# machine learning and flare forecasting

michele piana

the MIDA group

dipartimento di matematica, università di genova

CNR – SPIN, genova

an AI prelude

# data at the core

*"information consumes the attention of its recipients. hence a wealth of information creates a poverty of attention and a need to allocate that attention efficiently among the overabundance of information sources that might consume it"*

(herbert simon, nobel prize for economy)

from big data to
- rich data
- meaningful data
- understood data
- interpreted data

focus on computation
- data simulation
- data analysis
  - inverse problems
  - machine learning

**all this is artificial intelligence**

# data simulation

at disposal:
- a mathematical model mimicking the data formation process
- a set of input parameters for the model
- a numerical method for the solution of the model equations
- objective to accomplish: the set of simulated data


example: simulation of flaring emission
- model: MHD equations + standard model + bremsstrahlung equation
- input parameters: properties of the propagation medium
- numerical method: FEM, FDM, BEM,…
- objective to accomplish: the evolution of the flaring emission along time and spectral energy

# data analysis – inverse problems

at disposal:
- a mathematical model mimicking the data formation process
- a set of experimental measurements
- a numerical method for the solution of the inverse problem
- a statistical model to exploit for formulating the inversion method
- objective to determine: input parameters in the model

example: *"most people, if you describe a train of events to them, will tell you what the result would be...there are few people, however, who, if you tell them a result, would be able to evolve from their inner consciousness what the steps where which led up to that result. this power is what i mean when i talk of reasoning backwards, or analytically... there are fifty who can reason synthetically for one who can reason analytically..."*

(sherlock holmes in 'a study in scarlet')

# data analysis – machine learning

at disposal:
- a historical set of physical parameters (features) with corresponding labels describing the occurrence of a specific condition
- a set of un-labelled incoming features
- a numerical method able to generalize
- objective to determine: the set of labels associated to the set of incoming features

example: flare forecasting:
- historical data: a set of feature vectors extracted from AR magnetic images by means of pattern recognition + X-ray data stating flare occurrence and corresponding class
- incoming data: set of images of a new AR
- objective to determine: probability of occurrence of a flare generated by the new AR and corresponding class

# simulation vs analysis: a math perspective

simulation:
- well-posedness: stable and unambiguous problems
- crucial issues:
  - approximation accuracy
  - computational burden

analysis
- ill-posedness: unstable and ambiguous problems
- crucial issues:
  - how to restore uniqueness and stability
  - reconstruction/generalization accuracy

# a tentative general scheme

$A: X \to Y$      map mimicking the image formation process

$f \in X$      input parameters

$g \in Y$      experimental data set

simulation:    $A(f)$

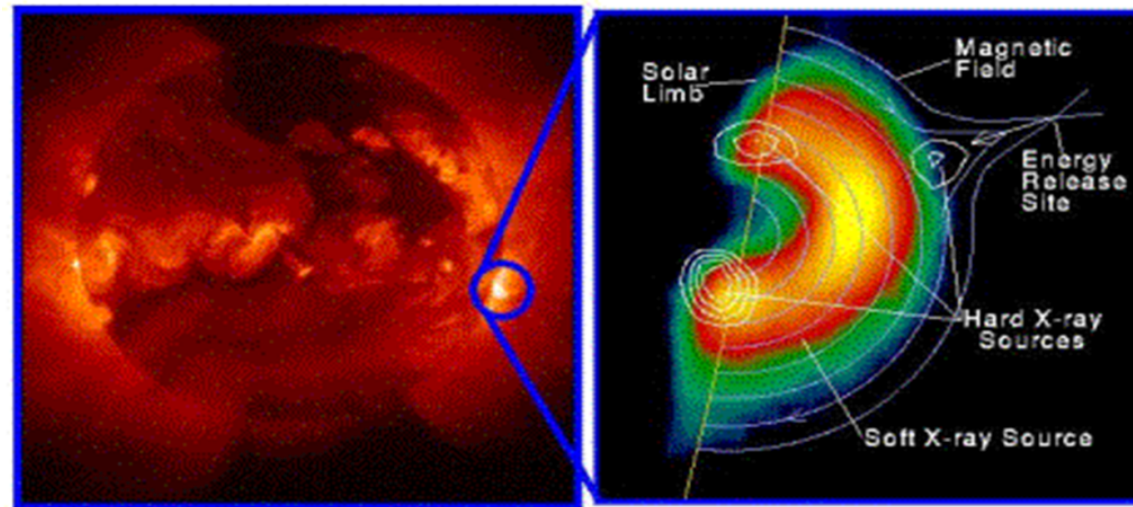inverse problems and machine learning:    $V(f,g) + \lambda \|B(f)\|_p^p = minimum$

- inverse problems: $V(f,g)$ measures how much accurately the candidate solution fits the experimental mesurements through the model; $\|B(f)\|_p^p$ realizes stability

- machine learning: $V(f,g)$ measures how much accurately the candidate predictor would reproduce the label of the historical set; $\|B(f)\|_p^p$ realizes generalization

# the flare problem: a machine learning perspective

- flare forecasting
  - data: SDO/HMI data of active regions (ARs); GOES flare observations and classification (for labelling)
  - unknowns: binary prediction with corresponding flare class
  - method: regularization networks

- flaring source reconstruction
  - data: hard X-ray visibilities measured by STIX in solar orbiter
  - unknowns: shape and physics parameters of the hard X-ray source
  - method: deep neural networks

a physics prelude
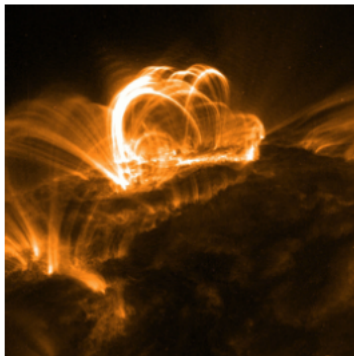
# solar flares: phenomenology



Yohkoh X-ray Image of a Solar Flare, Combined Image in Soft X-rays (left) and Soft X-rays with Hard X-ray Contours (right). Jan 13, 1992.

- generate from ARs
- extend over 10,000 kilometers
- release more than $10^{32}$ ergs in 10-100 seconds
- accelerate billion tons of material to more than a million km per hour
- produce electromagnetic radiation at all wavelengths
- **are the main trigger of space weather (connections with CMEs, SEPs, solar wind)**

# the flare paradox



- inductance: $10^{-6}$ henry
- voltage: 220 V
- light-up time (estimated): $10^{-9}$ s
- li ght-up time (observed): instantaneous



- inductance: 10 henry
- voltage: $10^6$ V
- light-up time (estimated): $3 \times 10^5$ years
- light-up time (observed): minutes

# flare-related data

- vector magnetograms:
    - information on ARs and their productivity
    - SDO/HMI  (looking ahead: PSI in solar orbiter)


- EUV maps:
    - flare morphology
    - SDO/AIA (looking ahead: EUI in solar orbiter)


- hard X-ray visibilities:
    - acceleration mechanisms
    - RHESSI (looking ahead: STIX in solar orbiter)

flare forecasting

# the data

point-in-time SDO/HMI images:
- time range: 09/14/2012 – 04/30/2016
- four issuing times: 00:00 UT – 06:00 UT – 12:00 UT – 18:00 UT
- cadence: 24 hours

features (for each AR):
- 171 features identified in each active region:
  - 167 extracted with a specific pattern recognition algorithms
  - longitude and latitude of the AR
  - binary label encoding the presence of a flare in the past
  - flare class (if occured)
- overall 4442 sets of 171-dimension feature vectors (one AR may last for more than one HMI image)

# training set and test set

we consider supervised learning methods: we need to construct a labeled training set for each issuing time:

1. 66% active regions (ARs) are randomly extracted from the overall set of ARs
2. feature vectors (FVs) associated to each AR are labeled by annotating whether a flare with class at least C1 occured in the next 24 hours
3. the labelling process is performed by using GOES data
4. the set of remaining FVs is not labeled and is used as test set for experiments

# the problem

given a set of 171 features extracted from an AR in the test set, we want to:

1. predict whether an at least C1 flare occurred in the next 24 hours
2. determine which features among the 171 ones mostly impacted the prediction (i.e., compute the weights with which the features contributed to the prediction task and rank them)

# the algorithms

- hybrid LASSO
- hybrid logit
- support vector machine for classification
- random forest

**the routines for the four methods (and for many more) are available at flarecast.eu**

# hybrid LASSO – first step

- X is an NxF matrix with N=4442, F=171:
  - each row contains a feature vector
  - X is the training set
- y is an Nx1 vector made of binary labels
- $\beta$ is an Fx1 vector made of feature weights

compute:

1. $\hat{\beta} = \text{argmin}_\beta(\|y - X\beta\|_2^2 + \lambda\|\beta\|_1)$
2. $\hat{y} = X\hat{\beta}$

# hybrid LASSO – second step

3. apply an unsupervised clustering algorithm to $\hat{y} = X\hat{\beta}$: the outcome is a partition of $\hat{y}$ in two classes (which corresponds to determine a data-adaptive threshold)

4. when a new feature vector x arrives compute the number $x^t\hat{\beta}$ and and assign it to the closest class

retrieved information:
- flare prediction
- set of feature weights $\hat{\beta}$ computed against the training set

# flare prediction: outcome

- a real number which is a probability measure for the (GOES class labelled) flare occurrence

- a binary prediction based on the probability measure

- some skill scores explaining the reliability of the prediction

# flare prediction: assessment

skill scores against the test set:

$$TSS = \frac{TP}{TP+FN} - \frac{FP}{FP+TN}$$ (true skill statistic)

$$HSS = \frac{2 \cdot (TP \cdot TN - FN \cdot FP)}{(TP+FN) \cdot (FN+TN) + (TP+FP) \cdot (FP+TN)}$$ (heidke skill score)

$$ACC = \frac{TP+TN}{TP+TN+FP+FN}$$ (accuracy)  $\qquad$  $$FAR = \frac{FP}{TP+FP}$$ (false alarm ratio)

$$POD = \frac{TP}{TP+FN}$$ (probability of detection)

# results: about training and scores

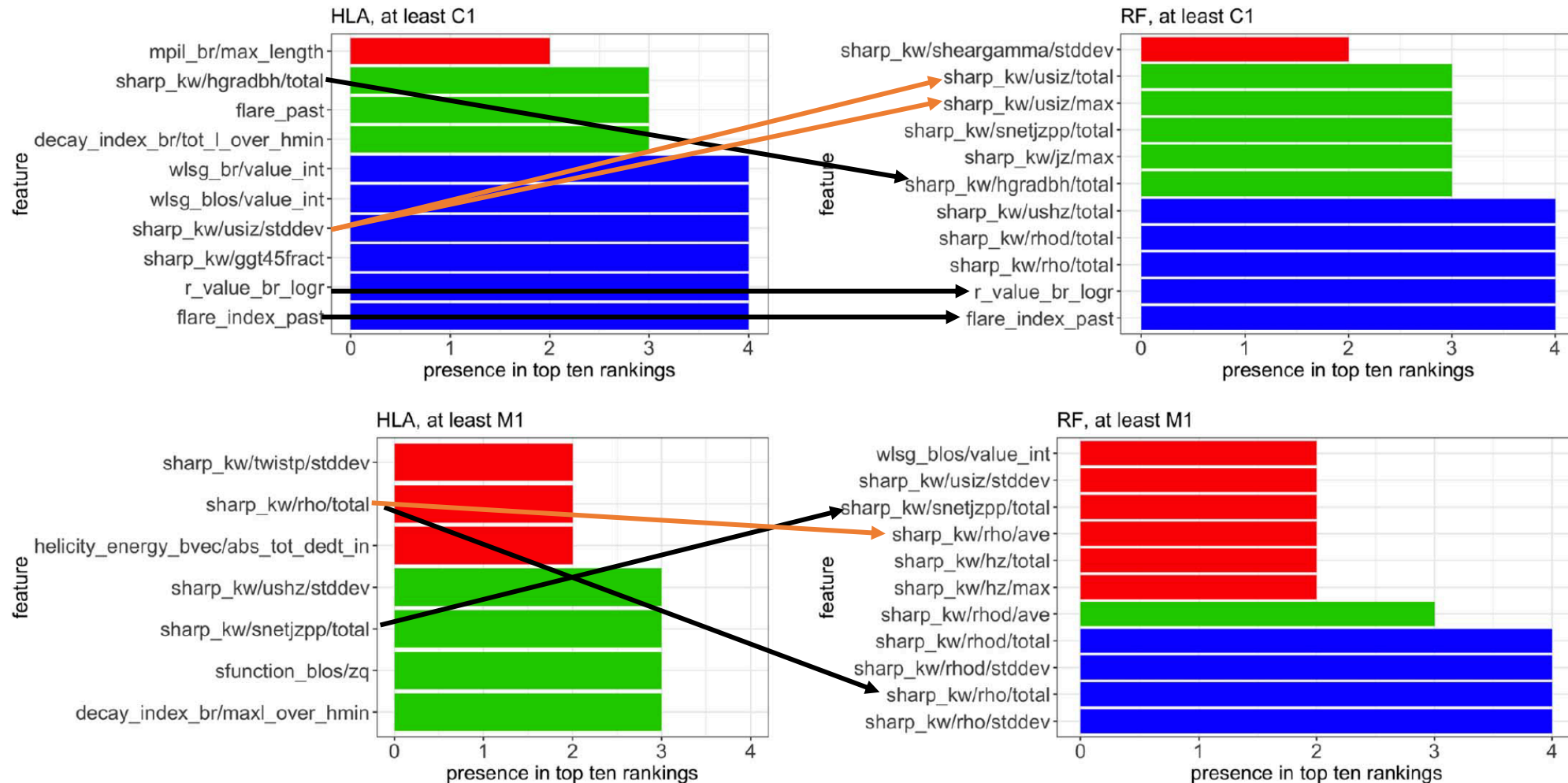| | Test Set-C1+ | Test Set-C1+ | Test Set-M1+ | Test Set-M1+ |
|---|---|---|---|---|
| 00:00:00UT | TSS | HSS | TSS | HSS |
| HLA | $0.48 \pm 0.06$ | $0.51 \pm 0.05$ | $0.56 \pm 0.14$ | $0.27 \pm 0.06$ |
| RF | $0.53 \pm 0.05$ | $0.52 \pm 0.04$ | $0.48 \pm 0.14$ | $0.33 \pm 0.09$ |
| 06:00:00UT | TSS | HSS | TSS | HSS |
| HLA | $0.53 \pm 0.03$ | $0.54 \pm 0.03$ | $0.67 \pm 0.05$ | $0.35 \pm 0.04$ |
| RF | $0.54 \pm 0.03$ | $0.54 \pm 0.03$ | $0.49 \pm 0.08$ | $0.42 \pm 0.06$ |
| 12:00:00UT | TSS | HSS | TSS | HSS |
| HLA | $0.51 \pm 0.04$ | $0.54 \pm 0.03$ | $0.66 \pm 0.06$ | $0.38 \pm 0.04$ |
| RF | $0.53 \pm 0.03$ | $0.53 \pm 0.03$ | $0.51 \pm 0.09$ | $0.43 \pm 0.06$ |
| 18:00:00UT | TSS | HSS | TSS | HSS |
| HLA | $0.54 \pm 0.04$ | $0.55 \pm 0.03$ | $0.64 \pm 0.07$ | $0.39 \pm 0.04$ |
| RF | $0.55 \pm 0.03$ | $0.55 \pm 0.03$ | $0.53 \pm 0.09$ | $0.43 \pm 0.06$ |

training according to active regions

| | Test Set-C1+ | Test Set C1+ | Test Set-M1+ | Test Set-M1+ |
|---|---|---|---|---|
| | TSS | HSS | TSS | HSS |
| HLA | $0.58 \pm 0.01$ | $0.51 \pm 0.01$ | $0.70 \pm 0.02$ | $0.31 \pm 0.03$ |
| RF | $0.61 \pm 0.01$ | $0.56 \pm 0.02$ | $0.71 \pm 0.03$ | $0.39 \pm 0.02$ |
| Florios et al. (2018) | $0.60 \pm 0.01$ | $0.59 \pm 0.01$ | $0.74 \pm 0.02$ | $0.49 \pm 0.01$ |
| Bobra & Couvidat (2015) | ... | ... | $0.76 \pm 0.04$ | $0.52 \pm 0.04$ |

training according to features

# results: top-ten rankings



number of times each feature is selected in the top-10 rankings,
on average over 100 random realizations of the test set, for all issuing times
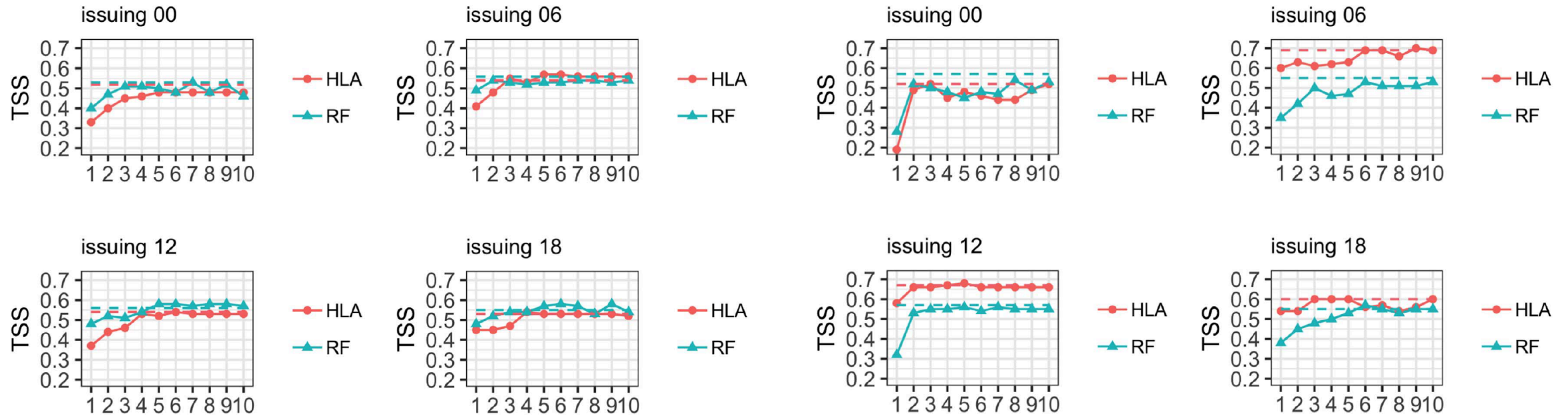
# feature ranking: results - 3

**issuing time: 12:00:00**

| | at least C1 flares | | | | | |
|---|---|---|---|---|---|---|
| | Hybrid Lasso | Hybrid Logit | SVC | Random Forest | average | std |
| flare_index_past | 13,98 | 28,84 | 19,91 | 3,51 | 16,56 | 10,63 |
| sharp_kw/hgradbh/total | 3,47 | 37 | 18,59 | 16,57 | 18,95 | 13,87 |
| wlsg_br/value_int | 3,74 | 14,43 | 22,86 | 43,14 | 21,04 | 16,68 |
| sharp_kw/jz/max | 26,05 | 28 | 16,94 | 18,58 | 22,27 | 5,29 |
| sharp_kw/usiz/max | 24,2 | 36,75 | 34,79 | 18,37 | 28,53 | 8,73 |
| wlsg_blos/value_int | 3,96 | 45 | 46,61 | 25,39 | 30,24 | 20,00 |
| r_value_br_logr | 3,52 | 2,81 | 128,91 | 7,47 | 35,68 | 62,19 |
| sharp_kw/ggt45fract | 14,99 | 32 | 57,6 | 49,3 | 38,35 | 19,00 |
| sharp_kw/usiz/stddev | 17,15 | 49,46 | 54,26 | 45,89 | 41,69 | 16,72 |
| sharp_kw/gamma/total | 61,65 | 20,76 | 52,95 | 34,67 | 42,51 | 18,35 |

**issuing time: 00:00:00**

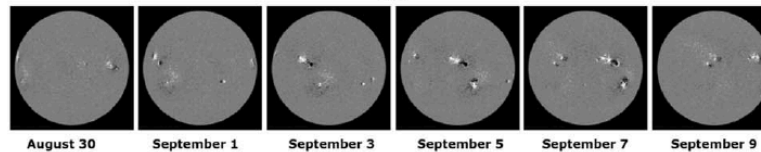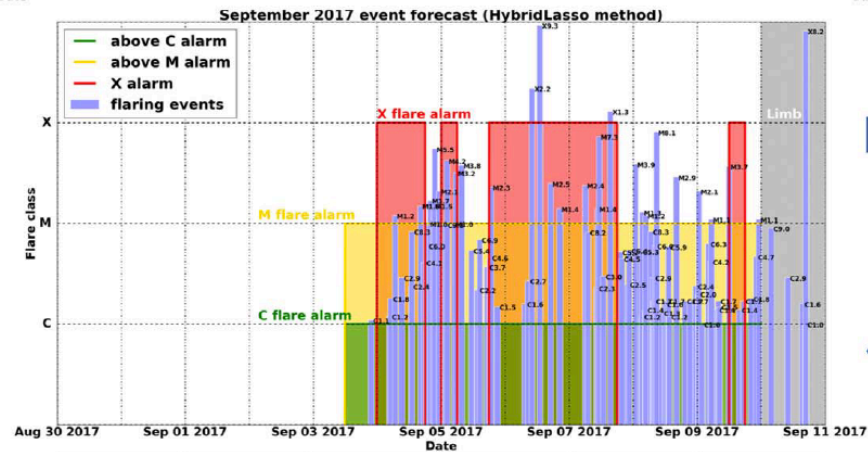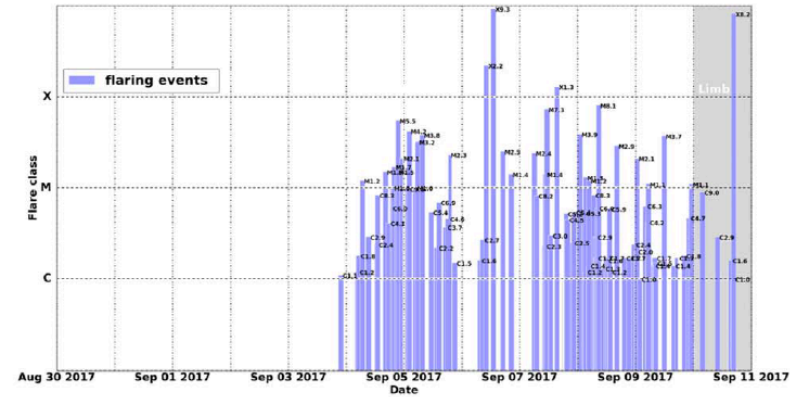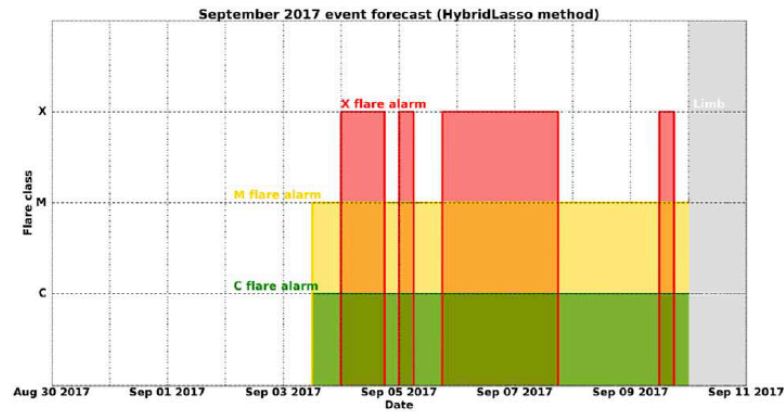| | at least C1 flares | | | | | |
|---|---|---|---|---|---|---|
| | Hybrid Lasso | Hybrid Logit | SVC | Random Forest | average | std |
| wlsg_br/value_int | 3,23 | 5 | 23,95 | 30,02 | 15,55 | 13,45 |
| flare_index_past | 13,89 | 31,13 | 33,92 | 5,47 | 21,10 | 13,68 |
| sharp_kw/usiz/total | 36,44 | 15,8 | 11,39 | 26,1 | 22,43 | 11,19 |
| sharp_kw/hgradbh/total | 29,06 | 7,55 | 24,87 | 39,45 | 25,23 | 13,29 |
| sharp_kw/ggt45fract | 5,25 | 17,14 | 50,68 | 28,25 | 25,33 | 19,33 |
| ising_energy_br/ising_energy | 24,14 | 19,67 | 26,49 | 55,1 | 31,35 | 16,08 |
| wlsg_blos/value_int | 12,83 | 35,08 | 41,12 | 51,11 | 35,04 | 16,21 |
| r_value_br_logr | 6,52 | 4,13 | 126,06 | 16,87 | 38,40 | 58,70 |
| sharp_kw/usiz/stddev | 4,67 | 22,41 | 69,63 | 57,61 | 38,58 | 30,21 |
| sharp_kw/usiz/max | 31,77 | 44,87 | 80,57 | 19,66 | 44,22 | 26,33 |

# results: redundancy of information



TSS scores obtained by using just the 10 top-ten features added one at a time

# machine learning as a warning machine
# forecasting of the september 2017 flaring storm

# some references

- poggio t and girosi f 1990 networks for approximation and learning, *proc* IEEE 78 (9), 1481-1497
- bertero m, poggio t and torr v 1988 ill-posed problems in early vision, *proc* IEEE 76 869
- bertero m 1989 linear inverse and ill-posed problems *adv electron electr phys* 75 1
- rosasco l, de vito e, caponnetto a, piana m and verri a 2004 are loss functions all the same? *neural comput* 16 1063

- benvenuto f, campi c, massone a m and piana m 2020 machine Learning as a flaring storm warning machine: was a warning machine for the 2017 september solar flaring storm possible? *ApJL* 904 L7
- campi c et al 2019 feature ranking of active region source properties in solar flare forecasting and the uncompromised stochasticity of flare occurrence *ApJ* 883 150
- benvenuto f, piana m, campi c and massone a m 2018 a hybrid supervised/unsupervised machine learning approach to solar flare prediction *ApJ* 853 90

- giordano s, pinamonti n, piana m and massone a m 2015 the process of data formation for the spectrometer/telescope for imaging X-rays (STIX) in solar orbiter *SIAM J Imag Sci* 8 1315