful discriminator between the positive and negative dataset. The flare peak intensity has been used as a feature to predict the occurrence and properties (peak proton intensity, event duration, and etc.) of SEP events (Laurenza et al., 2009; Balch, 2008). The four intervals present in the intensity graph ($[-7, -6]$, $[-6, -5]$, $[-5, -4]$ and $[-4, +\infty]$) represent the four different GOES X-ray classes B, C, M and X respectively. The predictive power difference between the flare peak intensity and the SMARP Predictors on the top of Figure 2 has a big impact when comparing the SEP prediction capability with and without SMARP data. Flare intensity is the predictor that has the least overlap between positive and negative. In the

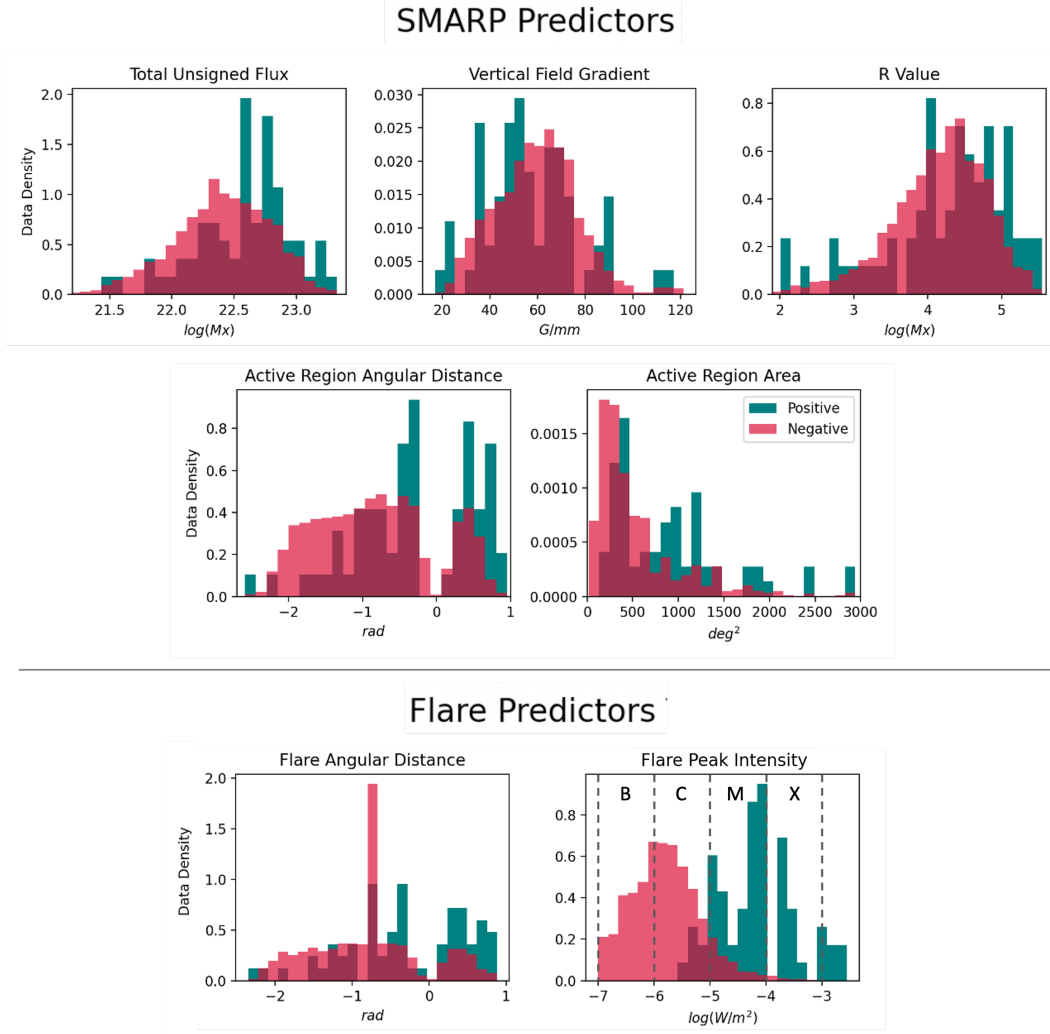


Figure 2. The probability density values are given in the histograms for the SMARP and Flare Predictors such that the area under each histogram integrates to 1. Both the SEP (green, positive samples) and flare-only (red, negative samples) data are split into 25 bins.

247 positive parameter spaces where there is no overlap, the associated SEP events are very
 248 intense and contain higher energy protons. Although larger flare intensity leads to higher
 249 proton fluxes, this is not necessarily the rule for all events as the problem is multidimen-
 250 sional (i.e. the correlation coefficient between the two physical quantities is not 1).

251 Moderate distinction between the events that led to an SEP and those that did not
 252 can be identified in the predictors acquired using the SMARP active region coordinates
 253 (Active Region Angular Distance and Area). Large active regions increase the likelihood
 254 of an SEP event occurrence. The total unsigned flux is related to the X-ray, the EUV
 255 emissions from the sun and the particle acceleration, which reflects the energy stored in
 256 an active region (Gurman et al., 1974; van Driel-Gesztelyi et al., 2003; Ugarte-Urra et
 257 al., 2015), therefore the SEP events are connected to higher flux values. The flux and
 258 intensity distributions show similar trends but with the former having less predictive power.
 259 The Vertical Field Gradient distribution of the Positive dataset aligns well with that of

260 the Negative dataset, making it the least powerful predictor along with the R Value which
 261 shows the same trend.

262 The predictive power differences observed in Figure 2 are consistent with a t-test
 263 (Kim, 2015) performed to all variables. For the most discriminating predictors, such as
 264 the Flare Peak Intensity and the AR Area, the Statistic values are high (21.3 and 6.35
 265 respectively) whereas for the weakest predictors (R Value and Vertical Field Gradient)
 266 the Statistic values do not exceed the 1.58 value.

267 4 Machine Learning Methods

268 To investigate whether the SOHO (SMARP AR data) or the GOES (flare erup-
 269 tion information) dataset can predict better the response variable of the two classes de-
 270 fined above, we use two popular groups of machine learning algorithms provided by the
 271 scikit-learn software package v0.24.2 for Python: different variations of the Support Vec-
 272 tor Machine (SVM; Cortes & Vapnik, 1995) and two Regression Models.

273 4.1 Support Vector Machine

274 SVMs were initially designed and have been used to solve binary classification prob-
 275 lems (Shao et al., 2014). In the most general case, the SVM is fitted to the data using
 276 a set of vector-target pairs (x_i, y_i) where $i = 1, 2, \dots, n$. The target for positive and neg-
 277 ative observations respectively is $y_i \in \{1, 0\}$ and the corresponding physical character-
 278 istics feature vector is $x_i = (f_{i1}, f_{i2}, \dots, f_{ip})$. For all tests performed, our training data
 279 length is $n = 116$ and the maximum feature vector length is $p = 7$, where all calcu-
 280 lated predictors are used. Each different SVM method maps the input feature vector x_i
 281 to a higher dimension space using an unknown function ϕ dependent on the user-defined
 282 kernel K . Given a regularization parameter $C > 0$ it solves an optimization problem
 283 to obtain the SVM trained weight vector w (Hsu & Lin, 2002; Inceoglu et al., 2018). The
 284 regularization parameter C controls the scaling of the SVM loss function and compen-
 285 sates for the change in the number of samples between the main problem and the smaller
 286 problems within the folds of the cross-validation. During testing, prediction is done by
 287 multiplying the trained vector w to the projected input feature vector $\phi(x_i)$ with an ad-
 288 dition of a bias term. A more detailed study on how to solve the SVM optimization equa-
 289 tions is out of the scope of this research and can be found elsewhere (Cortes & Vapnik,
 290 1995; Vapnik, 1998).

291 The kernel function K is defined as the inner product of data pairs that correspond
 292 to different observations i and j , $K(x_i, x_j) = \phi(x_i)^T \phi(x_j)$. In this study we train four
 293 different variations of the SVM (Amari & Wu, 1999). One uses the Linear kernel $K =$
 294 $\langle x_i, x_j \rangle$, two use Polynomial kernels $K = \gamma \langle x_i, x_j \rangle^d$ where $d \in \{2, 3\}$ (second and third
 295 degree) and the last one uses the Gaussian Radial Basis Function (RBF) kernel $K =$
 296 $\exp(-\gamma \|x_i - x_j\|^2)$ which has been used in similar studies (Inceoglu et al., 2018). To
 297 mitigate the risk of overfitting that is anticipated with the RBF kernel, the model pa-
 298 rameters are being chosen using the results that the ~~validation~~ testing dataset yields. Us-
 299 ing this cross-validation method we get similar results for the RBF kernel as we do for
 300 the rest models, which proves that excessive overfitting is being avoided. The weight-
 301 ing factor γ is user-defined and controls the influence a single training example has on
 302 the classification task. The different kernels help the prediction model deal with com-
 303 plex datasets such as the physical features of solar active regions by transforming the
 304 input into any desired form.

305 4.2 Linear Models

306 The observed physical properties of a SMARP AR can be also processed for the
 307 purpose of prediction by linear models: regression methods in which the target value is

308 expected to be a linear combination of the input features. Assuming a model function
 309 $f(x) = w^T x + b$ where w is a set of coefficients acquired during fitting, every feature's
 310 (x_i) predicted target y_i is 1 if $f(x_i) \geq 0$ and 0 if otherwise. In this case study linear
 311 models such as Ridge and Logistic regression are being used.

312 The ridge regression is one of the simplest machine learning algorithms and works
 313 well for small datasets while being computationally inexpensive. To fit the coefficients
 314 w to the training data, the ridge regression minimizes an ordinary Least Squares loss func-
 315 tion with an additional term that penalizes the size of the coefficients, as given in (1).

$$J_{Ridge} = \|w^T x - y\|_2^2 + \alpha \|w\|_2^2 \quad (1)$$

316 Between different training runs we vary the complexity parameter α in order to con-
 317 trol the amount of shrinkage and find the value that produces the most robust predic-
 318 tions. More specifically, α corresponds to $1/(2C)$ in the Logistic Regression or SVM and
 319 serves as a regularization parameter which improves the conditioning of the problem and
 320 reduces the variance of the estimates. Larger values specify stronger regularization. We
 321 do not adopt a cross-validation procedure for selecting the tuning parameter α due to
 322 considerations of sample sparsity and because the randomized picking process of the train-
 323 ing data leads to non-significant selection bias. Although it is a model often adopted when
 324 the response y takes real numbers, we chose ridge regression because it reduces overfit-
 325 ting, guarantees that we can find a solution and offers a different approach for binary
 326 classification compared to other competing models.

327 The dichotomous nature of Logistic Regression makes it a great candidate for the
 328 binary SEP prediction task. We use the default Logistic Regression module provided by
 329 the Scikit-Learn library in Python (Pedregosa et al., 2011) which includes the l_2 regu-
 330 larization as a penalty and the Limited-memory BroydenFletcherGoldfarbShanno (L-
 331 BFGS) optimization algorithm (Saputro & Widyaningsih, 2017) as a solver. The L-BFGS
 332 solver fits our application as it is robust and recommended for small dataset prediction
 333 tasks. To calculate the optimal w coefficients, Logistic Regression minimizes the cost func-
 334 tion J for w and c .

$$J_{LR} = \frac{1}{2} w^T w + C \sum_{i=1}^n \log(e^{-y_i(x_i^T w + c)} + 1) \quad (2)$$

335 The constant C controls the regularization strength of the model. Although nor-
 336 malization is applied to the flare data before the fitting process, the C constant is also
 337 varied throughout different training runs in order to find the value that produces the most
 338 numerically stable prediction.

339 4.3 Training and Tuning the ML Models

340 The scarcity of the SEP events along with the mission duration of the MDI/SOHO
 341 limits the size of the Positive dataset and leads to difficulties in separating the data into
 342 training and testing subsets in a reasonable way. To overcome this problem, every model
 343 is trained on 90% of the Positive events (58) and an equal number of Negative events.
 344 The training of each algorithm is followed by a similarly balanced testing on the remain-
 345 ing 10% Positive (7) and an equal number of Negative events. This balanced training
 346 and testing procedure is repeated k number of times to provide uncertainty assessment
 347 of the random selection of events. In our work, k is chosen to be equal to 100. Each time,
 348 a different batch of Negative events is randomly selected from the pool of 6,510 flare erup-
 349 tion events that did not lead to an SEP. This means that the Negative events are selected
 350 without replacement (in every one of the k different runs, all negative train and test sam-
 351 ples are different) while the Positive class is selected with replacement due to the low

352 number of available events. Similarly, in every run a different split between training and
353 testing occurs for the Positive dataset. It has to be highlighted that as the described method-
354 ology dictates, there is at no point any overlap between training and testing, neither in
355 the positive nor in the negative dataset. The training dataset is independent from the
356 testing therefore throughout the paper we only present the testing results. Furthermore,
357 the procedure that we describe above for assessing the uncertainty of the testing results
358 follows the idea of bootstrap (Efron, 1979; Efron & Tibshirani, 1994), which uses ran-
359 dom Monte Carlo samples to assign measures of accuracy, i.e. variance in our case. Across
360 the many ($k=100$) resampling of the training and testing sets, it is with probability al-
361 most 1 that we cover the best and worst cases in terms of testing accuracy and other met-
362 rics. Therefore, the uncertainty assessment that we give in this paper is crucial in un-
363 derstanding and interpreting our model performance.

364 It is important to acknowledge that the prediction of such rare events is a very dif-
365 ficult task, therefore our work is focused on using the SMARP dataset to underline the
366 contrast between the flares that lead to SEPs and flares that do not. Using the afore-
367 mentioned machine learning models we aim to understand the mechanism of SEP for-
368 mation and explore prediction capabilities rather than presenting an apparatus ready
369 for operation. More specifically, the results presented below aim to evaluate the poten-
370 tial the SMARP dataset has on predicting SEPs while showing which physical param-
371 eters are the most important for SEP prediction.

372 The very important limitations this work faces should also be noted. A very small
373 number of positive data along with the large uncertainties inherited by the SMARP dataset
374 are problems addressed by running the aforementioned Monte Carlo experiments (re-
375 peated random sampling) and by quantifying the uncertainties using box plots, statis-
376 tical spreads and standard deviations as seen in Section 5. Different solar cycles have dif-
377 ferent numbers of events and properties, therefore another limitation is that we are only
378 able to use data from solar cycle 23. Models that are ready for operational use and are
379 trained on highly unbalanced datasets need an increased weight to the penalty of the loss
380 function for the rare classes of data. In the models presented in this work such weight-
381 ing is not included, therefore it is expected that the model -if deployed as is- will have
382 an increased number of false alarms, as discussed in Section 5.4.

383 5 Results

384 For each one of the SVMs and Linear Models, we follow the same training proce-
385 dure, aiming to predict whether an AR that produces a flare will lead to an SEP event.
386 The goal is to illustrate how useful the SMARP dataset is for this particular task, we
387 therefore train the ML models using two separate sets of features, one that uses SMARP
388 information and one that uses flare information (see Section 2 for the detailed descrip-
389 tions of the two sets of features). The number of features vary from 2 to 5 and the ma-
390 chine learning algorithms are tested on a number of different predictor combinations.

391 The comparison between the different types of predictors and algorithms is done
392 using three metrics that characterize and quantify the predictive power of classifiers: the
393 Accuracy (ACC), the True Skill Statistics (TSS; Hanssen & Kuipers, 1965) and the Hei-
394 dke Skill Score (HSS; Heidke, 1926). The TSS and the HSS can take any values between
395 -1 (all incorrect) and 1 (all correct) while a value of 0 indicates a random forecast. Sim-
396 ilarly, the ACC can take values between 0 and 1 with 0.5 being the score of a random
397 forecast. More information about the metrics, their equations and statistical meaning
398 can be found in the works of Inceoglu et al. (2018) and Florios et al. (2018). For every
399 set of k different runs a cumulative contingency table (Figure 3, 4 and 5) is obtained based
400 on the results from the raw SVM and Linear Model outputs. Each row and column in
401 the contingency table represents the number of instances in an actual class and in a pre-

dicted class respectively. In an ideal case, the main diagonal of a contingency table gets high values while the rest of the matrix gets small values.

The ACC can be artificially high in the rare event where a model always predicts the majority class. In flare and SEP prediction such naive cases are common due to the data imbalance, but in this paper all models are trained on a one-to-one positive-negative ratio, so these rare cases are not a concern. Therefore, this study's basic prediction quality metric is the ACC, with the TSS and HSS being presented too as auxiliary metrics.

5.1 SEP Prediction with SMARP Predictors

The cumulative contingency table in Figure 3 shows that out of 1400 testing instances, Third Degree Polynomial SVM correctly classifies 552 as being Positive and 449 as being Negative when using the SMARP Flux and the AR Distance. This is the cumulative information obtained from 100 different runs, each of which has 14 testing points, therefore out of the 1400 total flares examined, there were 148 false alarms (flares wrongfully being classified as they will lead to an SEP) and 251 missed events (flares that lead to SEPs but the model predicted they do not). The ratio of all correctly classified positive and negative events over the total number of events is the accuracy (ACC). In this case, the mean accuracy suggests that 72% of the times ($\pm 12\%$ for a single run) the Third Degree Polynomial SVM algorithm can predict whether a flare will lead or not to an SEP using the its AR SMARP features. The comparison between the probability of detection and the probability of false detection lead to an average TSS level of 0.47 ± 0.24 for the same 100 runs. Similarly, the HSS measures a fractional forecast improvement over a random forecast of 0.44 ± 0.25 .

The results show that the Linear Models can predict whether a flare will be accompanied by SEPs with ACC values $\geq 0.70 \pm 0.12$ for a number of SMARP Predictor combinations (Table 2 in Appendix). The maximum corresponding TSS and HSS values for these SMARP Predictor runs are above the 0.40 levels. Similar to row 3 of Table 3, the Polynomial models that use the AR Distance and Area in Table 2 fail to produce a meaningful decision boundary yielding ACC values $\leq 0.52 \pm 0.04$, TSS values $\leq 0.15 \pm 0.28$ and HSS values $\leq 0.05 \pm 0.09$. Note here that a zero TSS or HSS value means that the method has no skill over the random forecast, therefore the these specific Polynomial examples do not show any predictive power at all.

Although the quality of the results cannot be judged based on the variance, the scores indicate that the better a model's predictive power is, the lower the variance between the different runs is. Thus, the flare peak intensity based models (Table 3, row 6-10) have an ACC standard deviation of ≤ 0.09 while in SMARP examples, where prediction quality is inferior, the ACC standard deviation is ≥ 0.10 . This pattern is even more evident when considering the TSS (or the HSS) for which the standard deviation can be as high as 0.28 at the SMARP Predictors exclusive runs in Table 2. Potential reasons about this behavior of variance is the small number of Positive data which allows for low quality runs to not converge at all.

Both SVM and Linear models are affected by user-defined constants such as the α and C in Equations 1 and 2. An embedded grid search -using the testing dataset to avoid overfitting- is employed for each experiment, where we vary each hyper-parameter on a range between 0.05 to 20. The parameter that produced the highest-quality and most consistent results was used for the examples presented in this study.

The negative dataset includes SMARP data only from the active regions that produced at least one flare throughout their lifespan. In a real-world application, the forecaster does not have prior knowledge of whether an AR is going to produce flares or not. In this work we chose to disregard quiet active regions (ARs that did not flare) not only because they are missing the flare predictors we compare SMARPs with, but mainly be-

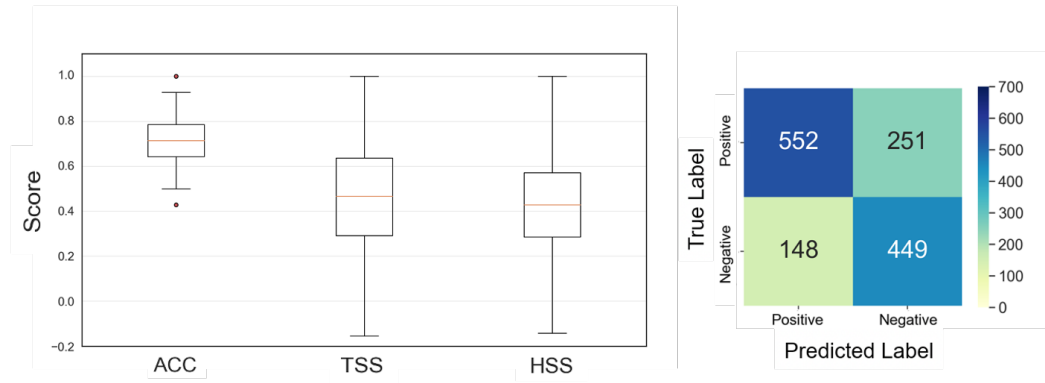


Figure 3. The distribution of $k = 100$ different ACC, TSS and HSS values are shown in the box plots (left). The values were obtained using the USFLUX & AR Distance on a Third Degree Polynomial SVM and constitute the best SEP prediction the SMARP data can achieve. The box range shows the interquartile range, the red line inside it the median value, the whiskers show the results range and the two red dots show two outlier values. The range of the y-axis is kept the same with Figure 5 for comparison. Adding all the individual TP, TN, FP and FN values respectively we produce a cumulative contingency table for the 100 different runs (right). For the runs where only SMARP predictors are used, the contingency table yields a False Alarm Ratio (FAR) of 0.211 and a Probability of Detection (POD) of 0.687.

452 cause they are very easily distinguishable compared to the ARs that produced SEPs. This
 453 claim is demonstrated in Figure 4, where the results of the binary classification between
 454 Positive SMARP events and SMARPs chosen randomly from quiet active regions are pre-
 455 sented.

456 The quiet dataset was created using the 3901 SMARPs that do not produce flares.
 457 Out of these 3901 ARs, almost half have missing data (especially RVALUE and USFLUX
 458 values that are zero). From the ones that have an adequate amount of data, we randomly
 459 select a SMARP data point, creating a quiet dataset of 1529 different predictor vectors.

460 In the histograms of Figure 4 we observe that the overlap between quiet and positive
 461 datasets for the USFLUX, ARAREA and RVALUE predictors is insignificant. The
 462 results under the histograms show that the classification task between positive and quiet,
 463 when using the USFLUX, the ARAREA and the RVALUE as predictors on a Polyno-
 464 mial SVM model, yield mean values of 0.96 ± 0.05 , 0.93 ± 0.09 and 0.91 ± 0.10 for the
 465 ACC, TSS and HSS respectively. These numbers very well demonstrate that distinguish-
 466 ing between Positive and Quiet active regions is a trivial problem for our machine learn-
 467 ing models, yielding classification accuracy errors of less than 4%.

468 5.2 SEP Prediction with Flare Predictors

469 Similar to the results presented before, the prediction quality metrics for the flare-
 470 only cases are calculated based on the contingency tables obtained from each different
 471 run. The cumulative contingency table in Figure 5 shows that out of 1400 testing instances,
 472 Ridge Regression correctly classifies 626 as being Positive and 651 as being Negative when
 473 trained on flare peak intensity and distance. This is the cumulative information obtained
 474 from $k = 100$ different runs, each of which was tested on 14 data points. The mean ac-
 475 curacy suggests that 91% of the times ($\pm 8\%$ for a single run) the Ridge algorithm can
 476 predict whether a flare will lead or not to an SEP using its physical characteristics. The
 477 comparison between the probability of detection and the probability of false detection

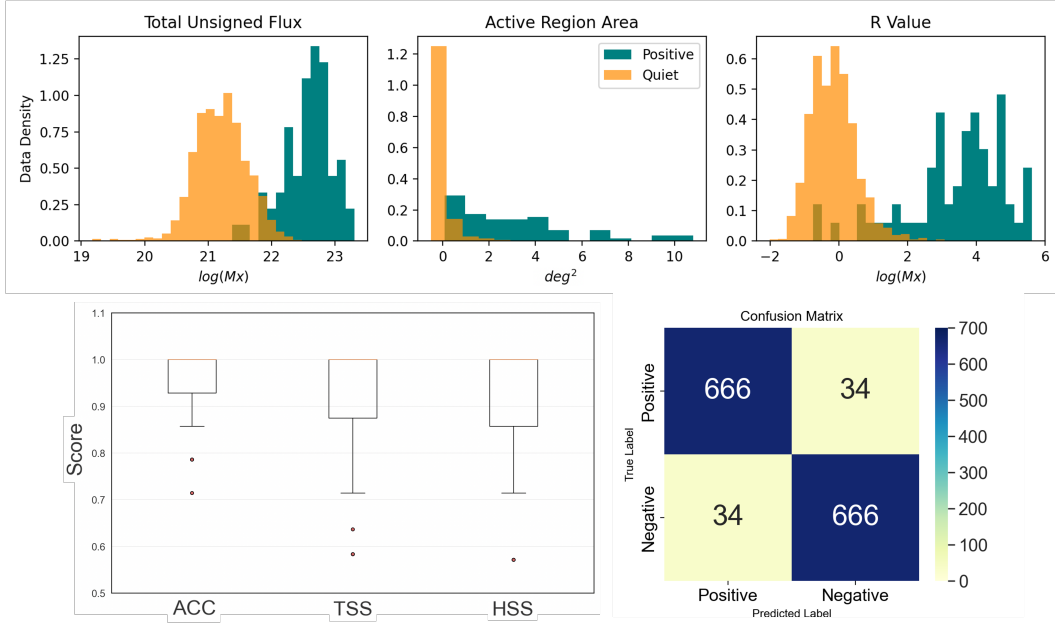


Figure 4. On the top panels of the figure, the probability density values are given in the histograms for three SMARP predictors (USFLUX, ARAREA and RVALUE) such that the area under each histogram integrates to 1. Both the SEP (green, positive ARs) and non-flaring (orange, quiet ARs) data are split into 25 bins. On the bottom part of the figure, the boxplots of the skill scores along with the cumulative contingency table for 100 different runs are presented. We observe that the median values for all three skill scores (ACC, TSS and HSS) are 1, which means no-error classification. For the runs that concern the quiet active regions, the contingency table yields a False Alarm Ratio (FAR) of 0.049 and a Probability of Detection (POD) of 0.951.

478 lead to an average TSS level of 0.84 ± 0.12 for the 100 runs. Similarly, the HSS mea-
 479 sures a fractional forecast improvement over a random forecast of 0.82 ± 0.14 .

480 Using explicitly the Flare Predictors (first row of Table 3) all six models produce
 481 similar results. The TSS and HSS show higher standard deviation values (varying from
 482 0.13 to 0.17) compared to the ACC. The predictive power of flare peak intensity is demon-
 483 strated when comparing the first two box plots in Figure 6 with the third and fourth one,
 484 where predictors (such as USFLUX, RVALUE, ARDIST, MEANGBZ etc.) other than
 485 Intensity are being used instead.

486 The ACC, TSS and HSS values range from 0.88 ± 0.09 to 0.92 ± 0.07 (values marked
 487 red and green in Table 3 of the Appendix), 0.78 ± 0.17 to 0.86 ± 0.13 and 0.76 ± 0.18 to
 488 0.84 ± 0.15 respectively for the runs that include Intensity accompanied with a SMARP
 489 Predictor (row 6-10). These results show that all models, when using the Flare Peak In-
 490 tensity, can successfully predict $\leq 92\%$ of the times if a flare will be accompanied with
 491 an SEP. When using the SMARP Predictors along with the Flare Distance (row 2-5) in-
 492 stead, the ACC values range from 0.60 ± 0.09 to a maximum of 0.71 ± 0.10 , the TSS
 493 from 0.36 ± 0.3 to 0.846 ± 0.2 and the HSS from 0.19 ± 0.18 to 0.42 ± 0.2 . This proves
 494 that when Intensity is not involved in the prediction process, all models yield inferior re-
 495 sults, losing at the best case 0.17 ± 0.09 from the accuracy metric. We only test our mod-
 496 els on a mix of SMARP and Flare Predictors to verify the prediction power of the In-
 497 tensity, as in real-life applications the two groups of predictors cannot be used together

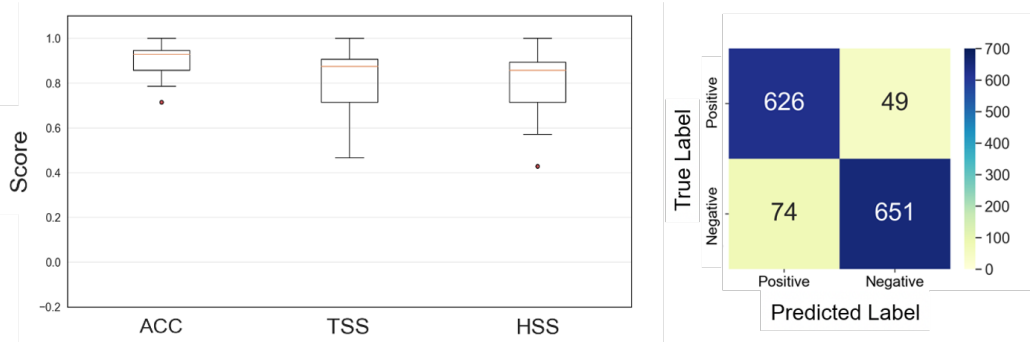


Figure 5. The distribution of $k=100$ different ACC, TSS and HSS values are shown in the box plots (left). The values were obtained using the Flare Peak Intensity & Flare Distance on a Ridge Regression model and constitute the best SEP prediction the flare data can achieve. The box range shows the interquartile range, the red line inside it the median value, the whiskers show the results range and the two red dots show two outlier values. The range of the y-axis is kept the same with Figure 3 for comparison. Adding all the individual TP, TN, FP and FN values respectively we produce a cumulative contingency table for the 100 different runs (right). For the runs where only flare predictors are used, the contingency table yields a False Alarm Ratio (FAR) of 0.110 and a Probability of Detection (POD) of 0.927.

498
499

due to the leading time difference (i.e. SMARP data is being acquired before the flare peak time).

500
501
502
503
504
505
506

Although each SVM or Linear model performs differently when trained on the same predictors, the variance between the models is of high significance only for some cases where the second and third degree Polynomial SVMs encounter convergence difficulties. The accuracy difference between the best and the worst performing models in Table 3 is always $\leq 0.03 \pm 0.13$, regardless the predictors combination, except for the extreme case of Flare Distance & ARAREA. This consistency gives us confidence to our results while suggesting that our models are not overfitting.

507
508
509
510
511
512
513

The maximum accuracy achieved on each one of the four main categories of predictor combinations is presented in Figure 6. The resulting ACC, TSS and HSS values show that regardless the machine learning model, the Flare Predictors generally perform better than the SMARP data because of the better predictive power of the flare peak intensity. Although the SMARP data cannot provide SEP forecast of quality similar to the flare peak intensity, it provides us with a larger leading time compared to the Flare Predictors as the flares precede in time the SMARP data points.

514

5.3 Comparison with Results in Literature

515
516
517
518
519
520
521
522
523
524

Inceoglu et al. (2018) used data provided by the SHARPs, GOES, and DONKI databases to train SVMs that forecast both CME and SEP events with maximum TSS and HSS of 0.92 ± 0.09 and 0.92 ± 0.08 . Anastasiadis et al. (2017) use the SDO/Helioseismic and Magnetic Imager (HMI) full-disk magnetograms and the flare information from the SOHO/MDI database on the prediction tool they call Forecasting Solar Particle Events and Flares (FORSPEF). They achieve Heidke Skill Scores (HSS) of 0.37 ± 0.011 and 0.67 ± 0.007 when using solar flare data and CME data respectively. While we only use GOES data to forecast exclusively SEP events (not CMEs), the best TSS and HSS our SVM implementations achieve are 0.84 ± 0.12 and 0.82 ± 0.14 , results that are comparable to both aforementioned studies.

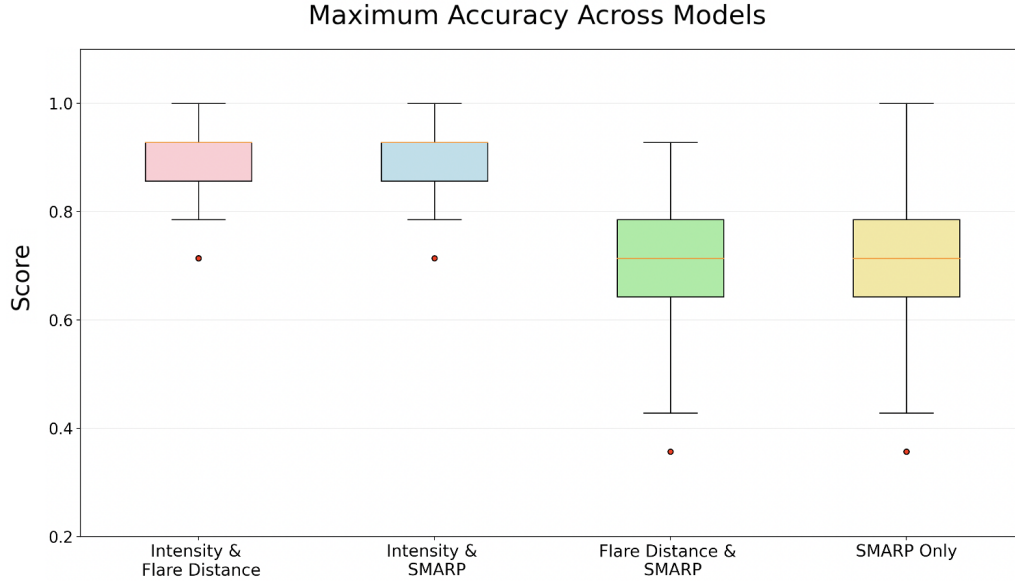


Figure 6. A cumulative box plot for the four main categories of predictor combinations outlined in the Appendix Tables. More specifically, the first plot (pink) corresponds to row 1 in Table 3, the second (blue) to rows 2-5, the third (green) to 6-10 and the fourth (yellow) corresponds to Table 2. The plot makes evident the superiority of the flare peak intensity over the SMARP data.

525 On the other hand, Papaioannou et al. (2018) perform a principal component analysis (PCA) on a set of six solar variables obtained from GOES and LASCO in order to
 526 calculate a decision boundary for their logistic regression. They classify events as SEP
 527 versus non-SEP and achieve a maximum POD (TSS + POFD) of 77.78%. Based on flare
 528 prediction, the warning tool García-Rigo et al. (2016) present provides long-term warn-
 529 ings of possible SEP event occurrence with POD scores of up to 58.3%. Núñez (2011)
 530 presents a dual-model system called UMASEP that has a POD of all (well and poorly
 531 connected with flares) SEP events of 80.72%. The SMARP data in Figure 3 achieves a
 532 POD 78.8%, similar to the works of Papaioannou and Nunez. If intensity gets involved
 533 in our logistic regression model, we can achieve POD scores of up to 90%.

535 All the results we report are using a probability threshold $p_t = 0.5$ because it is
 536 the one that is going to yield the optimal results for our models as demonstrated in Fig-
 537 ure 7. It is possible for works that use different models and datasets to slightly increase
 538 their prediction statistics as Anastasiadis et al. (2017) show in their work.

539 To demonstrate how the skill scores change when varying the decision threshold,
 540 we show how the combination of ARDIST and USFLUX (one of the best performing SMARP
 541 predictor combinations of Table 2) performs on the Logistic Regression model. On the
 542 left panel of Figure 7, a graph of the contingency table values is presented as the deci-
 543 sion threshold T_d changes from 0.3 to 0.7 in 0.05 increments. On the right one, the three
 544 different skill scores are being calculated and their box plots are graphed in order to present
 545 the uncertainties of each different decision threshold. Different models show similar trends
 546 with this specific case when the decision threshold is varied.

547 Lastly, it is important to note that none of the skill scores presented in this work
 548 can be immediately compared to the ones in the cited literature because the underly-
 549 ing class ratio and the train and testing data is not the same. It is therefore impossible

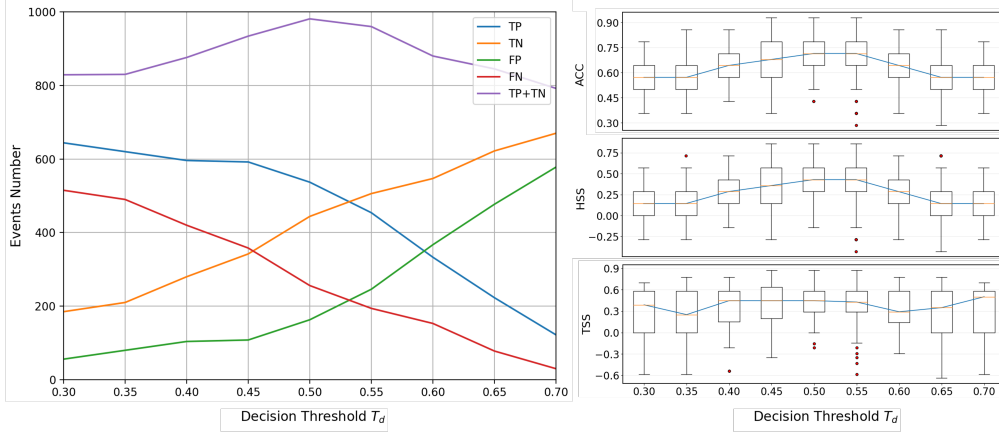


Figure 7. The left panel presents the contingency table values (TP, TN, FP and FN) along with the total number of correct predictions (TP+TN) for a number of different decision thresholds T_p . We observe that as the decision threshold increases, the number of false alarms (in green) increases while the number of missed events (in red) decreases. The total number of correct predictions is maximum where $T_d = 0.5$, therefore the results presented in the paper are calculated using this specific value. The right panel includes three different graphs for every one of the skill scores discussed in the Section 5 (ACC, TSS and HSS). The ACC and HSS mean values are at maximum when T_d is 0.50 and 0.55, while the TSS is mainly constant around 0.4 except a couple of decision threshold cases (0.35 and 0.6).

550 to do a fair comparison as SMARP is a newly published dataset and this is the first work
 551 that uses it for SEP prediction. Although we acknowledge that the comparisons with the
 552 literature should not be immediate, this is not a problem as the aim of this research is
 553 not to push the SEP prediction accuracy boundaries but rather to evaluate the SMARP
 554 dataset and its different predictors. This is the reason extra comparison runs using the
 555 GOES X-ray data are performed.

556 5.4 Connections to Operational Forecasting

557 Although the SMARP results presented above show that the dataset has some po-
 558 tential predictive capability, given its limitations (low resolution, limited positive events
 559 etc.), the SMARP data series could not be used alone to reliably forecast SEP events.
 560 Although the aim of the study is not to propose an operational forecasting apparatus,
 561 it is of interest to explore how the SMARP predictor models would perform if they were
 562 to be deployed. Three additional experiments (Experiment B, C and Operational in Ta-
 563 ble 5.4) are performed using the Support Vector Machine model of Figure 3 where the
 564 unsigned flux and the active region distance SMARPs are chosen to be the predictors.
 565 The differences between the experiments presented in Table 5.4 lie on the way the SMARP
 566 dataset is utilized during training and testing. The first experiment is the one discussed
 567 in Section 5.1 where a balanced number of positive and negative events is used during
 568 both training and testing. One way for handling imbalanced classes in SVMs is by in-
 569 troducing a hyperparameter C which determines the penalty for misclassification. In the
 570 last three experiments of Table 5.4, C is weighted in such a way (1:100 imbalance ratio)
 571 that it increases the penalty for misclassifying positive samples (the minority class) to
 572 prevent them from being overwhelmed by the negative ones (majority class).

573 In the second experiment (B), although during testing the same number (7) of posi-
 574 tive and negative samples is used (balanced testing), the SVM is trained using the en-

Experiment	Training	Ratio	Testing	Ratio	ACC	TSS	HSS	POD	FAR
Original	Balance	1:1	Balance	1:1	0.70 ± 0.12	0.43 ± 0.25	0.39 ± 0.23	0.78	0.31
B	Imbalance	1:100	Balance	1:1	0.69 ± 0.12	0.42 ± 0.27	0.39 ± 0.25	0.78	0.33
C	Imbalance	1:100	Imbalance	1:10	0.52 ± 0.04	0.01 ± 0.01	0.01 ± 0.01	0.80	0.98
Operational	Imbalance	1:100	Unknown	1:X	0.52 ± 0.05	0.01 ± 0.01	0.01 ± 0.01	0.81	0.98

Table 1. The four different experiments performed using the USFLUX and the AR Distance predictors on a Third Degree Polynomial Support Vector Machine model. The ACC, TSS, HSS results are given in the format of mean \pm standard deviation.

575 tire dataset mentioned in Section 2.2 comprised of 65 positive and 6,510 negative sam-
576 ples. In Experiments C and Operational, we further mimic a forecasting apparatus in
577 operational use in terms of having unbalanced samples: in Experiment C, the test set
578 is made up of 70 negative samples 7 positive samples; whereas in Experiment Operational,
579 the test set is randomly sampled from the full list of positive and negative samples. Again,
580 for all the experiments, we follow the idea of the bootstrap procedure as described in Sec-
581 tion 4.3 to assess the variability of the results, as given by the standard deviations of ACC,
582 TSS, and HSS in Table 5.4. The results in Table 5.4 show that although the SMARP
583 dataset has some ability to make distinctions between flares that produce SEPs and flares
584 that do not, when put in an operational setting where the testing is performed in a vastly
585 imbalanced set of samples, the forecasting model fails to produce meaningful results. It
586 is important to note that for the models which imitate an operational forecaster, the Prob-
587 ability of Detection (POD) is higher than the more experimental models, but the increase
588 in False Alarms results to low ACC, TSS and HSS scores. This is yet another proof of
589 the initial hypothesis that although the SMARP dataset includes meaningful informa-
590 tion which can be proven useful for SEP forecasting, it cannot be used by itself as a fore-
591 casting dataset.

592 6 Conclusions

593 To predict SEP events we use the newly published Space-Weather MDI Active Re-
594 gion Patches (SMARPs) dataset which includes observations of the solar magnetogram
595 that were made during the active Solar Cycle 23. Point data selected from the SMARP
596 time series is used on a variety of machine learning algorithms such as a different Sup-
597 port Vector Machines and Linear Regression models. The purpose of this study is to eval-
598 uate the power of this new data product for SEP forecast. Our results (Table 3 & 2) show
599 that SMARP can accomplish this task as it can identify correctly 72% of the times whether
600 an Active Region that produces a flare will lead to an SEP or not. Although the pre-
601 diction results for the SMARP dataset are worse than the ones produced using the flare
602 peak intensity and location, we demonstrate that not only SMARP data produces bet-
603 ter results compared to earlier SEP prediction works, but it also provides a better lead-
604 ing time than other datasets.

605 The task of SEP prediction using SMARP data is subject to inherent limitations
606 such as data uncertainties and a vastly unbalanced set of datasets which only includes
607 a limited amount of positive events. To overcome these difficulties, a Monte Carlo method
608 of random sampling, i.e. a bootstrap procedure, was employed to quantify of the results'
609 uncertainties. It is important to note that the results presented in this paper should not
610 be considered as an estimate of the accuracy that a prediction apparatus would yield if
611 deployed but should rather be viewed as an effort to quantify and compare the predic-
612 tion capability of the flare and SMARP predictors. In conclusion, although the SMARP
613 dataset is constructed from the MDI data set, which includes only the line-of-sight com-
614 ponent of the surface magnetic field at a relatively long 96-minute cadence, it can pro-

615 duce competitive prediction results for SEPs while providing a longer leading time than
616 using Flare Predictors.

617 **Acknowledgments**

618 This work was supported by NASA DRIVE Science Center grant 80NSSC20K0600.
619 The SEP event list we use in this work is documented in the NOAA Space Environment
620 Service Center website ([https://www.ngdc.noaa.gov/stp/satellite/goes/doc/SPE](https://www.ngdc.noaa.gov/stp/satellite/goes/doc/SPE.txt)
621 [.txt](https://www.ngdc.noaa.gov/stp/satellite/goes/doc/SPE.txt)), the SMARP data along with full-disk Helioseismic and Magnetic Imager (HMI)
622 magnetograms is available on the Joint Science Operations Center database at [http://](http://jsoc.stanford.edu/)
623 jsoc.stanford.edu/ and the NOAA solar X-ray flare dataset can be found at [https://](https://www.ngdc.noaa.gov/stp/solar/solarflares.html)
624 www.ngdc.noaa.gov/stp/solar/solarflares.html. All codes and data are included
625 in our Github repository at https://github.com/skasapis/SEP_Prediction_Using_SMAPR.

Table 2. Maximum ACC, TSS and HSS Values for the SVM and Linear Models using SMARP Predictors

SMARP Predictors	SVMs				Linear Models		Score
	Linear	RBF	Polynomial 2	Polynomial 3	Logistic Reg.	Ridge	
1. USFLUXL & ARDIST	0.67 ± 0.12	0.67 ± 0.13	0.70 ± 0.12	0.72 ± 0.12	0.70 ± 0.12	0.71 ± 0.12	ACC
	0.39 ± 0.28	0.38 ± 0.28	0.43 ± 0.25	0.47 ± 0.24	0.43 ± 0.23	0.47 ± 0.25	TSS
	0.34 ± 0.24	0.34 ± 0.25	0.39 ± 0.23	0.44 ± 0.25	0.40 ± 0.24	0.42 ± 0.24	HSS
2. USFLUXL & ARAREA	0.65 ± 0.11	0.67 ± 0.12	0.65 ± 0.12	0.67 ± 0.12	0.69 ± 0.11	0.65 ± 0.12	ACC
	0.35 ± 0.27	0.38 ± 0.28	0.34 ± 0.27	0.36 ± 0.26	0.37 ± 0.23	0.30 ± 0.27	TSS
	0.30 ± 0.23	0.35 ± 0.24	0.30 ± 0.24	0.33 ± 0.24	0.34 ± 0.22	0.27 ± 0.23	HSS
3. ARDIST & ARAREA	0.69 ± 0.11	0.69 ± 0.11	0.52 ± 0.04	0.51 ± 0.03	0.67 ± 0.12	0.70 ± 0.12	ACC
	0.42 ± 0.25	0.42 ± 0.23	0.15 ± 0.28	0.10 ± 0.22	0.36 ± 0.25	0.42 ± 0.26	TSS
	0.37 ± 0.23	0.38 ± 0.22	0.05 ± 0.09	0.03 ± 0.06	0.34 ± 0.24	0.40 ± 0.25	HSS
4. ARDIST & RVALUE	0.65 ± 0.13	0.68 ± 0.10	0.58 ± 0.11	0.60 ± 0.11	0.67 ± 0.11	0.66 ± 0.12	ACC
	0.33 ± 0.28	0.38 ± 0.22	0.18 ± 0.29	0.26 ± 0.29	0.36 ± 0.23	0.35 ± 0.25	TSS
	0.31 ± 0.26	0.35 ± 0.20	0.15 ± 0.22	0.21 ± 0.23	0.34 ± 0.21	0.33 ± 0.23	HSS
5. USFLUXL, ARDIST & ARAREA	0.67 ± 0.13	0.68 ± 0.11	0.70 ± 0.11	0.67 ± 0.10	0.70 ± 0.13	0.69 ± 0.10	ACC
	0.36 ± 0.28	0.38 ± 0.24	0.42 ± 0.23	0.37 ± 0.21	0.42 ± 0.27	0.41 ± 0.22	TSS
	0.34 ± 0.26	0.35 ± 0.22	0.39 ± 0.22	0.34 ± 0.19	0.39 ± 0.26	0.38 ± 0.21	HSS
6. All SMARP Predictors	0.68 ± 0.12	0.68 ± 0.13	0.66 ± 0.13	0.69 ± 0.12	0.67 ± 0.11	0.69 ± 0.13	ACC
	0.40 ± 0.24	0.32 ± 0.27	0.35 ± 0.28	0.42 ± 0.25	0.36 ± 0.24	0.40 ± 0.27	TSS
	0.37 ± 0.23	0.35 ± 0.25	0.33 ± 0.27	0.39 ± 0.24	0.33 ± 0.23	0.38 ± 0.26	HSS

Note. The ACC values ≥ 0.70 are marked in bolt. In green and red are marked the higher and lower accuracy values respectively.

626

Appendix

Table 3. Maximum ACC, TSS and HSS Values for the SVM and Linear Models using Different Predictors

Flare Predictors	Linear	SVMs			Linear Models		
		RBF	Polynomial 2	Polynomial 3	Logistic Reg.	Ridge	
1. Intensity & Flare Distance	0.90 ± 0.08	0.91 ± 0.07	0.90 ± 0.08	0.90 ± 0.08	0.90 ± 0.08	0.91 ± 0.07	ACC
	0.82 ± 0.16	0.84 ± 0.13	0.82 ± 0.15	0.80 ± 0.16	0.83 ± 0.16	0.84 ± 0.12	TSS
	0.80 ± 0.17	0.82 ± 0.14	0.80 ± 0.16	0.78 ± 0.17	0.80 ± 0.17	0.82 ± 0.14	HSS
SMARP & Flare Predictors							
2. Flare Distance & USFLUXL	0.71 ± 0.10	0.68 ± 0.13	0.70 ± 0.11	0.70 ± 0.11	0.71 ± 0.11	0.71 ± 0.12	ACC
	0.46 ± 0.20	0.39 ± 0.28	0.44 ± 0.24	0.42 ± 0.24	0.45 ± 0.22	0.46 ± 0.25	TSS
	0.42 ± 0.20	0.36 ± 0.26	0.40 ± 0.23	0.39 ± 0.23	0.42 ± 0.21	0.41 ± 0.23	HSS
3. Flare Distance & RVALUE	0.67 ± 0.14	0.69 ± 0.13	0.61 ± 0.14	0.61 ± 0.12	0.69 ± 0.12	0.70 ± 0.10	ACC
	0.36 ± 0.29	0.41 ± 0.27	0.25 ± 0.32	0.30 ± 0.32	0.41 ± 0.25	0.43 ± 0.22	TSS
	0.34 ± 0.27	0.38 ± 0.26	0.22 ± 0.27	0.23 ± 0.23	0.38 ± 0.24	0.40 ± 0.21	HSS
4. Flare Distance & ARAREA	0.66 ± 0.13	0.68 ± 0.11	0.60 ± 0.09	0.62 ± 0.10	0.70 ± 0.12	0.69 ± 0.11	ACC
	0.35 ± 0.27	0.40 ± 0.25	0.36 ± 0.30	0.39 ± 0.29	0.43 ± 0.24	0.42 ± 0.24	TSS
	0.32 ± 0.26	0.36 ± 0.22	0.19 ± 0.18	0.24 ± 0.21	0.41 ± 0.23	0.39 ± 0.22	HSS
5. Flare Distance, USFLUXL & ARAREA	0.69 ± 0.13	0.69 ± 0.13	0.67 ± 0.12	0.66 ± 0.11	0.69 ± 0.13	0.69 ± 0.11	ACC
	0.41 ± 0.27	0.42 ± 0.27	0.35 ± 0.30	0.36 ± 0.25	0.35 ± 0.26	0.41 ± 0.24	TSS
	0.38 ± 0.26	0.37 ± 0.26	0.32 ± 0.28	0.32 ± 0.23	0.39 ± 0.25	0.38 ± 0.22	HSS
6. Intensity & USFLUXL	0.88 ± 0.09	0.89 ± 0.07	0.90 ± 0.08	0.90 ± 0.08	0.88 ± 0.09	0.89 ± 0.09	ACC
	0.80 ± 0.16	0.80 ± 0.13	0.82 ± 0.14	0.80 ± 0.17	0.78 ± 0.17	0.79 ± 0.17	TSS
	0.77 ± 0.18	0.77 ± 0.15	0.80 ± 0.15	0.79 ± 0.18	0.76 ± 0.18	0.77 ± 0.17	HSS
7. Intensity & RVALUE	0.89 ± 0.08	0.91 ± 0.07	0.91 ± 0.08	0.89 ± 0.09	0.91 ± 0.07	0.90 ± 0.08	ACC
	0.79 ± 0.16	0.83 ± 0.14	0.83 ± 0.15	0.80 ± 0.16	0.84 ± 0.13	0.83 ± 0.15	TSS
	0.77 ± 0.17	0.82 ± 0.15	0.81 ± 0.16	0.78 ± 0.17	0.81 ± 0.14	0.80 ± 0.16	HSS
8. Intensity & ARDIST	0.91 ± 0.07	0.91 ± 0.07	0.91 ± 0.07	0.90 ± 0.07	0.92 ± 0.07	0.91 ± 0.07	ACC
	0.84 ± 0.13	0.84 ± 0.13	0.83 ± 0.14	0.83 ± 0.13	0.86 ± 0.13	0.84 ± 0.14	TSS
	0.82 ± 0.14	0.81 ± 0.14	0.81 ± 0.15	0.81 ± 0.14	0.84 ± 0.15	0.82 ± 0.15	HSS
9. Intensity, USFLUXL & ARDIST	0.91 ± 0.07	0.91 ± 0.08	0.90 ± 0.08	0.91 ± 0.07	0.92 ± 0.08	0.90 ± 0.08	ACC
	0.83 ± 0.14	0.84 ± 0.15	0.82 ± 0.15	0.83 ± 0.13	0.85 ± 0.15	0.81 ± 0.16	TSS
	0.82 ± 0.15	0.82 ± 0.16	0.80 ± 0.16	0.82 ± 0.14	0.83 ± 0.16	0.80 ± 0.17	HSS
10. Intensity & MEANGBL	0.90 ± 0.08	0.91 ± 0.07	0.91 ± 0.08	0.90 ± 0.08	0.91 ± 0.08	0.90 ± 0.08	ACC
	0.82 ± 0.15	0.85 ± 0.13	0.84 ± 0.15	0.82 ± 0.15	0.84 ± 0.14	0.82 ± 0.16	TSS
	0.80 ± 0.17	0.83 ± 0.14	0.83 ± 0.15	0.80 ± 0.16	0.82 ± 0.15	0.80 ± 0.16	HSS

Note. The ACC values ≥ 0.91 with standard deviation ≤ 0.07 are marked in bold. In green and red are marked the higher and lower accuracy values respectively for each one of the three predictor groups.

627 **References**

- 628 Amari, S.-i., & Wu, S. (1999). Improving support vector machine classifiers by mod-
 629 ifying kernel functions. *Neural Networks*, *12*(6), 783–789.
- 630 Anastasiadis, A., Papaioannou, A., Sandberg, I., Georgoulis, M., Tziotziou, K.,
 631 Kouloumvakos, A., & Jiggins, P. (2017). Predicting flares and solar energetic
 632 particle events: The forspet tool. *Solar Physics*, *292*(9), 1–21.
- 633 Bain, H., Brea, P., & Adamson, E. (2018). Using machine learning techniques to
 634 forecast solar energetic particles. In *Agu fall meeting abstracts* (Vol. 2018, pp.
 635 SM31D–3530).
- 636 Balch, C. C. (2008). Updated verification of the space weather prediction center’s
 637 solar energetic particle prediction model. *Space Weather*, *6*(1).
- 638 Barnes, W. T., Bobra, M. G., Christe, S. D., Freij, N., Hayes, L. A., Ireland, J., ...
 639 others (2020). The sunpy project: Open source development and status of the
 640 version 1.0 core package. *The Astrophysical Journal*, *890*(1), 68.
- 641 Bobra, M. G., Wright, P. J., Sun, X., & Turmon, M. J. (2021). Smarps and sharps:
 642 Two solar cycles of active region data.
- 643 Brueckner, G., Howard, R., Koomen, M., Korendyke, C., Michels, D., Moses, J., ...
 644 others (1995). The large angle spectroscopic coronagraph (lasco). In *The soho*
 645 *mission* (pp. 357–402). Springer.
- 646 Chen, Y., Manchester, W. B., Hero, A. O., Toth, G., DuFumier, B., Zhou, T., ...
 647 Gombosi, T. I. (2019). Identifying solar flare precursors using time series of
 648 sdo/hmi images and sharp parameters. *Space Weather*, *17*(10), 1404–1426.
- 649 Choudhary, D. P., Gosain, S., Gopalswamy, N., Manoharan, P., Chandra, R., Uddin,
 650 W., ... others (2013). Flux emergence, flux imbalance, magnetic free energy and
 651 solar flares. *Advances in Space Research*, *52*(8), 1561–1566.
- 652 Cortes, C., & Vapnik, V. (1995). Support-vector networks. *Machine learning*, *20*(3),
 653 273–297.
- 654 Efron, B. (1979). Bootstrap methods: Another look at the jackknife. *Annals of*
 655 *Statistics*, *7*(1), 1–26.
- 656 Efron, B., & Tibshirani, R. J. (1994). *An introduction to the bootstrap*. CRC press.
- 657 Florios, K., Kontogiannis, I., Park, S.-H., Guerra, J. A., Benvenuto, F., Bloomfield,
 658 D. S., & Georgoulis, M. K. (2018). Forecasting solar flares using magnetogram-
 659 based predictors and machine learning. *Solar Physics*, *293*(2), 1–42.
- 660 Forbush, S. E. (1946). Three unusual cosmic-ray increases possibly due to charged
 661 particles from the sun. *Physical Review*, *70*(9-10), 771.
- 662 García-Rigo, A., Núñez, M., Qahwaji, R., Ashamari, O., Jiggins, P., Pérez, G., ...
 663 Hilgers, A. (2016). *Prediction and warning system of sep events and solar flares*
 664 *for risk estimation in space launch operations*. EDP Sciences.
- 665 Gurman, J. B., Withbroe, G. L., & Harvey, J. W. (1974). A comparison of euv spec-
 666 troheliograms and photospheric magnetograms. *Solar Physics*, *34*(1), 105–111.
- 667 Hanssen, A., & Kuipers, W. (1965). *On the relationship between the frequency of*
 668 *rain and various meteorological parameters*. Koninklijk Nederlands Meteorologisch
 669 Instituut.
- 670 Heidke, P. (1926). Berechnung des erfolges und der güte der windstärkevorhersagen
 671 im sturmwarnungsdienst. *Geografiska Annaler*, *8*(4), 301–349.
- 672 Horne, R., Glauert, S., Meredith, N., Boscher, D., Maget, V., Heynderickx, D., &
 673 Pitchford, D. (2013). Space weather impacts on satellites and forecasting the
 674 earth’s electron radiation belts with spacecast. *Space Weather*, *11*(4), 169–186.
- 675 Hsu, C.-W., & Lin, C.-J. (2002). A comparison of methods for multiclass support
 676 vector machines. *IEEE transactions on Neural Networks*, *13*(2), 415–425.
- 677 Inceoglu, F., Jeppesen, J. H., Kongstad, P., Marcano, N. J. H., Jacobsen, R. H., &
 678 Karoff, C. (2018). Using machine learning methods to forecast if solar flares will
 679 be associated with cmes and seps. *The Astrophysical Journal*, *861*(2), 128.

- 680 Kim, T. K. (2015). T test as a parametric statistic. *Korean journal of anesthesiol-*
681 *ogy*, 68(6), 540.
- 682 Laurenza, M., Cliver, E., Hewitt, J., Storini, M., Ling, A., Balch, C., & Kaiser, M.
683 (2009). A technique for short-term warning of solar energetic particle events based
684 on flare location, flare size, and evidence of particle escape. *Space Weather*, 7(4).
- 685 Lavasa, E., Giannopoulos, G., Papaioannou, A., Anastasiadis, A., Daglis, I., Aran,
686 A., & Pacheco, D. (2021). Assessing the predictability of solar energetic particles
687 with the use of machine learning techniques.
- 688 Lopez, R. E., Baker, D. N., & Allen, J. (2004). Sun unleashes halloween storm. *Eos,*
689 *Transactions American Geophysical Union*, 85(11), 105–108.
- 690 McCracken, K., & Ness, N. (1966). The collimation of cosmic rays by the interplane-
691 tary magnetic field. *Journal of Geophysical Research*, 71(13), 3315–3318.
- 692 Mumford, S. J., Christe, S., Pérez-Suárez, D., Ireland, J., Shih, A. Y., Inglis, A. R.,
693 ... others (2015). Sunpypython for solar physics. *Computational Science &*
694 *Discovery*, 8(1), 014009.
- 695 Núñez, M. (2011). Predicting solar energetic proton events (e_i 10 mev). *Space*
696 *Weather*, 9(7).
- 697 Papaioannou, A., Anastasiadis, A., Kouloumvakos, A., Paassilta, M., Vainio, R.,
698 Valtonen, E., ... Abunin, A. (2018). Nowcasting solar energetic particle events
699 using principal component analysis. *Solar Physics*, 293(7), 1–23.
- 700 Papaioannou, A., Sandberg, I., Anastasiadis, A., Kouloumvakos, A., Georgoulis,
701 M. K., Tziotziou, K., ... Hilgers, A. (2016). Solar flares, coronal mass ejections
702 and solar energetic particle event characteristics. *Journal of Space Weather and*
703 *Space Climate*, 6, A42.
- 704 Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ...
705 others (2011). Scikit-learn: Machine learning in python. *the Journal of machine*
706 *Learning research*, 12, 2825–2830.
- 707 Reames, D. V. (2013). The two sources of solar energetic particles. *Space Science*
708 *Reviews*, 175(1-4), 53–92.
- 709 Richardson, I., Mays, M., & Thompson, B. (2018). Prediction of solar energetic
710 particle event peak proton intensity using a simple algorithm based on cme speed
711 and direction and observations of associated solar phenomena. *Space Weather*,
712 16(11), 1862–1881.
- 713 Rodriguez, J., Krosschell, J., & Green, J. (2014). Intercalibration of goes 8–15 solar
714 proton detectors. *Space Weather*, 12(1), 92–109.
- 715 Sadykov, V., Kosovichev, A., Kitiashvili, I., Oria, V., Nita, G. M., Illarionov, E., ...
716 Ali, A. (2021). Prediction of solar proton events with machine learning: Com-
717 parison with operational forecasts and” all-clear” perspectives. *arXiv preprint*
718 *arXiv:2107.03911*.
- 719 Saputro, D. R. S., & Widyaningsih, P. (2017). Limited memory broyden-fletcher-
720 goldfarb-shanno (l-bfgs) method for the parameter estimation on geographically
721 weighted ordinal logistic regression model (gwolr). In *Aip conference proceedings*
722 (Vol. 1868, p. 040009).
- 723 Scherrer, P. H., Bogart, R. S., Bush, R., Hoeksema, J.-A., Kosovichev, A., Schou, J.,
724 ... others (1995). The solar oscillations investigationmichelson doppler imager. In
725 *The soho mission* (pp. 129–188). Springer.
- 726 Schou, J., Scherrer, P. H., Bush, R. I., Wachter, R., Couvidat, S., Rabello-Soares,
727 M. C., ... others (2012). Design and ground calibration of the helioseismic and
728 magnetic imager (hmi) instrument on the solar dynamics observatory (sdo). *Solar*
729 *Physics*, 275(1), 229–259.
- 730 Schrijver, C., Beer, J., Baltensperger, U., Cliver, E., Güdel, M., Hudson, H., ... oth-
731 ers (2012). Estimating the frequency of extremely energetic solar events, based
732 on solar, stellar, lunar, and terrestrial records. *Journal of Geophysical Research:*
733 *Space Physics*, 117(A8).

- 734 Shao, Y.-H., Chen, W.-J., & Deng, N.-Y. (2014). Nonparallel hyperplane support
735 vector machine for binary classification problems. *Information Sciences*, *263*, 22–
736 35.
- 737 Shea, M., & Smart, D. (1995). History of solar proton event observations. *Nuclear*
738 *Physics B-Proceedings Supplements*, *39*(1), 16–25.
- 739 Ugarte-Urra, I., Upton, L., Warren, H. P., & Hathaway, D. H. (2015). Magnetic
740 flux transport and the long-term evolution of solar active regions. *The Astrophys-*
741 *ical Journal*, *815*(2), 90.
- 742 van Driel-Gesztelyi, L., Démoulin, P., Mandrini, C. H., Harra, L., & Klimchuk, J.
743 (2003). The long-term evolution of ar 7978: The scalings of the coronal plasma
744 parameters with the mean photospheric magnetic field. *The Astrophysical Jour-*
745 *nal*, *586*(1), 579.
- 746 van Driel-Gesztelyi, L., & Green, L. M. (2015). Evolution of active regions. *Living*
747 *Reviews in Solar Physics*, *12*(1), 1–98.
- 748 Vapnik, V. (1998). The support vector method of function estimation. In *Nonlinear*
749 *modeling* (pp. 55–85). Springer.
- 750 Webb, D. F., & Allen, J. H. (2004). Spacecraft and ground anomalies related to the
751 october-november 2003 solar activity. *Space Weather*, *2*(3).
- 752 Wild, J., Smerd, S., & Weiss, A. (1963). Solar bursts. *Annual Review of Astronomy*
753 *and Astrophysics*, *1*, 291.