# Machine and deep learning in solar flare forecasting

Sabrina Guastavino[1], Michele Piana[1,2], Federico Benvenuto[1], Anna Maria Massone[1,2], Francesco Marchetti[3] and Cristina Campi[1]

[1]MIDA, Dipartimento di Matematica, Università di Genova, Genova, Italy, [2]CNR-SPIN, Genova, Italy, [3]Dipartimento di Matematica, Università di Padova, Padova, Italy

www.unige.it

methods for image and data analysis
http://mida.dima.unige.it

## Introduction

Solar flare forecasting can be realized by means of either the numerical solution of model equations or the analysis of magnetic data by means of artificial intelligence techniques. Here we consider this second approach and discuss how neural network optimization and forecasting assessment by means of skill scores are deeply intertwined. Applications will be concerned with time series of magnetograms recorded by SDO/HMI.



The active region is localized

Feature extraction

Will solar flare be originated?

YES

NO

## Feature extraction

Solar flares originate from magnetically active regions (ARs) but not all solar ARs give rise to a flare. Flare forecasting may rely on features extracted from magnetogram images of Active Regions (ARs) e.g. the ones recorded by the Helioseismic and Magnetic Imager (HMI) on-board Solar Dynamics Observatory (SDO). Two approaches:

► features, representing numerical approximation of physical parameters, can be previously extracted from HMI images and then some machine learning techniques are used for the prediction (as SVM, RF, Hybrid Lasso..) [1];
► features can be automatically extracted by some deep learning methods as Convolutional Neural Networks (CNNs) [2].

We focus on the second approach and, to our knowledge, we use for the first time a deep neural network suitable for video classification to predict the occurrence of solar flares from time series of HMI magnetogram images.

## References

[1] Campi, C. et al. (2019), ApJ, [2] Li, X. et al. (2020), ApJ, [3] Marchetti, F. et al. (2021), AirXiv:2103.15522.

## Data

In this work, we consider SDO/HMI images recorded in the time range between 2012 September 14 and 2017 September 30.

► For each AR, we consider the HMI magnetogram frames associated to it and we collect them in 24 hour long time series.
► We have a collection of samples, each one represented by a video of HMI magnetogram frames associated to an AR.
► Data from the past: we know if a flare occurred w.r.t. each time series, therefore we can label each time series with 0 if no flare occurred and with 1 if flares occurred.

## Deep neural network model

We consider a Long-term Recurrent Convolutional Network (LRCN), which combines a Convolutional Neural Network (CNN) to extract sequences of features from HMI images and a Long Short-Term Memory (LSTM) network to analyse the temporal aspect of such sequences (the architecture is shown in Figure 1).



Figure: Deep neural network model. LRCN takes in input time series of HMI images: the CNN extracts features from each HMI image and the LSTM analyses the time series of the extracted features. The output is the probability that a flare occurs in the next 24 hours after the time of the last frame of the time series.

## Gap between loss minimization and score maximization

Two crucial points in machine learning:

► choice of the loss function to be minimized in the training phase;
► choice of the skill scores to evaluate the predictions.

**Idea**: define a loss function which maximizes the desired skill score.
**Ingredients**:

► define a probabilistic version of the confusion matrix;
► define scores over the entries of the probabilistic matrix.

These lead to the definition of Score-oriented Loss (SOL) functions [3].
**Advantage**: no need of a posteriori optimization of the desired skill score.

## Results

We construct 10 different triple of training, validation and test sets in order to assess the statistical robustness in results so that

► AR separation in training, validation and test sets is respected
► the training, validation and test sets are reliable and they are uniform w.r.t. the variety of samples.

In Figure 2 we report the True Skill Statistic (TSS) computed on the 10 test sets for the prediction of solar flares of class C and above (C+) and of class M and above (M+) (blue boxplots). The green boxplots represent the TSS computed on the 10 test sets by excluding the 0-labelled samples corresponding to ARs which originate C+ flares.