

# Short-Term Solar Flare Prediction Using Predictor Teams

Xin Huang · Daren Yu · Qinghua Hu · Huaning Wang · Yanmei Cui

Received: 18 November 2009 / Accepted: 8 March 2010 / Published online: 3 April 2010  
© Springer Science+Business Media B.V. 2010

**Abstract** A short-term solar flare prediction model is built using predictor teams rather than an individual set of predictors. The information provided by the set of predictors could be redundant. So it is necessary to generate subsets of predictors which can keep the information constant. These subsets are called predictor teams. In the framework of rough set theory, predictor teams are constructed from sequences of the maximum horizontal gradient, the length of neutral line and the number of singular points extracted from SOHO/MDI longitudinal magnetograms. Because of the instability of the decision tree algorithm, prediction models generated by the C4.5 decision tree for different predictor teams are diverse. The flaring sample, which is incorrectly predicted by one model, can be correctly forecasted by another one. So these base prediction models are used to construct an ensemble prediction model of solar flares by the majority voting rule. The experimental results show that the predictor team can keep the distinguishability of the original set, and the ensemble prediction model can obtain better performance than the model based on the individual set of predictors.

**Keywords** Active regions · Magnetic fields · Flares, forecasting · Ensemble learning

## 1. Introduction

There are two main aspects to improve the performance of the short-term solar flare prediction. One is to construct more informative predictors. The other is to build more powerful prediction models.

---

X. Huang (✉) · D. Yu · Q. Hu  
Harbin Institute of Technology, Power Engineering, Harbin, Heilongjiang Province, China  
e-mail: [huangxinhit@yahoo.com.cn](mailto:huangxinhit@yahoo.com.cn)

H. Wang  
Key Laboratory of Solar Activity, National Astronomical Observatories of Chinese Academy of Sciences, Beijing, China

Y. Cui  
Center for Space Science and Applied Research, Chinese Academy of Sciences, Beijing, China

On the one hand, many predictors of solar flares have been proposed. McIntosh (1990) applied the McIntosh classification of sunspots to reflect morphological characteristics of active regions. Bornmann and Shaw (1994) pointed out that the McIntosh parameters act as proxies for magnetic properties of an active region. More recently, many predictors based on the magnetic observation of active regions were proposed. Leka and Barnes (2003) extracted a series of parameters from the magnetic vector field. McAteer, Gallagher, and Ireland (2005) presented a fractal dimension measure to quantify the magnetic complexity of active regions. Schrijver (2007) considered the total unsigned flux  $R$  within  $\sim 15$  Mm of strong-field, high-gradient polarity separation lines as a characteristic appearance of magnetic fibrils carrying electrical currents when they emerge through the photosphere. Georgoulis and Rust (2007) defined the effective connected magnetic field ( $B_{\text{eff}}$ ) to measure the flaring potential in active regions. Cui *et al.* (2006) proposed three physical measures, the maximum horizontal gradient, the length of the neutral line, and the number of singular points, to describe the nonpotentiality and the complexity of the photospheric magnetic field of active regions. Next, Yu *et al.* (2009, 2010a) introduced the sequence of these predictors for the short-term solar flare prediction.

On the other hand, many methods have been developed to build the solar flare prediction model. McIntosh (1990) developed an expert system (Theo) to forecast solar flares. Wheatland (2004) proposed a Bayesian approach to flare prediction using the flaring records of an active region. Leka and Barnes (2007) applied discriminant analysis to determine which properties are associated with flare eruption. Wang *et al.* (2007) trained a neural network for solar flare prediction. Li *et al.* (2007) presented a method combining the support vector machine and the  $k$ -nearest neighbors to construct a solar flare forecasting model. Based on the work of Qahwaji and Colak (2007) and Colak and Qahwaji (2008, 2009) presented an automated hybrid computer platform (ASAP) for the realtime prediction of solar flares. Yu *et al.* (2010b) proposed the Bayesian network approach for the short-term solar flare prediction.

No single predictor was found to dramatically distinguish flaring from flare-quiet active regions (Barnes *et al.*, 2007). So we construct predictor teams from the sequence of predictors. On the one hand, predictor teams are more informative than the single predictor proposed by Cui *et al.* (2006). On the other hand, predictor teams reduce redundancy in the sequence of predictors proposed by Yu *et al.* (2009). Based on these predictor teams, base prediction models are built. The result of the proposed prediction model is the combination of all the outputs of base prediction models by the voting strategy (Kittler *et al.*, 1998).

The paper is organized as follows. In Section 2, the data are introduced. In Section 3, predictor teams are constructed. In Section 4, the ensemble prediction model based on predictor teams is proposed. In Section 5, the performance of the proposed prediction model is analysed. Finally, conclusions and future work are discussed in Section 6.

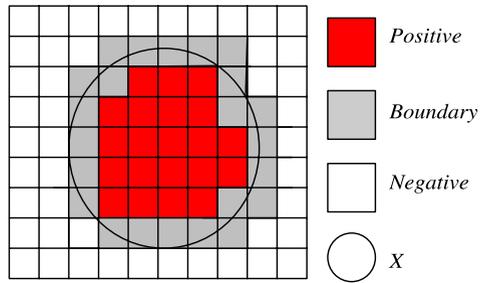
## 2. Data

The data have been introduced in Cui *et al.* (2006) and Yu *et al.* (2009), that contain 870 active regions and 3509 flares from 15 April 1996 to 10 January 2004.

The importance of a solar flare is conventionally described by its GOES class, for example, C, M or X. Within a certain interval, more than one flare may occur. The importance of these flares is summed up with weights. The total importance of flares is defined as follows:

$$I_{\text{tot}} = \sum C + 10 \times \sum M + 100 \times \sum X. \quad (1)$$

**Figure 1** Positive region, boundary and negative region in Pawlak’s rough set. Positive region is also called lower approximation. Summation of positive region and boundary is called upper approximation. Negative region is the complement of the upper approximation. Figure adopted from Hu, Liu, and Yu (2008).



Equation (1) considers the influence of all the flares within the forward-looking period. A forecasting model usually includes flares whose significance is larger than a certain threshold. Here, the threshold of  $I_{tot}$  is supposed to be 10.

Three predictors, the maximum horizontal gradient ( $|\nabla_h B_z|_m$ ), the length of the neural line ( $L$ ) and the number of singular points ( $\eta$ ), are extracted from SOHO/MDI full disk longitudinal magnetograms with a pixel size of  $2''$  and the noise level of 20G (Cui *et al.*, 2006). Yu *et al.* (2009) introduced the sequence of predictors to reflect the evolutionary information of the photospheric magnetic field by the sliding window method:

$$\mathbf{x}(t - W\Delta t) \cdots \mathbf{x}(t - \Delta t) \mathbf{x}(t) I_{tot}(t + F), \tag{2}$$

where  $\mathbf{x}(t)$  is the vector of predictors at time  $t$ .  $\mathbf{x} = \{|\nabla_h B_z|_m, L, \eta\}$ .  $\Delta t$  is 96 minutes which is the interval between the successive magnetograms.  $W$  is the length of the sequence.  $F$  is the forecasting time.  $I_{tot}(t + F)$  is the total importance of flares within  $F$ . When  $I_{tot} = 10$  and  $F = 48$  (hr), the appropriate length of sequence is 45 ( $W = 45$ ). So the predictors  $|\nabla_h B_z|_m(t - 45\Delta t), \dots, |\nabla_h B_z|_m(t)$  are numbered from 1 to 46,  $L(t - 45\Delta t), \dots, L(t)$  is numbered from 47 to 92, and  $\eta(t - 45\Delta t), \dots, \eta(t)$  is numbered from 93 to 138.

### 3. Predictor Teams

The information provided by the sequence of predictors is redundant for the short-term solar flare prediction. So it is necessary to find the subset of predictors which can maintain the discernibility of all the predictors. The rough set theory (Pawlak, 1991) can be used to characterize the discernibility of the set of predictors. In Pawlak’s rough set model, the samples ( $x$ ) with the same values of predictors ( $B$ ) are drawn together to form an equivalence class ( $[x]_B$ ). Due to the inconsistency in the data, the set  $X$  cannot be precisely described by the equivalence class. The lower approximation ( $\underline{B}X$ ) and the upper approximation ( $\overline{B}X$ ) approximately describing the set  $X$  are defined as follows:

$$\underline{B}X = \{x|[x]_B \subseteq X\}, \tag{3}$$

$$\overline{B}X = \{x|[x]_B \cap X \neq \emptyset\}. \tag{4}$$

As shown in Figure 1 (Hu, Liu, and Yu, 2008), the lower approximation is used to reflect the discernibility of the set of predictors. When the discernibility between a set and its subset is the same, the information provided by the two sets is equal. In order to decrease the redundant information, the reduct set, defined by the minimal subset of predictors which

keeps the discernibility of the original dataset, should be found. However, the problem of searching a minimal reduct set is computationally complicated. Therefore, the genetic algorithm is used to find the approximation of the optimal solutions (Wroblewski, 1998). In the genetic algorithm, the combination of predictors is called individual. The individual is encoded by bit strings: 1 for the selected predictor, contrariwise, the non-selected predictor is encoded by 0. We randomly choose the set of individuals as an initial population, and calculate the fitness of each individual which is defined as the difference between the combination of selected predictors and the set of all the predictors. The individuals with high value of fitness are selected, and they are bred through crossover and mutation operations to give birth to offspring (Liu *et al.*, 2009). This means that the old population is replaced by the new population with the higher value of fitness. The above process is repeated until the satisfactory individuals are selected. Reduct sets of the original predictors are generated. They are called predictor teams in the solar flare prediction method.

The established 11 predictor teams with 21 predictors are shown in Table 1. These predictor teams are a combination of three types of predictors, which are the maximum horizontal gradient at different times ( $|\nabla_h B_z|_m(t - i\Delta t)$ ,  $i = 0, \dots, 45$ ), the length of the neural line at different times ( $L(t - i\Delta t)$ ,  $i = 0, \dots, 45$ ) and the number of singular points at different times ( $\eta(t - i\Delta t)$ ,  $i = 0, \dots, 45$ ).

“The state of the photospheric magnetic field at any single time has limited bearing on the occurrence of solar flares”, as concluded by Leka and Barnes (2007), so predictors at different time are together used to forecast the eruption of solar flares (Yu *et al.*, 2009). However, the sequential information is redundant and it is necessary to decrease the redundant information, hence some predictors are selected to form the predictor team maintaining the information of the original set. But combinations of predictors are not unique, thus sequences of the maximum horizontal gradient ( $|\nabla_h B_z|_m(t - i\Delta t)$ ,  $i = 0, \dots, 45$ ), the length of the neural line ( $L(t - i\Delta t)$ ,  $i = 0, \dots, 45$ ) and the number of singular points ( $\eta(t - i\Delta t)$ ,  $i = 0, \dots, 45$ ) are grouped into 11 predictor teams by the rough set method. Using 21 predictors within each predictor team generated by the genetic algorithm, the positive region of the original set with 138 predictors is maintained. This means that the distinguishability of each predictor team is the same as the original set of predictors in the framework of the rough set theory. The rationality of the selection of predictors is explained from two aspects. From the geometrical points of view (Wang, Cui, and He, 2009), the maximum horizontal gradient is a measure of the local scale. It is influenced by the length of the neural line because the high gradient of the magnetic field usually appears in the vicinity of the neutral line, and the length of the neural line is a measure of the line scale. Both the maximum horizontal gradient and the length of the neural line are influenced by the number of singular points, which indicates the complexity of the topological structures of the magnetic field. The number of singular points is a surface scale measurement. For different geometrical scales, the information provided of these types of predictors is complementary, so they are grouped into a predictor team. From the sequential points of view, current predictors are more informative, thus in all the 11 predictor teams, the predictors nearby the current time  $t$  ( $|\nabla_h B_z|_m(t)$ ,  $|\nabla_h B_z|_m(t - \Delta t)$ ,  $|\nabla_h B_z|_m(t - 2\Delta t)$ ,  $L(t)$ ,  $L(t - \Delta t)$ ,  $\eta(t)$ ,  $\eta(t - \Delta t)$  and  $\eta(t - 2\Delta t)$ ) are taken. In order to maintain the information with the original set of predictors, some predictors at another time are randomly added to provide more comprehensive information.

**Table 1** Predictor teams generated by the rough set method. PT<sub>1</sub>–PT<sub>11</sub> stands for 11 predictor teams, respectively. The number  $i$  ( $i = 1, \dots, 138$ ) stands for the  $i$ th selected predictor. See Section 2 for an explanation of the predictors and their number ranges.

PT <sub>1</sub>	PT <sub>2</sub>	PT <sub>3</sub>	PT <sub>4</sub>	PT <sub>5</sub>	PT <sub>6</sub>	PT <sub>7</sub>	PT <sub>8</sub>	PT <sub>9</sub>	PT <sub>10</sub>	PT <sub>11</sub>
1	4	3	7	7	3	7	1	7	7	7
26	26	19	26	39	33	20	28	34	34	39
31	35	36	35	40	41	35	36	39	39	40
34	42	41	42	42	42	37	41	42	40	42
42	44	42	44	44	44	42	42	44	42	44
44	45	44	45	45	45	44	44	45	44	45
45	46	45	46	46	46	45	45	46	45	46
46	47	46	48	48	48	46	46	48	46	48
48	48	48	50	53	53	48	48	53	48	53
53	51	53	53	55	55	53	53	55	53	55
55	53	55	87	61	91	55	55	66	55	61
91	91	91	91	91	92	61	91	91	61	91
92	92	92	92	92	93	91	92	92	91	92
97	97	94	95	93	112	92	99	93	92	97
105	107	113	104	110	113	97	115	112	97	110
115	113	122	106	121	116	113	122	121	113	121
121	122	128	114	124	128	121	124	129	121	124
135	132	132	135	126	132	133	132	133	126	126
136	136	136	136	136	136	136	136	136	136	136
137	137	137	137	137	137	137	137	137	137	137
138	138	138	138	138	138	138	138	138	138	138

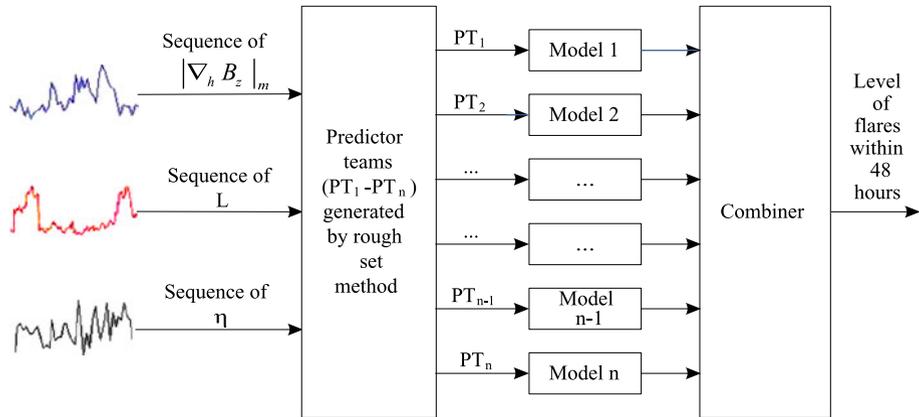
#### 4. Ensemble Prediction Model Using Predictor Teams

The short-term solar flare prediction is a complex problem; it is difficult for an individual model to precisely forecast the eruption of solar flares. Thus it is necessary to combine the results of multiple prediction models built on different predictor teams. The general schematic view of the ensemble prediction model of solar flares using predictor teams is shown in Figure 2.

Each predictor team can obtain the same discernibility of the original predictors. Based on these predictor teams, base prediction models are built respectively. Combining the same prediction models will not yield any improvement (Kuncheva, 2004), so the decision tree, which is an unstable algorithm, is used to learn the diverse prediction models based on different predictor teams. Finally, the voting strategy (Kittler *et al.*, 1998) is applied to combine results of the base prediction models.

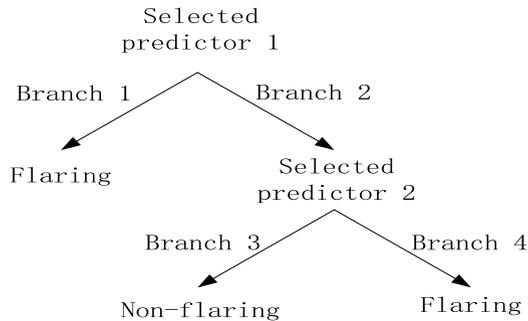
##### 4.1. Base Prediction Model

A decision tree can be constructed from a dataset using the divide and conquer strategy, that is, the best predictor is selected in each test node to split samples into smaller subsets. The top-down and recursive splitting technique is used to produce the following subtrees. When samples in the subset have the same classification or every possible test has the same



**Figure 2** General schematic view of the ensemble prediction model of solar flares using predictor teams.

**Figure 3** Example of decision tree. The selected predictor is used to split samples into subsets. Each branch corresponds to the value of the predictor used in its selected predictor. Each leaf node assigns the samples in the subset into a certain classification.



class distribution, the leaf node is generated. An example of a decision tree is shown in Figure 3. For different predictor teams, selected predictors are different, thus the generated trees are diverse among themselves. This is the instability of the decision tree algorithm, and the improvement of the ensemble prediction model comes from the diversity among the base prediction models. So the decision tree is well-suited to construct the base prediction model for the ensemble prediction model. Here, the C4.5 decision tree algorithm (Quinlan, 1993), which evaluates the quality of predictors using information gain ratios, is selected to generate base prediction models.

4.2. Combination of Base Prediction Models

After base prediction models have been built, the rule of the majority vote is used to combine outputs of base decision trees. The sample ( $x$ ) is to be assigned to one of the two possible classes ( $\omega_j, j = 1, 0$ ). If the  $i$ th base prediction model ( $M_i, i = 1, \dots, n$ ) labels  $x$  in  $\omega_j, d_{i,j}$  equals to one. Otherwise  $d_{i,j}$  equals to 0. The sample is assigned to the class which is the same as the result of most prediction models.

$$\omega = \begin{cases} \omega_1 & \text{if } \sum_{i=1}^n d_{i,1} > \sum_{i=1}^n d_{i,0}, \\ \omega_0 & \text{otherwise,} \end{cases} \tag{5}$$

where  $n$  is the number of the base prediction models.

**Table 2** Different outcomes of two-class prediction.

	Predicted positive class	Predicted negative class
Actual positive class	True Positive	False Negative
Actual negative class	False Positive	True Negative

## 5. Experimental Results

### 5.1. Performance Evaluation

Results of the proposed prediction model are grouped into “flaring” or “non-flaring”. The flaring samples are considered as positive class. Otherwise, they are considered as negative class. In this case, the prediction model has four different possible outcomes as shown in Table 2.

Samples correctly classified as positive are defined as True Positive (TP), while samples correctly classified as negative are defined as True Negative (TN). On the other hand, samples wrongly predicted as positive are defined as False Positive (FP) and samples wrongly predicted as negative are defined as False Negative (FN). The prediction performance is measured using the TP rate and TN rate (Witten and Frank, 2005).

The TP rate is defined as the ratio of the number of positive class samples predicted as positive to the number of actual positive class samples:

$$\text{TPrate} = \frac{\text{TP}}{\text{TP} + \text{FN}}. \quad (6)$$

The TN rate is defined as the ratio of the number of negative class samples predicted as negative to the number of actual negative class samples:

$$\text{TNrate} = \frac{\text{TN}}{\text{TN} + \text{FP}}. \quad (7)$$

TP rate and TN rate are used to evaluate the accuracy of flaring and non-flaring predictions, respectively. Furthermore, the Heidke skill score (HSS) is used to generally evaluate the performance of the proposed method (Jolliffe and Stephenson, 2003).

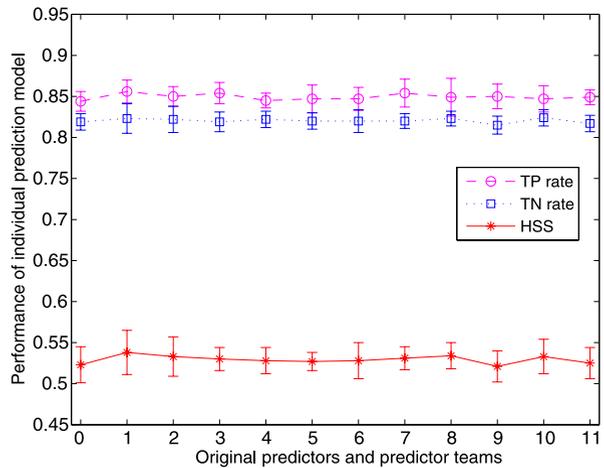
$$\text{HSS} = \frac{\text{PC} - \text{E}}{1 - \text{E}}, \quad (8)$$

where  $N = \text{TP} + \text{TN} + \text{FP} + \text{FN}$ ,  $\text{PC} = \frac{\text{TP} + \text{TN}}{N}$ , and  $\text{E} = \frac{(\text{TP} + \text{FN})(\text{TP} + \text{FP})}{N^2} + \frac{(\text{TN} + \text{FP})(\text{TN} + \text{FN})}{N^2}$ . E is PC for a random forecast, so this version of HSS shows the increase in predictive power over that of a random chance.

### 5.2. Results and Analyses

The sampling interval of the magnetogram is 96 minute. Each sampling interval is treated as a data point. For the 870 active regions, 8612 flaring samples and 39 732 non-flaring samples from 15 April 1996 to 10 January 2004 are selected to form the dataset. The ten-fold cross-validation is used to validate the performance of the proposed prediction model. The dataset is divided into ten folds and then nine folds are used for training and the remaining fold for

**Figure 4** Performance of C4.5 decision trees trained on original set of predictors or each predictor team ( $PT_1, \dots, PT_{11}$ ). 0 stands for the original set of predictors and 1–11 stands for the  $i$ th ( $i = 1, \dots, 11$ ) predictor team, respectively.



testing. This process is then repeated ten times, until each of the ten subsets is used once as the validation data. The performances of the prediction model are given as the mean of the ten testing results, and the uncertainty of the prediction model is estimated by the standard deviation of the ten testing results.

### 5.2.1. Performance of Prediction Model Using Single Predictor Team

Predictor teams generated by the rough set method keep the lower approximation consistent with the original dataset; this means that these predictor teams own the same amount of information as the original dataset in the framework of rough set theory. The base models are built by the C4.5 decision tree algorithm using these predictor teams. The performance of the base decision trees is shown in Figure 4. It is concluded that the information provided by the original set of predictors is redundant, and the predictor teams with 21 predictors basically maintain the distinguishability of the original 138 predictors.

### 5.2.2. Performance of Ensemble Model Using Predictor Teams

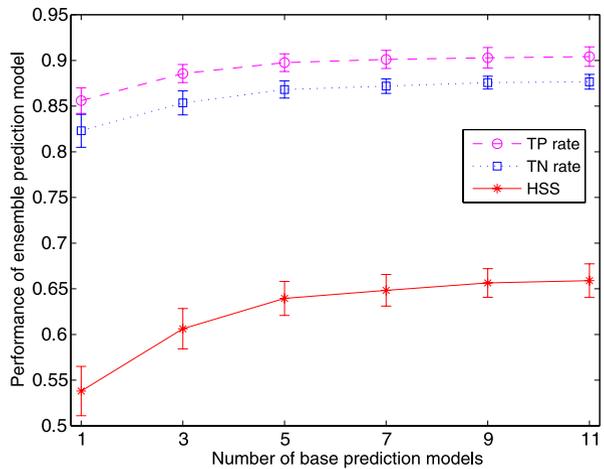
The performance of the ensemble model based on different numbers of base models is shown in Figure 5. The performance of the ensemble model based on up to five base decision trees is rapidly increased, because the complementary nature of these base decision trees is strong.

There are two main reasons that using an ensemble prediction model can achieve a better performance than using an individual prediction model (Ranawana and Palade, 2005). Firstly, the prediction correct rate of each individual model is greater than 0.5. Secondly, errors made by each individual model are uncorrelated. The performance of each base model in the present work is greater than 0.5, and the C4.5 decision trees learned from the different predictor teams are diverse. So the performance of the ensemble model for the short-term solar flare prediction is improved.

## 6. Conclusions

One single property of photospheric magnetic field has limited bearing on the flare productivity, so predictor teams are proposed to combine supplemental predictors as a whole.

**Figure 5** Performance of ensemble prediction model versus the number of base decision trees.



Meanwhile, predictor teams generated by the rough set method can preserve the prediction information of flares with a subset of the predictors, and the redundancy between predictors is decreased. However, it is difficult to obtain the optimum combination of predictors, so a series of approximate optimal combination are found. Base prediction models based on these predictor teams are combined into an ensemble model to make full use of the complementary information. These models learned from the different predictor teams by decision tree algorithm are diverse. This means that flares will not be correctly predicted by one model, while they will be predicted by another, so the prediction power of the ensemble model using predictor teams is enhanced and the performance of the ensemble model is improved. In conclusion, from the geometrical and sequential points of view, the generated predictor teams are reasonable and the ensemble model based on predictor teams is advised to be applied for the short-term solar flare prediction.

In the future, other quality measures of predictors (Kira and Rendell, 1992; Yu and Liu, 2003) can be used to generate the predictor teams, and the relationships between these predictor teams need to be studied. Other types of ensemble learning methods (Freund and Schapire, 1996; Ho, 1998) should be applied for the short-term solar flare prediction.

**Acknowledgements** This work is supported by the National Natural Science Foundation of China (NSFC) through Grant Nos. 10673017, 10733020, 10978011, the National Basic Research Program of China (973 Program) through Grant No. 2006CB806307, and the Program for Changjiang Scholars and Innovative Research Team in University (PCSIRT). We thank the SOHO/MDI consortium for the data. SOHO is a project of international cooperation between ESA and NASA. This paper has benefited from comments of the anonymous reviewer.

## References

- Barnes, G., Leka, K.D., Schumer, E.A., Della-Rose, D.J.: 2007, *Space Weather* **5**, S09002.
- Bormmann, P.L., Shaw, D.: 1994, *Solar Phys.* **150**, 127.
- Colak, T., Qahwaji, R.: 2008, *Solar Phys.* **248**, 277.
- Colak, T., Qahwaji, R.: 2009, *Space Weather* **7**, S06001.
- Cui, Y.M., Li, R., Zhang, L.Y., He, Y.L., Wang, H.N.: 2006, *Solar Phys.* **237**, 45.
- Freund, Y., Schapire, R.E.: 1996, In: *Proceedings of the Thirteenth International Conference on Machine Learning*, 148.
- Georgoulis, M.K., Rust, D.M.: 2007, *Astrophys. J.* **661**, 109.

- Ho, T.K.: 1998, *IEEE Trans. Pattern Anal. Mach. Intell.* **20**, 832.
- Hu, Q.H., Liu, J.F., Yu, D.R.: 2008, *Knowl.-Based Syst.* **21**, 294.
- Jolliffe, I.T., Stephenson, D.B.: 2003, *Forecast Verification: A Practitioner's Guide in Atmospheric Science*, Wiley, New York.
- Kira, K., Rendell, L.A.: 1992, In: *Proceedings of the Ninth International Workshop on Machine Learning*, 249.
- Kittler, J., Hatef, M., Duin, R.P.W., Matas, J.: 1998, *IEEE Trans. Pattern Anal. Mach. Intell.* **20**, 226.
- Kuncheva, L.I.: 2004, *Combining Pattern Classifiers: Methods and Algorithms*, Wiley-Interscience, New York.
- Leka, K.D., Barnes, G.: 2003, *Astrophys. J.* **595**, 1277.
- Leka, K.D., Barnes, G.: 2007, *Astrophys. J.* **656**, 1173.
- Li, R., Wang, H.N., He, H., Cui, Y.M., Du, Z.L.: 2007, *Chin. J. Astron. Astrophys.* **7**, 441.
- Liu, H., Abraham, A., Li, Y., Dalian, C.: 2009, *Rough Set Res., Adv. Theory Appl.* **174**, 261.
- McAteer, R.T.J., Gallagher, P.T., Ireland, J.: 2005, *Astrophys. J.* **631**, 628.
- McIntosh, P.S.: 1990, *Solar Phys.* **125**, 251.
- Pawlak, Z.: 1991, *Rough Sets: Theoretical Aspects of Reasoning About Data*, Kluwer Academic, Dordrecht.
- Qahwaji, R., Colak, T.: 2007, *Solar Phys.* **241**, 195.
- Quinlan, J.R.: 1993, *C4.5: Programs for Machine Learning*, Morgan Kaufmann, San Francisco.
- Ranawana, R., Palade, V.: 2005, *Neural Comput. Appl.* **14**, 122.
- Schrijver, C.J.: 2007, *Astrophys. J.* **655**, 117.
- Wang, H.N., Cui, Y.M., He, H.: 2009, *Res. Astron. Astrophys.* **9**, 687.
- Wang, H.N., Cui, Y.M., Li, R., Zhang, L.Y., He, H.: 2007, *Adv. Space Res.* **42**, 1464.
- Wheatland, M.S.: 2004, *Astrophys. J.* **609**, 1134.
- Witten, I.H., Frank, E.: 2005, *Data Mining: Practical Machine Learning Tools and Techniques*, Morgan Kaufmann, San Francisco.
- Wróblewski, J.: 1998, *Rough Sets Knowl. Discov.* **2**, 471.
- Yu, L., Liu, H.: 2003, In: *Proceedings of the Twentieth International Conference on Machine Learning*, 856.
- Yu, D.R., Huang, X., Wang, H.N., Cui, Y.M.: 2009, *Solar Phys.* **255**, 91.
- Yu, D.R., Huang, X., Hu, Q.H., Zhou, R., Wang, H.N., Cui, Y.M.: 2010a, *Astrophys. J.* **709**, 321.
- Yu, D.R., Huang, X., Wang, H.N., Cui, Y.M., Hu, Q.H., Zhou, R.: 2010b, *Astrophys. J.* **710**, 869.