

The Higgs boson Machine Learning Challenge

Emanuele Nevali, Matteo Suez, Leonardo Bruno Trentini
Machine Learning and Optimization Laboratory, EPFL, Switzerland

Abstract—In this project, carried out as part of the 2022 Machine Learning course at EPFL, we applied machine learning techniques to actual CERN particle accelerator data to recreate the process of “discovering” the Higgs particle.

I. INTRODUCTION

Rarely, collisions between protons at high speed can produce a Higgs boson, an elementary and unstable particle in the Standard Model of particle physics theory[1]. Since this boson decays rapidly into other particles, it is not possible to observe it directly, but rather measure its “decay signature”, or the products that result from its decay process. The aim of this project is to estimate the likelihood that a given event’s signature was the result of a Higgs boson (signal) or some other process/particle (background). After analysis and preprocessing of the dataset, we introduce the model built with the 6 demanded algorithms and we explain how we managed to choose the best hyperparameters. In conclusion, we present our results.

II. DATA ANALYSIS AND PREPROCESSING

A. Exploratory data analysis

Our dataset contains 250000 points for training and 568238 for testing with their corresponding binary label (“s” for “signal” and “b” for “background”, which have to be predicted for the test set values). All of these points are characterized by 30 different features.

B. Categorization

First of all, we notice that there is only one categorical feature, i.e. *Pri jet num*: it shows the number of jets (see [2] for a physical explanation), and it is an integer that ranges from 0 to 3.

We noticed that some features are meaningless for some values of jets, therefore we have split the dataset into 4 subclasses, each one characterized by a different *Pri jet num*.

C. Missing Values

Thanks to the documentation, we noticed that each “-999” value in the dataset consists in a missing value: we decided to replace all of them with the median of the feature. In particular, we chose to use the median instead of the mean because of its nature of robust estimator to outliers.

Furthermore, for subclasses marked by *Pri jet num* 0 and 1 we decided not to consider columns with too many missing values, in order to avoid singular matrices.

D. Absolute transform of Symmetrical Features

For features whose distribution is symmetric around 0 and for which the proportion of categories is not too even, we replaced them with their absolute value.

This increases the gap between the categories, and makes them more detectable by the model.

E. Outliers erasure

After a search for outliers in the data, we set $\alpha = 0.1$ and decided to edge the extreme values of each feature to the α -quantile (for the lower tail) and to the $(1 - \alpha)$ -quantile (for the upper tail).

F. Angles

We noticed that some of the features represents measurements of angles. For this reason, in order to adapt our model to the periodicity, we tried to replace these columns with sine and cosine transformations. However, only the latter appeared to be effective, and it was the only one kept in the final model.

G. Correlated features

We observed that the empirical distribution of some of the features was too similar, index of an high correlation between them. By not considering these features, we slightly improved our performance.

H. Heavy-tailed features

Noticing that many of the features were heavy-tailed, we computed for these a logarithmic transformation: $x_k^* = \log(1 + x_k)$ (where k is the feature column). These transformation helps in reducing or removing the skewness of our original data.

I. Standardization

In order to ensure a better functioning of our algorithms, we decided to standardize the data by subtracting from each feature its mean and dividing by its standard deviation: $Z = \frac{X - \mu}{\sigma}$. The same process was applied both to the train set and to the test set, using on the latter μ and σ of the former.

J. Polynomial Expansion

We managed to improve our predictions by expanding the features. In particular, we augmented features by implementing a polynomial expansion for each column of the form $\phi : [X] \rightarrow [X, X^2, \dots, X^D]$.

The polynomial degree D is one of the hyper-parameters of our model: we found its optimal value by performing a 4-fold cross validation on each subset.

In addition, we decided to add for each couple of features a new feature with their product. This last operation improved significantly the accuracy of the model.

III. METHODS

After the pre-processing on the dataset, we managed to find a model for our problem by implementing six different methods:

- **Gradient Descent**
- **Stochastic Gradient Descent**
- **Least squares** using normal equations
- **Ridge regression** using normal equations
- **Logistic regression** using gradient descent or stochastic gradient descent ($y \in \{0, 1\}$)
- **Regularized logistic regression** using gradient descent or stochastic gradient descent ($y \in \{0, 1\}$, with regularization term $\lambda \|\omega\|^2$)

All functions return: (w, loss), which is the last weight vector of the method, and the corresponding loss value (or cost function).

Optimal parameters λ (regularization to reduce overfitting), γ (step size for gradient and stochastic gradient descent) and D (degrees for polynomial expansion) were chosen through a 4-fold cross validation for each subset.

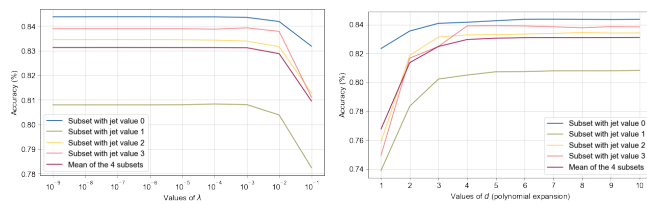


Figure 1: Accuracy for different values of λ and D

IV. RESULTS

Due to the fact that we were facing a binary classification problem, we first focused on the Regularized Logistic Regression. However, in the meanwhile we also carried out multiple tests with Least Squares and Ridge Regression, and remarkably the latter turned out to be the most effective method.

For this reason, we focused on calibrating Ridge Regression

hyperparameters. We did it by using grid-search method and 4-fold cross-validation, looking for the best compromise between underfitting and overfitting. We adjusted the hyperparameters for each of the four subclasses.

Method	Subset 1	Subset 2	Subset 3	Subset 4	Average
Gradient Descent	70.7%	65.4%	67.9%	68.3%	68.2%
Stochastic Gradient Descent	69.8%	66.6%	67.5%	68.5%	68.3%
Least squares	84.6%	81.4%	84.6%	84.2%	83.6%
Ridge regression	84.7%	81.3%	84.7%	84.5%	83.7%
Logistic regression	77.1%	68.6%	67.1%	69.8%	71.8%
Regularized logistic	76.6%	69.1%	66.7%	69.8%	71.6%
Variance	$\pm 10^{-3}$				

Table I: Accuracies for each method and subset. The average was weighted on the number of points for each subset

At the end, we obtained a model that achieved an accuracy of 0.837 on AICrowd platform.

V. CONCLUSIONS

The best performance on the Higgs boson challenge, with an accuracy on the test set of 0.837, was obtained with a Ridge Regression.

It is important to underline that data analysis, data cleaning and feature preprocessing played a key role in this task, enabling us to achieve massive improvements in model's accuracy.

REFERENCES

- [1] M. Goulette, "What should we know about the higgs particle?" *Atlas experiment - CERN*, 2013.
- [2] C. Adam-Bourdariosa, G. Cowanb, C. Germain, I. Guyond, B. Kegl, and D. Rousseau, "Learning to discover: the higgs boson machine learning challenge," *Higgs Challenge*, 2014.