

Improving the performance of Machine Learning models with features selection techniques in the context of signal/background discrimination for the AGILE telescope

Leonardo Baroncelli (PhD student in Data Science and Computation), Andrea Bulgarelli, Nicolò Parmiggiani (INAF/OAS Bologna)

leonardo.baroncelli@inaf.it

Deep Learning @ Inaf, F2F meeting, 16-19 September 2019



The task: signal/background discrimination

AGILE (Astrorivelatore Gamma ad Immagini LEggero) is a scientific mission of the Italian Space Agency (ASI), launched on April 23, 2007 for Gamma Ray Astrophysics.

The AGILE-GRID data acquisition system is composed of three main detectors:

- a Tungsten-Silicon Tracker designed to detect and image photons in the 30 MeV-50 GeV energy band,
- a Mini-Calorimeter that detects gamma-rays and charged particles energy deposits between 300 keV and 100 MeV,
- an anti-coincidence (AC) system that surrounds the Silicon Tracker and the Mini-Calorimeter.

On-ground signal-to-background discrimination: **FM3.199** filter.

The FM3.119 uses AdaBoost and a **subset** of discriminant variables (57 features) selected by a domain expert (manual feature selection).

 \Rightarrow Will automatic features selection techniques bring to better results?

⇒ Which is the best features selection method for this task? How many features?





The dataset

- It has been generated from Monte Carlo simulations of the AGILE on-fly data acquisition systems.
- It describes the particle interactions with the AGILE instruments (silicon tracker, calorimeter, anti-coincidence system).
- It is a supervised dataset: the particles are divided into two classes: gamma photons and background particles.
- It is composed by 169.813 rows (particles) and 260 columns (interaction features) in csv format.





Dataset preprocessing

Random (uniform) subsampling of the gamma examples to obtain a balanced dataset.



From 169.185 to 67.952 examples.

Dropped ~60% of the dataset.





Dataset exploration

For each feature:

- Counts histogram
- Density histogram
- Counts histogram in log scale
- Box plot



NCLUSTX density estimation





NCLUSTX outliers detection







Dataset splitting

The dataset has been random splitted in three parts.





Dimensionality reduction

Why?

- Data storage is reduced.
- Machine learning models are simplified which can lead to increased generalization capability.
- Computational complexity for training and testing machine learning models is reduced.

Dimensionality reduction techniques are generally divided into two categories:

- Features extraction ⇒ linear or non-linear projection of data into a lower-dimensional subspace.
- **Features selection** ⇒ subset of the original features based on some performance criterion.



Dimensionality reduction workflow



..and many more techniques:

Bachu, Venkatesh & Anuradha, J.. (2019). A Review of Feature Selection and Its Methods. Cybernetics and Information Technologies. 19. 3. 10.2478/cait-2019-0001.



Dimensionality reduction workflow





Dimensionality reduction workflow



..and many more techniques:

Bachu, Venkatesh & Anuradha, J.. (2019). A Review of Feature Selection and Its Methods. Cybernetics and Information Technologies. 19. 3. 10.2478/cait-2019-0001.



Automatically selected features vs FM3.119 features

First 57 most discriminant features		
Method	Intersection with FM3.119	
PEARSON CORR. COEFF.	26.31%	
MUTUAL INFORMATION	38.6 %	
LASSO	35.09 %	
EXTRA TREES	33.33 %	





Training of the ML algorithms



Features	subsets	Algorithm	Hyper-parameters space
Features selection method	Number of features	Random Forest	'criterion': ['gini', 'entropy'], 'n_estimators': [900, 1100, 1300, 1500, 1700],
FM3.119	5/		'max_depth': [15, 20, 25, 30, 35], 'max_features': ['sqrt', 'log2']
ALL features	241		⇒ 100 mode
Pearson	30, 60 ,120, 180	Gradient Boosting	'learning_rate': [1, 0.1],
Mutual Information	30, 60 ,120, 180		'n_estimators': [300 ,500, 700, 900, 1100], 'max_depth': [5, 10, 15, 20, 25], 'max_features': ['sqrt', 'log2']
Lasso	30, 60 ,120		
Extra Trees	30, 60 ,120, 180		⇒ 100 models
⇒	17 features subsets	Total models t	to train = 17 * (100+100) = 3.400

























Choosing the final model

Algorithm	Fs method	Number of features	F1 score ↓	ROCAUC
Gradient Boosting	ALL features	239	0.909491	0.911164
Gradient Boosting	Mutual Information	180	0.908302	0.910746
Gradient Boosting	Lasso	120	0.907402	0.909958

Algorithm	Fs method	Number of features	F1 score ↓	ROCAUC
Random Forest	Mutual Information	180	0.881477	0.885892
Random Forest	Extra Trees	180	0.880734	0.885109
Random Forest	Lasso	120	0.879837	0.883844
				27

$\frac{\widehat{H}}{2}$ BG rejection vs Signal efficiency

Conclusions

- A multi-collinearity analysis is required to completely remove the correlation among the features.
- Features selection techniques MI, Lasso and Extra Trees selected a subset of features that generated a better model with respect to manual selection.
- The best performances were NOT achieved using a subset of automatic selected features but instead they were achieved training the Gradient Boosting model using all the features (Gradient Boosting performs features selection implicitly).
- With ~50% of the total features (120) selected by Lasso, it has been possible to obtain performance really close to the best results. Data storage requirement is reduced.
- With a ~12% of the total features (30) selected by Extra Trees, it has been possible to obtain a quite good 0.85 F1 score and potentially an easy interpretable model (I need the domain expert, now!).

Thank you for your attention! Any questions?

EXTRA: CVGridSearch mt scalability

Feature selection hyper-parameters space ("mi", 10)

Random forest hyper-parameters space

CPU(s)	192
Thread(s) per core	8
Core(s) per socket	6
Socket(s)	4

Elapsed time (min) vs. Core per CPU

 \Rightarrow 96 will be used to run the Grid Search