



The University of Bradford Institutional Repository

<http://bradscholars.brad.ac.uk>

This work is made available online in accordance with publisher policies. Please refer to the repository record for this item and our Policy Document available from the repository home page for further information.

To see the final version of this work please visit the publisher's website. Access to the published online version may require a subscription.

Link to original published version: <http://dx.doi.org/10.1007/s11207-011-9896-1>

Citation: Ahmed OW, Qahwaji RSR, Colak T, Higgins PAB, Gallagher P and Bloomfield S (2013) Solar flare prediction using advanced feature extraction, machine learning and feature selection. *Solar Physics*. 283(1): 157-175.

Copyright statement: © 2013 Springer Verlag. Full-text reproduced in accordance with the publisher's self-archiving policy.

Editorial Manager(tm) for Solar Physics
Manuscript Draft

Manuscript Number:

Title: Solar Flare Prediction using Advanced Feature Extraction, Machine Learning, and Feature Selection

Article Type: TI: Image Processing in the Petabyte Era

Keywords: "Active Regions, Magnetic Fields"; "Flares, Forecasting"; "Photosphere"; "Space Weather"; "Feature Extraction"; "Machine learning"; "Feature Selection"

Corresponding Author: Omar Wahab Ahmed, BSc

Corresponding Author's Institution: University of Bradford

First Author: Omar Wahab Ahmed, BSc

Order of Authors: Omar Wahab Ahmed, BSc;Rami Qahwaji, PhD;Tufan Colak, PhD;Paul Higgins, A.B.;Peter Gallagher, PhD;Shaun Bloomfield, PhD

Abstract: Novel machine-learning and feature-selection algorithms have been developed to study: (i) the flare prediction capability of magnetic feature (MF) properties generated by the recently developed Solar Monitor Active Region Tracker (SMART); (ii) SMART's MF properties that are most significantly related to flare occurrence. Spatio-temporal association algorithms are developed to associate MFs with flares from April 1996 to December 2010 in order to differentiate flaring and non-flaring MFs and enable the application of machine learning and feature selection algorithms. A machine-learning algorithm is applied to the associated datasets to determine the flare prediction capability of all 21 SMART MF properties. The prediction performance is assessed using standard forecast verification measures and compared with the prediction measures of one of the industry's standard technologies for flare prediction that is also based on machine learning - Automated Solar Activity Prediction (ASAP). The comparison shows that the combination of SMART MFs with machine learning has the potential to achieve more accurate flare prediction than ASAP. Feature selection algorithms are then applied to determine the MF properties that are most related to flare occurrence. It is found that a reduced set of 6 MF properties can achieve a similar degree of prediction accuracy as the full set of 21 SMART MF properties.

Solar Flare Prediction using Advanced Feature Extraction, Machine Learning, and Feature Selection

Omar W. Ahmed¹; Rami Qahwaji¹; Tufan Colak¹; Paul A. Higgins²; Peter T. Gallagher²; D. Shaun Bloomfield²

Scholl of Computing Informatics and Media, University of Bradford, Bradford, UK¹;
Astrophysics Research Group, School of Physics, Trinity College Dublin, Dublin²,
Ireland

o.w.ahmed@bradford.ac.uk; r.s.r.qahwaji@bradford.ac.uk; t.colak@bradford.ac.uk;
pohuigin@gmail.com; peter.gallagher@tcd.ie; shaun.bloomfield@tcd.ie;

ABSTRACT:

Novel machine-learning and feature-selection algorithms have been developed to study: (i) the flare prediction capability of magnetic feature (MF) properties generated by the recently developed Solar Monitor Active Region Tracker (SMART); (ii) SMART's MF properties that are most significantly related to flare occurrence. Spatio-temporal association algorithms are developed to associate MFs with flares from April 1996 to December 2010 in order to differentiate flaring and non-flaring MFs and enable the application of machine learning and feature selection algorithms. A machine-learning algorithm is applied to the associated datasets to determine the flare prediction capability of all 21 SMART MF properties. The prediction performance is assessed using standard forecast verification measures and compared with the prediction measures of one of the industry's standard technologies for flare prediction that is also based on machine learning – Automated Solar Activity Prediction (ASAP). The comparison shows that the combination of SMART MFs with machine learning has the potential to achieve more accurate flare prediction than ASAP. Feature selection algorithms are then applied to determine the MF properties that are most related to flare occurrence. It is found that a reduced set of 6 MF properties can achieve a similar degree of prediction accuracy as the full set of 21 SMART MF properties.

Keywords: *Active Regions, Magnetic Fields; Flares, Forecasting; Photosphere; Space Weather; Feature Extraction; Machine learning; Feature Selection;*

1. Introduction

Solar flares can have catastrophic effects on our infrastructure by degrading the Global Positioning System (GPS), interrupting power grids, and causing failures in communications satellites. Roughly half of Coronal Mass Ejections (CMEs) are associated with flares (Zhang et al., 2001), which produce magnetic storms and distort our ionosphere as they impact upon the Earth (Gopalswamy et al., 2005). This distortion renders sensitive GPS measurements highly inaccurate. Commercial airplanes rely on GPS to take off, navigate and land. Currents produced in the ionosphere by intense space weather events may generate huge currents in power grids, terminally damaging the

1 massive transformers that are integral to these systems. Finally, ionising particle radiation
2 produced by flares and CMEs may damage or even result in the loss of communications
3 satellites, as was the case with the Galaxy 15 satellite in April 2010¹. The accurate
4 prediction of solar flare occurrence is essential for operations teams to safely perform
5 their respective jobs in anticipation of damaging space weather (Committee on the Social
6 and Economic Impacts of Severe Space Weather Events, 2008): power grid operators
7 need to know when to expect ionospheric currents; pilots need to know when to divert
8 transpolar flights to lower latitudes; satellite operators need to know when to turn off
9 equipment; astronauts need to know when to seek cover in shielded areas.

10
11 To date, various systems and models designed to predict the occurrence of solar flares
12 have made significant progress, but the achieved prediction performances are far from
13 what is required by operations teams. Therefore, further investigations are needed to
14 enhance both the understanding of the physical causes of flares and the design of a more
15 accurate flare prediction system. Three main categories of prediction models exist –
16 expert-based (with human input), linear statistical, and non-linear statistical (including
17 machine learning). Recent prediction systems relying on non-linear methods, such as
18 Artificial Neural Networks (ANN), show the most promise (Messerotti et al., 2009).

19
20 There have been a number of attempts and proposed approaches to create an accurate
21 flare prediction system. One of the earliest systems is THEO (McIntosh 1990), an expert
22 system using subjective judgements and statistical correlations, that was adopted in 1987
23 by the Space Environment Center (SEC) at the National Oceanic and Atmospheric
24 Administration (NOAA). This system utilises a number of sunspot and magnetic field
25 properties to generate a prediction for the occurrence of various solar flare classes.
26 Gallagher et al. (2002) produced a linear prediction system that was adopted by
27 SolarMonitor, using the average flare rate for each human observed McIntosh sunspot
28 classification and Poisson statistics to calculate the flare probability for individual
29 classifications. Later systems used aspects of each of these, such as determining multiple
30 characteristics of sunspot groups and active regions, and using both linear statistical and
31 non-linear prediction methods.

32
33 Linear statistical studies have been performed by several authors that attempt to make
34 flare predictions by identifying the active region magnetic properties that are most
35 correlated with flare activity. Using line-of-sight magnetograms, Cui et al. (2006)
36 investigated active region maximum horizontal gradient, the length of the neutral lines,
37

38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
¹ See: Report of the NOAA Tiger Team (retrieved 15 March 2011):
http://ngdc.noaa.gov/stp/satellite/anomaly/2010_sctc/docs/1-2_WDenig.pdf

1 and the number of singular points using sigmoid analysis. They found that although there
2 were high correlations with flaring, these properties did not accurately predict flares. Jing
3 et al. (2006) found a positive correlation with flaring for the mean value of spatial
4 magnetic gradients along strong-gradient magnetic neutral lines, the length of strong-
5 gradient magnetic neutral lines, and the total magnetic energy. Leka and Barnes (2007)
6 calculated many properties from vector magnetograms and applied linear discriminant
7 analysis. They found that total magnetic flux or the combination of total vertical currents
8 with measures of the magnetic shear were best for predicting C-class flares and above,
9 while excess photospheric magnetic energy was best for M-class flares and above. Barnes
10 and Leka (2008) also use discriminant analysis to investigate total flux, total excess
11 energy, a measure of the amount of magnetic flux close to high gradient polarity-
12 separation lines, and the effective connected magnetic fields. They found that by using a
13 discriminant boundary none of the investigated properties were able to predict major
14 flares (i.e., M- or X- class) significantly better than always predicting that no flare would
15 occur. Song et al. (2008) use ordinal logistic regression with measures of total flux,
16 strong-gradient neutral line length, and magnetic energy dissipation (overall gradient
17 measure). Mason and Hoeksema (2010) use superposed epoch analysis with total
18 magnetic flux, primary inversion line length, effective separation, and gradient-weighted
19 inversion-line length.
20
21
22
23
24
25
26
27
28
29
30

31
32 Other studies take advantage of complex non-linear learning algorithms to train decision-
33 making systems using large samples of characterised sunspot group and active region
34 observations. Colak and Qahwaji (2009) implemented an automated near real-time hybrid
35 system, based on machine learning, called Automated Solar Activity Prediction (ASAP),
36 using measurements of sunspot area and automated McIntosh classifications. Yu et al.
37 (2009, 2010a, 2010b) use machine learning on neutral line properties determined from
38 magnetograms. Yuan et al. (2010) combine the methods in Song et al. (2008) with
39 machine learning.
40
41
42
43
44
45

46
47 To rigorously evaluate the performances of prediction systems and the physical properties
48 utilised by them, standard forecast verification measures such as the Heidke Skill Score
49 (HSS; Balch, 2008) must be adopted (Barnes and Leka, 2008). Some of the
50 abovementioned systems were validated using large data sets and report accurate
51 prediction results, but few if any have been tested by operationally predicting solar flares.
52 Running a system operationally corresponds to the way it would be used in a real-time
53 setting, implying that all features detected inside of some observational bounds are given
54 a prediction. While validation must be done operationally, the data sets used to train a
55 prediction system may be segmented using some selection criteria. A portion of the total
56
57
58
59
60
61
62
63
64
65

1 data set is removed with the intention of helping a system to discriminate between flaring
2 and non-flaring feature populations more clearly. In addition to using a segmented
3 training set, some authors also segment the testing set by applying selection criteria to the
4 data that are used for system validation. This results in prediction performances that
5 reflect the efficiency of training rather than how the system would perform in a truly
6 operational sense. We wish to emphasise the importance of determining skill scores by
7 performing validation in a manner as close as possible to how the system will actually be
8 used.
9

10
11
12
13
14 In this paper we apply machine learning and feature selection algorithms to a set of
15 magnetic feature (MF) properties to determine: (i) their overall flare prediction capability;
16 (ii) the properties that are most significantly related to flare occurrence. In this work we
17 also aim to improve on previous work in several important ways. Here we explore the
18 difference between operational and segmented validation for the first time. The flare
19 prediction system is tested against data in a segmented training format (i.e., defining its
20 training benchmark) as well as being tested against non-segmented data (i.e., defining its
21 operational prediction performance). In addition to realistic validation, magnetic features
22 are identified and extracted consistently using automated feature recognition to avoid any
23 selection bias, while previous studies have used NOAA visually identified features.
24
25
26
27
28
29
30

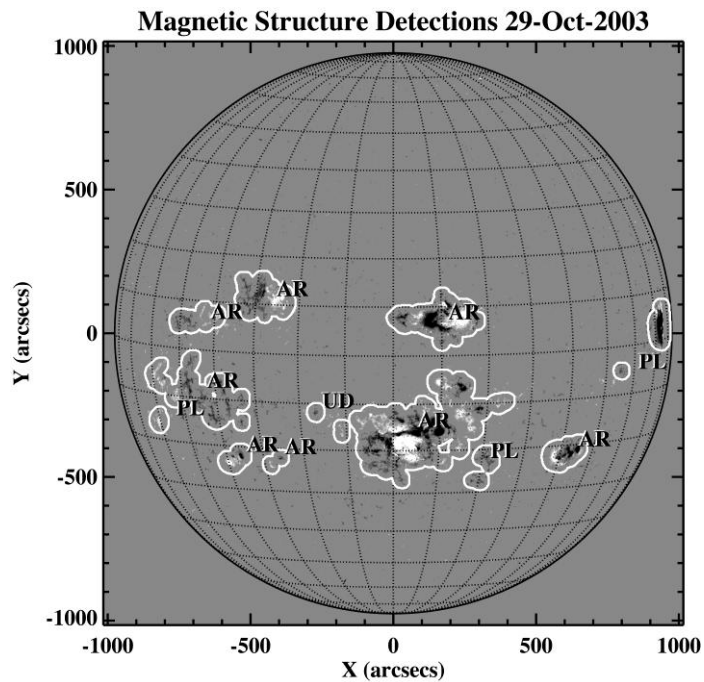
31
32 This paper is organised as follows. The data sources and their specifications are discussed
33 in Section 2. The methods are explained in Section 3, including the MF-flare association
34 algorithms (Section 3.1), machine learning (Section 3.2), and feature selection (Section
35 3.3). The results are presented in Section 4, including the prediction capability of the MF
36 properties studied (Section 4.1) and the MF properties that are most related to flare
37 occurrence (Section 4.2). Finally, some discussion and ideas for future work are presented
38 in Section 5.
39
40
41
42
43
44

45 **2. Data**

46
47 Solar flares are the impulsive release of large amounts of energy (up to $\sim 10^{27}$ J) in the
48 form of energetic particles and emission across the entire electromagnetic spectrum. The
49 common format for classifying these events uses the peak magnitude of soft X-ray flux as
50 observed by the Geostationary Operational Environmental Satellite (GOES) series.
51 Catalogues of flare events recorded by these satellites were obtained from the National
52 Geophysical Data Center (NGDC), which holds one of the most comprehensive public
53 databases for solar features and activity recorded by multiple observatories around the
54 world. Only those flare events with peak GOES magnitudes above the C1.0 level (i.e.,
55 10^{-6} W/m²) with known locations were included in this study.
56
57
58
59
60
61
62
63
64
65

1 In this paper, the flare prediction potential of MF properties generated by the Solar
 2 Monitor Active Region Tracker (SMART; Higgins et al. 2010) are evaluated for the first
 3 time. SMART is a recently developed feature extraction algorithm to detect, characterise,
 4 and catalogue MFs using 96-min *SoHO*/MDI line-of-sight magnetograms. MFs are
 5 detected in magnetograms by segmenting quiet-Sun and feature pixels using a
 6 combination of image processing techniques. SMART detects MFs automatically and is
 7 completely independent from NOAA active regions. Throughout this paper the term “MF
 8 detection” refers to an individual SMART MF detected in one MDI magnetogram (i.e., a
 9 single MF will be observed multiple times through its lifetime). An example of a set of
 10 SMART detections is shown in Figure 1. Of the various magnetic field properties
 11 determined by SMART, 21 are utilised in this paper and these are described in Table 1.

12
 13
 14
 15
 16
 17
 18
 19 The complete time range considered here for MF detections and flare events extends from
 20 April 1996 to December 2010. The data in this period have been investigated in several
 21 ways, according to the aim of each experiment. More details about the number of MF
 22 detections and flare events in these catalogues are given in Section 3.1.



23
 24
 25
 26
 27
 28
 29
 30
 31
 32
 33
 34
 35
 36
 37
 38
 39
 40
 41
 42
 43
 44
 45
 46
 47
 48
 49
 50
 51 Figure 1: An example of all SMART MF detections on 29 October 2003. AR (active region)
 52 denotes features classified as multipolar, while PL (plage) and UD (unipolar decaying) denote two
 53 different classes of unipolar feature.

Table 1: SMART MF properties.

Property ID	Property	Description
v1	Type-Polarity	AR polarity (Unipolar / Multipolar)
v2	Type-Size	AR size (Big / Small)
v3	Type-Evolution	AR evolution (Emerging / Decaying)
v4	A	Area of the region
v5	Φ	Total unsigned magnetic flux of the region
v6	Φ_+	Total positive flux in the region
v7	Φ_-	Total negative flux in the region
v8	Φ_{IMB}	Flux imbalance fraction in the region
v9	$\Delta\Phi/\Delta t$	Flux emergence rate
v10	B_{MIN}	Minimum magnetic field in the region
v11	B_{MAX}	Maximum magnetic field in the region
v12	B_{MEAN}	Mean magnetic field in the region
v13	L_{NL}	Neutral line length in the region
v14	L_{SG}	High gradient neutral line length in the region
v15	∇_{MAX}	Maximum gradient along the neutral line
v16	∇_{MEAN}	Mean gradient along the neutral line
v17	∇_{MEDIAN}	Median gradient along the neutral line
v18	R	Schrijver R value
v19	WL_{SG}	Falconer WL_{SG} value
v20	R^*	Schrijver R value with a lower threshold
v21	WL_{SG}^*	A modified version of WL_{SG}

3. Methods

To enable the investigation of SMART's MF detections in relation to flares, there is a need to establish the flaring and non-flaring MF detections. Two types of association algorithms have been adopted for this purpose (Section 3.1). The experiments conducted in this work aims to achieve two goals: (i) determine the flare prediction capability of SMART's MF properties (Section 3.2) and (ii) determine the MF properties that are most related to flaring (Section 3.3). Data preparation and the methods applied in this work are discussed in this section.

3.1. MF AND FLARE ASSOCIATION

Algorithms have been developed to associate SMART MF detections with flares from the NGDC catalogue for the complete time period under consideration (i.e., April 1996–

December 2010). The purpose of this association is to classify MF detections as *flaring* or *non-flaring*. NGDC-listed flares may already be associated with NOAA active regions, but a new association algorithm is necessary for this work because SMART's MF detections are independent of NOAA active regions – i.e., some SMART MF detections correspond to NOAA-numbered spotted regions, while the rest correspond to unspotted magnetic flux regions. MFs and flares are associated based on their location and the time difference between them. Flare locations provided in the NGDC catalogue are remapped to the times of MF detections using the method described in (Colak and Qahwaji, 2010) and the location compared to the MF spatial coverage. MF detections are then defined as flaring if a remapped flare location falls within the boundary of the SMART MF contour. In order to minimise the error in magnetic field properties caused by projection effects, only MF detections located within 45° from solar disk centre are considered for this work.

We consider two forms of association – *segmented* and *operational* – that are distinguished by differing criteria for a non-flaring MF detection. In the segmented form, a MF detection is classified as flaring if it produced at least one C-, M-, or X- class flare in the following 24-hour period, and non-flaring if did not cause any C-, M-, or X- class flares in the +/-48-hour period around its observation time. Coupled with the flaring definition, this means that some MF detections are discarded from the complete set by the segmented association algorithm (i.e., MF detections observed 24–48 hours prior to a flare are neither classified as flaring nor non-flaring). In effect, the data are segmented into MF detections that were observed very close in time before a flare and those that are at least two days removed from a flare.

In the operational form, MF detections are classified as flaring in exactly the same way that was used to classify flaring in the segmented association. However, the operational definition of non-flaring is that a MF did not produce any C-, M-, or X-class flares in the 24-hour period following its observation. This satisfies the primary requirement of a real-time operational prediction system, such that each MF detection must be given a prediction and so requires classification as either flaring or non-flaring. A summary of the input and output data in the association processes are shown in Tables 2 and 3.

The choice of which association set and data time range to use in the later sections of this paper is decided according to the experimental aim. To determine the flare prediction capability of SMART's MF properties, two experiments were considered. In the first experiment, we use both the segmented and the operational sets covering the entire period of April 1996–December 2010 (Sections 3.2.1 and 4.1.1). The number of flaring/non-flaring MF detections used in this experiment is detailed in Table 3. The machine-

learning algorithm is trained and tested on data chosen randomly from that time period. In the second experiment, the data covering April 1996–December 2000 and January 2003–December 2008 are used to train the machine-learning algorithm and the data covering January 2001–December 2002 and January 2009–December 2010 are used to test the system (Sections 3.2.2 and 4.1.2). Both segmented and operational data sets have been experimented with. Table 4 details the number of flaring/non-flaring MF detections used in this experiment. The time coverage of the training set was chosen so that the remaining testing set would contain MF detections and flare activity from periods around the maximum and minimum levels of solar activity.

To determine the MF properties that are most related to flares, we use data covering the entire time range of April 1996–December 2010 (Sections 3.3 and 4.2). The segmented set is used to determine the MF properties that are most related to flares, while both segmented and operational sets are used to determine the prediction capability of the selected MF properties. Once again, the number of flaring/non-flaring MF detections used in this experiment is detailed in Table 3.

Table 2: Number of flares and MF detections within 45° of disk centre input to the association algorithms, covering April 1996–December 2010.

MF Detections	Flare Events		
	C	M	X
521,578	7,319	1,072	107

Table 3: Number of flaring and non-flaring MF detections in the segmented and operational sets covering April 1996–December 2010.

Association Method	Flaring MF Detections	Non-Flaring MF Detections	Total MF Detections
Segmented	27,539	469,516	497,055
Operational	27,539	494,039	521,578

Table 4: Number of flaring and non-flaring MF detections in time independant training and testing sets, from each association output (data sets from segmented and operational association).

Association Method	Training Set (Apr1996–Dec2000, Jan2003–Dec2008)			Testing Set (Jan2001–Dec2002, Jan2009–Dec2010)		
	Flaring MF Detections	Non-Flaring MF Detections	Total MF Detections	Flaring MF Detections	Non-Flaring MF Detections	Total MF Detections
Segmented	16,864	300,306	317,170	10,675	169,210	179,885
Operational	16,864	315,561	332,425	10,675	178,478	189,153

3.2. MACHINE LEARNING

In this section, the flare prediction capability is investigated by applying a Cascade Correlation Neural Network (CCNN) machine-learning algorithm to the associated data sets. CCNN is a learning algorithm that is proven to provide efficient performance in applications that involve classification and time-series prediction (Frank et al. 2001). It has been shown in (Qahwaji and Colak, 2006) that CCNN is the optimal neural network learning algorithm for solar flare prediction using sunspot properties. A detailed description of CCNN and its application in flare prediction can be found in Qahwaji and Colak (2006, 2007).

The CCNN algorithm that is available in Matlab has been utilised to implement the experiments described in this section. The CCNN used here consists of several layers – an input layer, multiple hidden layers, and an output layer. The number of MF properties that are fed into the machine learning system determines the number of nodes in the input layer. The numbers of nodes in the hidden layers are determined automatically during the training, while the number of classes determines the number of nodes in the output layer (i.e., 1 node for flare/no-flare). It is essential to provide the machine learning with uniform data to enhance its learning and performance. Therefore, the input data has been normalised so that the measurements of each MF property are represented in the range 0.1–0.9 and the output classes are represented as 0.9 for flare and 0.1 for no-flare.

For each association dataset (segmented and operational), machine learning is applied twice. In the first instance, machine learning is applied using cross-validation in order to determine the overall prediction capability of the investigated dataset (Sections 3.2.1 and 4.1.1). In the second instance, the machine learning is applied using time-separated training and testing sets in order to determine the system's prediction capability on data from particular periods of time (Sections 3.2.2 and 4.1.2). Further details about the application of machine learning using each method are described below.

3.2.1. Machine Learning using Cross Validation

Cross-validation is a method that partitions the input data into subsets so that the learning algorithm can be trained on a subset and internally tested on a different subset. Cross-validation is a useful approach for analysing the prediction performance of machine learning, as it is important to avoid over-fitting. Over-fitting occurs when the learning algorithm performs very well on the training data, but not so well when provided with new data. Different forms of cross-validation method exist and repeated random sub-sampling validation is applied here. This method is based on randomly dividing the data into a number of subsets, which is repeated a number of times so that the learning

1 algorithm is trained and tested on different data. For each repetition, one subset is used
2 for training and the rest are used to evaluate the prediction performance by calculating a
3 number of forecast verification metrics. These measurements are then averaged in order
4 to provide an indication of the effectiveness of the machine learning on the training data
5 (Hall, 1999).
6
7

8
9 Cross-validation is applied separately to both the operational and segmented data sets for
10 the entire period covering April 1996- December 2010. For each investigated set, the data
11 are randomised and two separate portions of data are created: a training portion (60%)
12 and a testing portion (40%). The MF properties and their corresponding flare/no-flare
13 classifications from the training portion are fed into the learning algorithm for training
14 purposes. When the training process is completed, the learning algorithm is fed with the
15 MF properties from the testing portion. The learning algorithm attempts to predict their
16 flare/no-flare classifications, producing values in the range 0.1–0.9. A threshold value of
17 0.5 is used to categorise the generated prediction outputs as either flare (>0.5) or no-flare
18 (<0.5). These predicted outputs are compared with the testing portion's actual
19 classifications using standard forecast verification measures to evaluate the prediction
20 performance of the learning algorithm, such as True Positive Rate (TPR), False Positive
21 Rate (FPR), True Negative Rate (TNR), False Negative Rate (FNR), False Alarm Rate
22 (FAR), Mean Squared Error (MSE), Accuracy (ACC), and Heidke Skill Score (HSS).
23 Detailed information about each of these measures can be obtained from Fawcett (2006)
24 and Balch (2008). Among the prediction measures, HSS is one of the best indicators of
25 the overall performance of a prediction method since it accounts for correct chance
26 forecasts (Barnes and Leka, 2008). The cross-validation process is repeated 10 times and
27 the means of the prediction measures are calculated.
28
29
30
31
32
33
34
35
36
37
38
39
40
41

42 *3.2.2. Machine Learning with Time-separated Training and Testing Sets*

43 Machine learning has been applied by training and testing the system using data from
44 different time periods, with both the operational and segmented data sets investigated.
45 The data covering April 1996-December 2010 is divided into two sets for this purpose.
46 The training set includes the majority of the time coverage with the exclusion of two 2-
47 year periods (April 1996–December 2000 and January 2003–December 2008). This was
48 chosen so that the remaining data in the testing set would contain MF detections and flare
49 activity from time periods around the maximum and minimum levels of solar activity
50 (January 2001–December 2002 and January 2009–December 2010).
51
52
53
54
55
56
57
58

59 The machine-learning algorithm is trained by inputting the MF properties and their
60 corresponding flare/no-flare classifications from the training set time range. Then, when
61
62
63
64
65

1 the training process is completed, the predictions of the machine learning system is tested
2 against the training set time range and the prediction performance is evaluated by
3 following the same steps described in the previous subsection (Section 3.2.1).
4
5

6 **3.3. FEATURE SELECTION**

7
8
9 In this section, feature selection algorithms have been applied to the segmented data set
10 only (covering the entirety of April 1996– December 2010) to identify the most
11 significant SMART MF properties that are related to flare occurrence. Feature selection,
12 also known as variable selection or attribute selection, is the process of selecting a subset
13 of features according to certain criteria (Liu, 1998, 2008; Guyon, 2003). This process
14 enhances the efficiency and usability of a data set by removing features that are
15 irrelevant, redundant, and leading to noise (Liu and Motoda, 1998, 2007; Guyon and
16 Elisseff, 2003). Feature selection has been applied in many areas of research, such as
17 genomic analysis, text mining, and image retrieval. As far as the authors are aware, this is
18 the first time that feature selection has been applied to solar data. The application of
19 feature selection in this study should enable us to determine the MF properties that are
20 most relevant to flare occurrence, and thus enhance our understanding of the underlying
21 physics behind flare occurrence.
22
23
24
25
26
27
28
29
30

31 Two different feature evaluation algorithms have been applied here – Correlation-based
32 Feature Selection (CFS) and Minimum Redundancy Maximum Relevance (MRMR). The
33 reader is referred to the Appendix for a detailed description of these methods. These
34 evaluation methods can be applied in combination with differing search methods and data
35 types to create different feature selection processes. We consider multiple combinations
36 of each feature evaluation and search method to create 11 different feature selection
37 processes, presented in Table 8. Feature selection methods have been utilised using the
38 Waikato Environment for Knowledge Analysis open source package (Hall et al., 2009)
39 and the feature selection repository developed at the Data Mining and Machine Learning
40 Laboratory at Arizona State University².
41
42
43
44
45
46
47
48

49 The feature selection experiment is carried out as follows. Initially, cross-validation is
50 applied to the segmented dataset (consisting of all 21 MF properties) in order to select
51 50% of the data in a random manner and then feature selection is applied. This is repeated
52 20 times and the most common set of MF properties is recorded for each feature selection
53 process. It is important to note that CFS determines the best, but un-ranked, MF
54
55
56
57
58
59

60 ² Feature Selection at Arisona State University:
61 <http://featureselection.asu.edu/index.php>
62
63
64
65

properties while MRMR is set to rank the 10 most significant MF properties according to their importance. CCNN machine learning is then applied to determine the prediction capability of the selected MF properties, using both segmented and operational data sets, enabling direct comparison to the prediction capability of the full set 21 MF properties. Machine learning using cross-validation is applied for this purpose, as previously described in Section 3.2.1.

4. Results

In this section, the results achieved from applying the various methods described in the previous section are presented and discussed. The results are presented according to the experimental aims of this work: Section 4.1 discusses machine learning validation with respect to overall prediction capability, while Section 4.2 aims to determine the properties most related to flaring.

4.1. FLARE PREDICTION CAPABILITY OF SMART MF PROPERTIES

CCNN machine learning has been applied to determine the capability of the 21 MF properties generated by SMART to predict flares at and above the C1.0 level within the following 24-hour period. The performance of the machine-learning system has been investigated in two separate ways. Full data set cross-validation (Section 4.1.1) determines the overall effectiveness of the machine learning, and time-separated training and testing (Section 4.1.2) realistically validates the system.

4.1.1. Machine Learning using Cross Validation

Machine learning is applied using cross-validation to determine the overall prediction capability of the investigated data and to set a benchmark performance of the system using the methods described in Section 3.2.1. This process separately uses segmented and operationally associated MF-flare data covering the entire time range (April 1996–December 2010). The prediction measures achieved for both of the associated data sets are shown in Table 5. It can be seen that using the segmented data set provides higher prediction measures than using the operational data set. Segmentation thus allows the machine learning to more easily discriminate between flaring and non-flaring MFs. As outlined in Section 3.1, this is because the no-flare component of the segmented data set consists of MF detections that are clearly separated from flares (i.e., no flare occurs in a +/-48-hour period), while the no-flare component of the operational data set will consist of MF detections recorded just after flares (i.e., only requires that no flare occurs in the following 24-hour period). However, despite the reduced level of prediction measures

achieved, the operational data set results are regarded as the realistic capability of the system to provide flare prediction in a near real-time operational mode.

Table 5: Prediction measures achieved from applying machine learning with cross-validation on the segmented and operational data sets covering April 1996–December 2010.

Association Method	Forecast Verification Measures							
	MSE	TPR	FPR	TNR	FNR	FAR	ACC	HSS
Segmented	0.017	0.662	0.008	0.992	0.338	0.176	0.974	0.720
Operational	0.024	0.455	0.010	0.990	0.545	0.278	0.962	0.539

4.1.2. Machine Learning with Time-separated Training and Testing Sets

In contrast to the process of cross-validation, the flare prediction capability of the machine-learning algorithm is investigated by training and testing the system on data from completely separate time ranges (described in Section 3.2.2). This process makes use of data in the time range April 1996–December 2000 and January 2003–December 2008 for training, and data in the time range January 2001–December 2002 and January 2009–December 2010 for testing. The training and testing sets are obtained from both the segmented and operational MF-flare associated data sets detailed in Section 3.1. This approach of using common time ranges for training and testing is adopted to ensure that direct comparisons can be carried out between the different combinations of the training/testing data sets (i.e., segmented/segmented, operational/operational, and segmented/operational). The prediction measures achieved by these three training/testing combinations are given in Table 6. Once again, the highest prediction measures are achieved when the machine-learning algorithm is both trained and tested on segmented data. However, these measures do not reflect the actual capability of the system if it were run operationally because the data supplied to the prediction system does not contain all MF detections (as previously discussed in Section 3.1).

Table 6: Prediction measures achieved from applying machine learning on different combinations of time-separated segmented and operational data sets. Note that the training sets always cover the combined time range of April 1996–December 2000 and January 2003–December 2008, while the testing sets always cover the combined time range of January 2001–December 2002 and January 2009–December 2010.

Association Method		Forecast Verification Measures							
Testing Set	Training Set	MSE	TPR	FPR	TNR	FNR	FAR	ACC	HSS
Segmented	Segmented	0.016	0.677	0.006	0.994	0.323	0.118	0.976	0.754
Operational	Operational	0.024	0.523	0.011	0.989	0.477	0.258	0.963	0.595
Operational	Segmented	0.025	0.662	0.021	0.979	0.338	0.349	0.961	0.636

The prediction measures resulting from our new machine-learning system are compared in Table 7 to one of the industry’s standard flare prediction systems, ASAP (Qahwaji and Colak, 2009). This system also uses machine learning to predict flares at and above the C-class level within 24 hours, but ASAP was trained on data covering a longer time period, did not discard active regions further than 45° from solar disk centre, and was tested on a data set that contained less number of detections in comparison to the number detections in the testing sets used here, given in Table 4. Overall, Tables 6 and 7 indicate that the use of SMART MF properties in our machine learning system has achieved significantly improved flare prediction accuracy over that of ASAP. Further comparison to HSS values achieved by alternative prediction methods are presented and discussed in Section 5.

Table 7: Prediction measures achieved by ASAP.

Testing Period	Training Period	Total Detections	Forecast Verification Measures				
			TPR	ACC	FAR	MSE	HSS
Feb1999–Dec2002	Jan1982–Jan1999, Jan2003–Dec2006	40,534	0.814	0.805	0.301	0.146	0.512

4.2. MF PROPERTIES MOST RELATED TO FLARE OCCURRENCE

In the previous sections we determined that the segmented form of MF-flare association is capable of achieving the highest prediction performances within our CCNN machine-learning system. In this section we will use only segmented data, as we are interested in finding which MF properties are most related to predicting flare occurrence using the feature selection methods described in Section 3.3. The output from each of the 11 feature selection processes is presented in Table 8, where the MF properties are listed in terms of the property IDs of Table 1. The results from these feature selection processes were grouped into four categories to study the frequency of property selection: 1) the union of all CFS processes, 2) MRMR-MIQ, 3) MRMR-FCQ, and 4) MRMR-FCQ. Table 9 presents the selection frequency of MF properties, with individual rows indicating properties that appear in all 4, at least 3, at least 2, and at least 1 of these categories.

Table 8: The initial SMART MF properties that have been selected by each feature selection process.

Evaluation Method	Search Method	Data Type	Output Type	Selected SMART MF Property IDs
CFS	BestFirst Backward	Normalised	Subset	v7, v9, v13
CFS	BestFirst Bidirectional	Normalised	Subset	v7, v9, v13
CFS	BestFirst Forward	Normalised	Subset	v7, v9, v13
CFS	GreedyStepwise Backward	Normalised	Subset	v5, v6, v7, v13, v14, v15, v19, v20, v21
CFS	GreedyStepwise Forward	Normalised	Subset	v7, v9, v13
CFS	BestFirst Backward	Discretised	Subset	v9, v13, v14, v18, v19, v20
CFS	BestFirst Bidirectional	Discretised	Subset	v9, v13, v14, v18, v19, v20
CFS	BestFirst Forward	Discretised	Subset	v9, v13, v14, v18, v19, v20
MRMR-MIQ	Forward	Discretised	Weighted	v13, v21, v20, v19, v18, v17, v16, v15, v14, v11
MRMR-FCQ	Forward	Normalised	Weighted	v13, v4, v14, v15, v21, v7, v6, v20, v11, v19
MRMR-FCD	Forward	Normalised	Weighted	v13, v21, v14, v20, v15, v4, v5, v19, v7, v18

Table 9: Selection frequencies of SMART MF properties from the four feature selection process categories.

Group No.	No. of Properties	Selection Frequency	SMART Property ID																				
			v1	v2	v3	v4	v5	v6	v7	v8	v9	v10	v11	v12	v13	v14	v15	v16	v17	v18	v19	v20	v21
G1	6	4													X	X	X				X	X	X
G2	8	>=3							X						X	X	X			X	X	X	X
G3	12	>=2				X	X	X	X				X		X	X	X			X	X	X	X
G4	15	>=1				X	X	X	X		X		X		X	X	X	X	X	X	X	X	X

The machine learning and cross-validation method presented in Section 3.2 was applied to each of the four groups of MF properties listed in Table 9. In order to examine the significance of the prediction capabilities for each of the property groups, their prediction performances were compared to that of the complete set of MF properties (i.e., all 21 SMART properties under consideration) with the results presented in Tables 10 and 11 for the segmented and operational data sets, respectively. These findings show that a prediction capability comparable to that using all 21 MF properties can be achieved from the set of 6 properties that were selected most frequently. However, including a greater number of MF properties leads to a marginally higher prediction performance.

Table 10: Prediction capability measures of the four feature-selected property groups in Table 9, using segmented data set.

Benchmark – Machine Learning using Cross Validation on all 21 MF Properties								
	MSE	TPR	FPR	TNR	FNR	FAR	ACC	HSS
	0.017	0.662	0.008	0.992	0.338	0.176	0.974	0.720
Experiment – Machine Learning using Cross Validation on Selected MF Properties								
Group No	MSE	TPR	FPR	TNR	FNR	FAR	ACC	HSS
G1	0.018	0.610	0.007	0.993	0.390	0.163	0.972	0.690
G2	0.018	0.610	0.007	0.993	0.390	0.157	0.972	0.691
G3	0.017	0.650	0.008	0.992	0.350	0.168	0.973	0.716
G4	0.017	0.659	0.008	0.992	0.341	0.165	0.974	0.723

Table 11: Prediction capability measures of the four feature-selected property groups in Table 9, using operational data set.

Benchmark – Machine Learning using Cross-Validation on all 21 MF Properties								
	MSE	TPR	FPR	TNR	FNR	FAR	ACC	HSS
	0.024	0.455	0.010	0.990	0.545	0.278	0.962	0.539
Experiment – Machine Learning using Cross-Validation on Selected MF Properties								
Group No	MSE	TPR	FPR	TNR	FNR	FAR	ACC	HSS
G1	0.025	0.440	0.009	0.991	0.560	0.277	0.961	0.528
G2	0.025	0.454	0.010	0.990	0.546	0.292	0.961	0.533
G3	0.025	0.457	0.010	0.990	0.543	0.286	0.962	0.538
G4	0.024	0.467	0.011	0.989	0.533	0.288	0.962	0.545

The rank ordering of SMART MF properties towards flare prediction is, according to their frequency of selection in Table 9: **Error! Reference source not found.**

1. L_{NL} , L_{SG} , ∇_{MAX} , WL_{SG} , R^* , and WL_{SG}^* ,
2. Φ_- , R ,
3. A , Φ , Φ_+ , B_{MAX} ,
4. $\Delta\Phi/\Delta t$, ∇_{MEAN} , ∇_{MEDIAN} .

The first group lists the six MF properties chosen by each of the four feature selection categories. These properties are all related to magnetic neutral lines and are all extensive quantities (with the exception of ∇_{MAX}). These properties are commonly considered to be highly relevant to flaring (e.g., Cui et al. 2006; Schrijver 2007; Falconer et al. 2009) due to their indication of non-potentiality in magnetic field topology. The rank ordering of the less frequently selected MF properties in terms of significance for flare prediction is not surprising, with extensive measurements of total magnetic flux and area being generally

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

more relevant than intensive measurements such as statistical moments of the magnetic field distribution (Welsch et al. 2009).

5. Discussion and Future Work

We have compared the flare prediction accuracy of training a machine learning system using both segmented and non-segmented (i.e., operational) data sets to determine whether removing portions of the MF population from the machine-learning training results in an improvement in prediction. The system has been validated using a testing data set segmented in the same way and an operational data set to illustrate the important difference between segmented and operational testing. To the author’s knowledge, this is the first time that this form of comparison has been presented.

The results from machine learning using cross-validation show that segmented training and testing is more successful than non-segmented (i.e., operational) training and testing (see Tables 5 and 6). This comparison shows that the value of FAR increases when operational training and testing is used, indicating that the system over predicts flare from MFs that were observed within 24–48 hours before or 24 hours after a flare (i.e., the MF detections excluded from the segmented training and testing). This could be a result of the operational system including MF detections within 24 hours after a flare if the footprints of field topologies capable of producing low-magnitude flares do not significantly change in the photosphere over the course of a flare. MF detections observed shortly after a flare may then be predicted to flare, as they are similar to the machine-learned pre-flare state, resulting in an increased number of false positive predictions.

It is instructive to compare our average segmented HSS result of 0.72 (with a standard deviation of 0.01) from Table 5 to the results of previous studies that also use segmented data. Yu et al. (2009) achieve a mean HSS of 0.65 (X. Huang, private communication) in predicting the cumulative equivalent of at least ten C1.0 flares (or one M1.0 class flare) within 48 hours, but discard observations that do not produce at least one C-class flare. Using the same segmentation and prediction, Yu et al. (2010a, 2010b) achieve a maximum HSS of 0.77 and 0.69, respectively. Mason and Hoeksema (2010) attain a HSS of 0.69 in predicting X-class flares within 6 hours of an observation, but do not predict for features that fall between two thresholds of property evolution in the 40 hours prior to a particular observation. Our particular choice of segmentation may have other benefits than simply increasing flare prediction training accuracy. The segmented dataset used here discards MF detections that have exhibited any flare history in the previous 48 hours. This will likely lead to the more accurate prediction of flaring “all-clear” periods because

1 active regions that have flared in the past have a high potential to flare again in the future
2 (Wheatland, 2005).
3

4 In contrast to the segmented data set, the operational data set uses all MF detections
5 distinguished only by their flaring status in the following 24 hours. The results for the
6 operational data set thus show the actual capability of the system for flare forecasting in
7 an operational manner. We achieve an average HSS of 0.54 (with a standard deviation of
8 0.02) for machine learning using cross-validation on our operational data set (Table 5)
9 and compare this to other work that did not segment their data. Barnes and Leka (2008)
10 test an operational data set using discriminant analysis, achieving a maximum HSS of
11 0.15 in predicting at least one M- or X- class flare within 24 hours. The major departure
12 between these results is likely to come from the inclusion of predicting C-class flares in
13 our system, which are more common than M- or X- class flares. In addition, we train the
14 system on a large data set containing periods of minimum and maximum solar activity to
15 expose the system to the most complete and diverse magnetic property parameter range
16 that is possible. Colak and Qahwaji (2009) report a maximum HSS of 0.51 for ASAP in
17 predicting at least one C-, M-, or X- class flare within 24 hours. These are quite similar
18 results, with the marginal improvement offered by our system probably arising from the
19 use of many magnetic field properties.
20
21
22
23
24
25
26
27
28
29
30

31 The results from the machine learning using time-separated training and testing data sets
32 show that the highest prediction performance for operational testing is obtained when the
33 system is trained on segmented data (Table 6). The value of HSS reached by this method
34 (0.64) lies between the cross-validation results for the segmented data (0.72+/-0.01) and
35 the operational data (0.54+/-0.02). In addition, the combination of segmented training and
36 operational testing outperforms that of operational training and testing (HSS=0.59). This
37 indicates that the machine-learning system is capable of accurately applying the more
38 clearly separated flare/no-flare parameter distributions in the segmented training set to the
39 less distinct operational testing set. It is worth noting that this segmented training with
40 operational testing also outperforms the operational training and testing scheme of ASAP
41 (HSS=0.51).
42
43
44
45
46
47
48
49
50

51 Another aim for this work was the investigation of which MF properties are most
52 significantly related to flare occurrence. This should provide insight into the physical
53 relationship between photospheric magnetic fields and flare activity in the corona. The
54 prediction capabilities of feature-selected MF properties subsets were determined and it
55 was found that smaller sets of MF properties achieve equivalent prediction performances
56 to that achieved by all 21 SMART MF properties (Tables 10 and 11). The MF properties
57 that are related to the polarity separation line are seen to be the most significant (Tables 8
58
59
60
61
62
63
64
65

1 and 9). This is not surprising, as these properties are proxies for the degree of non-
2 potentiality within a MF. Non-potentiality is believed to be one of the most important
3 factors in enabling flares to occur as it allows suitable amounts of energy to be stored in
4 the magnetic field (Régnier and Priest, 2007). In addition, five of the six most significant
5 properties were found to be extensive properties, with the sixth being the maximum field
6 gradient that is an intensive property.
7
8
9

10 Previous flare prediction systems have been limited in a number of ways, including
11 automation, accuracy, and the ability to make a prediction for all magnetic features within
12 some observational limits. SMART-ASAP is designed to work in an operational setting –
13 it is completely automated, uses real-time data, and runs in a matter of minutes. The true
14 prediction capability of the system has been evaluated here using a number of verification
15 measures so its performance can be directly compared to that of other prediction systems.
16 However, this work is inherently limited in that it uses snapshot information about the
17 magnetic field at the photosphere to predict activity in the corona. It is our belief that to
18 surpass the present HSS barrier of 0.8, the evolution of magnetic field needs to be taken
19 into account when predicting solar flares. To improve the SMART-ASAP system, we
20 intend on adding more MF properties to the system, such as Ising energy (Ahmed, 2010),
21 in addition to investigating the difference in prediction capability for different peak flare
22 magnitudes.
23
24
25
26
27
28
29
30
31

32 To summarise, the main conclusions of the experiments presented in this paper are:
33
34

- 35 • CCNN machine learning managed to successfully classify SMART MF detections as
36 flaring or non-flaring with a HSS of 0.72 for our segmented data set and 0.54 for our
37 operational data set, using cross-validation.
38
- 39 • The highest HSS value achieved for the operational testing data (0.64) was achieved
40 when the system was trained using segmented data.
41
- 42 • A small set of SMART MF properties (i.e., 6) can achieve comparable prediction
43 performance to that provided by the full set of 21 MF properties. However, flare
44 predictions based on sets with higher numbers of MF properties result in marginally
45 higher prediction performance.
46
- 47 • The SMART MF properties that are most related to flare occurrence are those
48 involving neutral lines properties.
49
- 50 • SMART and machine learning systems are both automated. The execution time of
51 SMART is about 20-60 seconds and of the machine learning is about 5 seconds, on a
52
53
54
55
56
57
58
59
60
61
62
63
64
65

computer with 2.66 GHz Intel Core 2 Duo with 2GB of 800MHz DDR2 SDRAM.

Hence, both systems can be integrated to run in real-time.

Overall, the technologies and the findings that have been presented in this paper can work as a corner stone to develop accurate flare prediction systems and to provide an improved understanding of the underlying physics behind flare occurrence. The system presented here will be modified to use SDO/HMI magnetograms and run on both SolarMonitor.org and spaceweather.inf.brad.ac.uk in near real-time.

Acknowledgements

PAH is partially supported by an ESA/Prodex grant administered by Enterprise Ireland and a European Commission Framework Programme 7 grant (HELIO). DSB is supported by a Marie Curie Intra-European Fellowship (IEF) for Career Development funded under the European Commission's Seventh Framework Programme (FP7).

APPENDIX

Feature Selection Algorithms

Feature selection consists of two processes – search and evaluation. Feature search selects feature candidates and feeds them to the feature evaluation in order to determine their utility. This process is repeated so different subsets are evaluated, until the optimum subset of features is achieved. The best search strategy generates all possible combinations of feature subsets. However, this approach is exhaustive when the numbers of investigated features are large. Therefore, heuristic search methods are adopted. The common heuristic search approaches are: *forward* search, when the search starts with no features and successively adds features; *backward* search, when the search starts with all features and successively removes features; *bidirectional* search, when the search starts somewhere in the middle and moves outward from the starting point. Feature evaluation can be conducted using different methods: filters have been adopted in this work. Filter methods evaluate a subset of features using correlation methods. They are fast, efficient, and most frequently used in real world applications (Liu et al. 2010). Two feature evaluation filter methods have been applied in this work: Correlation-based Feature Selection (CFS) and Minimum Redundancy Maximum Relevance (MRMR).

A. CFS

CFS is a supervised feature evaluation method that selects a subset of features that are highly correlated with the class and uncorrelated with each other. Each feature is selected according to its ability to predict the class in areas that are not already predicted by other features. This algorithm can be applied to discrete or continuous data. CFS evaluates features using Equation 1,

$$M_s = \frac{k\overline{r_{cf}}}{\sqrt{k + k(k-1)\overline{r_{ff}}}} \quad \text{Equation 1}$$

where M_s is the heuristic “merit” of a feature subset containing k features, $\overline{r_{cf}}$ is the mean feature-class correlation, and $\overline{r_{ff}}$ is the average feature-feature inter-correlation. The correlation type is determined according to the class type when symmetrical uncertainty correlation is applied for discrete classes, while Pearson’s correlation is applied to continuous classes. CFS feature evaluation has been applied with two common heuristic search methods – greedy hill climbing (or greedy stepwise) and best-first. Greedy hill climbing adopts a forward or backward search approach to select feature candidates by searching the entire set of features as long as the feature evaluation does not degrade. Best-first adopts a forward, backward, or bidirectional search approach to select feature candidates. Best-first allows backtracking during the search so, when a certain path looks less promising, best-first can backtrack to a more promising previous subset and continue from there. However, a stopping criterion is applied if a limited number of fully expanded subsets (normally 5) result in no further improvement. More details about CFS can be obtained from (Hall, 1999).

B. MRMR

MRMR is a supervised feature evaluation method that selects features that are mutually dissimilar to each other, but highly related to the class. The selected features are ranked according to their importance, and the user determines the size of the selected features. MRMR can be applied to discrete or continuous data. For discrete data, the mutual information is used to calculate the level of similarity between the features to measure the minimum redundancy using Equation 2, and it is also used to calculate the discriminant power between the features and the class to measure the maximum relevance using Equation 3. For continuous data, the Pearson correlation coefficient is used to calculate the similarity between the features to measure the minimum redundancy using Equation 4, while F -test is used to calculate the maximum relevance between the features and the class using Equation 5.

$$\min W, W_I = \frac{1}{|S|^2} \sum_{i,j \in S} I(i,j) \quad \text{Equation 2}$$

$$\max V, V_I = \frac{1}{|S|} \sum_{i \in S} I(h,i) \quad \text{Equation 3}$$

$$\min W, W_c = \frac{1}{|S|^2} \sum_{i,j} |c(i,j)| \quad \text{Equation 4}$$

$$\max V, V_F = \frac{1}{|S|} \sum_{i \in S} F(i,h) \quad \text{Equation 5}$$

Where S is the set of features; $I(i, j)$ is the mutual information between features i and j ; $c(i, j)$ is the correlation between feature i and j ; $F(i, j)$ is the F-statistic between i and j ; and h is the target class. Features are ranked by optimising the minimum redundancy and maximum relevance for each feature by subtracting or dividing the two values, as shown in Equation 6 and 7. For discrete data, this is named as Mutual Information Difference (MID) or Mutual Information Quotient (MIQ). For continuous data, it is named as F-test Correlation Difference (FCD) or F-test Correlation Quotient (FCQ). However, MIQ and FCQ seem to provide better results (Ding and Peng, 2005).

$$\text{MID or FCD, } \max(V - W) \quad \text{Equation 6}$$

$$\text{MIQ or FCQ, } \max(V / W) \quad \text{Equation 7}$$

Where V is the minimum redundancy and W is the maximum relevance. MRMR uses heuristic forward search to add features according to their importance, which are measured using Equation 6 or 7 above. More details about MRMR can be obtained from (Ding and Peng, 2005).

Figures

Figure 1: An example of all SMART MF detections on 29 October 2003. AR (active region) denotes features classified as multipolar, while PL (plage) and UD (unipolar decaying) denote two different classes of unipolar feature.

References

- Ahmed, O., Qahwaji, R., Colak, T., Dudok De Wit, T., and Ipson, S.: 2010, The Visual Computer, 26, 385.
- Balch, C.C.: 2008, Space Weather, 6, 13.

- 1 Barnes, G. and Leka, K.D.: 2008, *Astrophys. J.*, 688, L107.
- 2 Colak, T. and Qahwaji, R.: 2009, *Space Weather*, 7, 12.
- 3
- 4 Colak, T. and Qahwaji, R.: 2010, *Automated Prediction of Solar Flares*, LAP LAMBERT
5 Academic Publishing, Saarbrücken, Germany, p. 74.
- 6
- 7 Committee on the Societal and Economic Impacts of Severe Space Weather Events: 2008, *Severe
8 Space Weather Events - Understanding Societal and Economic Impacts*, The National Academies
9 Press, Washington D.C., USA.
- 10
- 11 Cui, Y., Li, R., Zhang, L., He, Y., and Wang, H.: 2006, *Solar Phys.*, 237, 45.
- 12
- 13 Ding, C. and Peng, H.: 2005, *Journal of Bioinformatics and Computational Biology*, 3, 185.
- 14
- 15 Fawcett, T.: 2006, *Pattern Recognition Letters*, 27, 861.
- 16
- 17 Frank, R.J., Davey, N., and Hunt, S.P.: 2001, *Journal of Intelligent and Robotic Systems*, 31, 91.
- 18
- 19 Falconer, D. A., R. L. Moor, Gary, G. A., and Adams, M.: 2009, *Astrophysical J.*, 700(2): L166.
- 20
- 21 Gallagher, P.T., Moon, Y., and Wang, H.: 2002, *Solar Phys.*, 209, 171.
- 22
- 23 Gopalswamy, N., Barbieri, L., Lu, G., Plunkett, S.P., and Skoug, R.M.: 2005, *Geophys. Res. Lett.*,
24 32, L03S01.
- 25
- 26 Guyon, I. and Elisseeff, A.: 2003, *Journal of Machine Learning Research*, 3, 1157.
- 27
- 28 *Hall, M.A.: 1999, Correlation-based Feature Selection for Machine Learning*, PhD Thesis: The
29 University of Waikato, Hamilton, New Zealand.
- 30
- 31 Higgins, P.A., Gallagher, P.T., McAteer, R.T.J., and Bloomfield, D.S.: 2010, *Advances in Space
32 Res.*, in press.
- 33
- 34 Hall, M., Frank, E., Holmes, G., Pfahringer, B., Reutemann, P., and Witten, I.H.: 2009, *Open
35 Source Analytics*, 11, 10.
- 36
- 37
- 38 Jing, J., Song, H. Abramenko, V., Tan, C., and Wang, H.: 2006, *Astrophys. J.*, 644, 1273.
- 39
- 40 Liu, H., Motoda, H., Setiono, R., and Zhao, Z.: 2010, in Liu, H., Motoda, H., Setiono, R. Zhao, Z.
41 (eds.), *The Fourth Workshop on Feature Selection in Data Mining*, Hyderabad, India, p. 4.
- 42
- 43 Liu, H. and Motoda, H.: 2008, *Computational Methods of Feature Selection*, Chapman and
44 Hall/CRC, New York, U.S., p. 4.
- 45
- 46 Liu, H. and Motoda, H.: 1998, *Feature Selection for Knowledge Discovery and Data Mining*,
47 Kluwer Academic Publishers, Boston, USA, p. 17.
- 48
- 49
- 50 Leka, K.D. and Barnes, G.: 2007, *Astrophys. J.*, 656, 1173.
- 51
- 52 Mason, J.P. and Hoeksema, J.T.: 2010, *Astrophys. J.*, 723, 634.
- 53
- 54 McIntosh, P.S.: 1990, *Solar Phys.*, 125, 251.
- 55
- 56 Messerotti, M., Zuccarello, F., Guglielmino, S., Bothmer, V., Lilensten, J., Noci, G., Storini, M.,
57 and Lundstedt, H.: 2009, *Space Sci. Rev.*, 147, 121.
- 58
- 59 Qahwaji, R. and Colak, T.: 2006, in Chu, H.W., Aguilar, J, Rische, N., and Azoulay, J. (eds.), *3rd
60 International Conference on Cybernetics and Information Technologies*. Florida, USA, p. 192.
- 61
- 62
- 63
- 64
- 65

1 Qahwaji, R. and Colak, T.: 2007, *Solar Phys.*, 241, 195.

2 Régnier, S. and Priest, E.R.: 2007, *Astrophys. J.*, 669, L53.

3
4 Song, H., Tan, C., Jing, J., Wang, H., Yuchyshyn, V., and Abramenko, V.: 2008, *Solar Phys.*, 254,
5 101.

6
7 Schrijver, C. J.: 2007, *Astrophys J.*, 665, L117

8
9 Yu, D., Huang, X., Wang, H., and Cui, Y.: 2009, *Solar Phys.*, 255, 91.

10
11 Yu, D., Huang, X., Hu, Q., Zhou, X., Wang, H., and Cui, Y.: 2010a, *Astrophys. J.*, 709, 321.

12
13 Yu, D., Huang, X., Wang, H., Cui, Y., Hu, Q., and Zhou, R.: 2010b, *Astrophys. J.*, 710, 869.

14
15 Yuan, Y., Shih, F.Y., Jing, J., and Wang, H.-M.: 2010, *Res. Astron. Astrophys.*, 10, 785.

16
17 Welsch, B. T., Yan, L., Schuck, P. W., and Fisher, G. H.: 2009, *Astrophys. J.*, 705, 821.

18
19 Wheatland, M.S.: 2005, *Space Weather*, 3, 11.

20
21 Zhang, J., Dere, K.P., Howard, R.A., Kundu, M.R., and White, S.M.: 2001, *Astrophys. J.*, 559,
22 452.

23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65