

Tuning of logistic regression for Higgsboson ML challenge

Ettore Fincato, Hannah Sansford, Harry Tata

January 16, 2022

Abstract

We investigate the tuning of a logistic regression for finding evidence of the decay into *fermion* pairs of the Higgs boson, for the Higgs Boson Machine Learning Challenge [2]. With the aim of finding a logistic classification maximising the approximate median significance (AMS) objective function (5), we conduct a tuning based on a two-stage procedure from [3]. This tuning strategy relies on a first stage of fitting the logistic model coefficients, and on a second stage of tuning the classification threshold on the basis of AMS. By showing that AMS values based on the two-stage procedure are significantly larger than the ones based on the standard logistic classification threshold, we deduce that the coefficients from the logistic model are still of use. Logistic regression is a simple model with easily interpretable coefficients; this is convenient in settings where which the variables have a physical interpretation.

1 Introduction

1.1 Background

The Higgs boson can *decay* through various different processes, producing other particles in the process. In physics, one calls a decay into specific particles a *channel*. Until fairly recently, the Higgs boson had been seen only in boson pair decay channels. It is now of importance to seek evidence on the decay into *fermion* pairs, specifically *tau-leptons* or *b-quarks*, and to measure their characteristics [2]. The ATLAS experiment [1] was the first to report evidence of the Higgs boson to tau-tau channel and the goal of this report is to improve on this analysis.

1.2 Overview

The ATLAS experiment at CERN provided simulated data that was used to optimise the analysis of the Higgs boson. In the Large Hadron Collider, proton bunches are accelerated in both directions on a circular trajectory. This results in some of the protons colliding as the bunches cross the ATLAS detector (called an *event*), which produces hundreds of millions of proton-proton collisions per second. The particles resulting from each event are detected by sensors and, from this raw data, certain real-valued features are estimated [2].

Most of the uninteresting events (called the *background*) are discarded using a real-time multi-stage cascade classifier. However, many of the remaining events represent known processes that are also known as *background*. Our aim is to find the region of the feature space in which there is a significant excess of events compared to what known background processes can explain (called *signal*).

Once the region has been fixed, the significance of the excess is determined using a statistical test. If the probability that the excess has been produced by background processes falls below a pre-determined limit, the new particle is deemed to be discovered.

2 Problem Formulation

Let $\mathcal{D} = \{(\mathbf{x}_1, y_1, w_1), \dots, (\mathbf{x}_n, y_n, w_n)\}$ be the training set, where $\mathbf{x}_i \in \mathbb{R}^d$ is a d -dimensional feature vector, $y_i \in \{\text{b}, \text{s}\}$ is the label, and $w_i \in \mathbb{R}^+$ is a non-negative weight. Let $\mathcal{S} = \{i : y_i = \text{s}\}$ and

$\mathcal{B} = \{i : y_i = \text{b}\}$ be the index sets of signal and background events respectively, and let $n_s = |\mathcal{S}|$ and $n_b = |\mathcal{B}|$ be the number of simulated signal and background events.

The simulated dataset also includes importance weights for each event. Since the objective function (5) depends on the *unnormalised sum* of weights, in order to make the setup invariant to the *numbers* of simulated events n_s and n_b , the sum across each set (test/training) and each class (signal/background) is set to be fixed, i.e.

$$\sum_{i \in \mathcal{S}} w_i = N_s \quad \text{and} \quad \sum_{i \in \mathcal{B}} w_i = N_b. \quad (1)$$

These normalisation constants N_s and N_b are simply the *expected total number* of signal and background events, respectively, during the time interval of the data taking. The individual weights are then proportional to the conditional densities,

$$p_s(\mathbf{x}_i) = p(\mathbf{x}_i | y = \text{s}) \quad \text{and} \quad p_b(\mathbf{x}_i) = p(\mathbf{x}_i | y = \text{b}),$$

divided by the instrumental densities $q_s(\mathbf{x}_i)$ and $q_b(\mathbf{x}_i)$, i.e.

$$w_i \propto \begin{cases} p_s(\mathbf{x}_i)/q_s(\mathbf{x}_i), & \text{if } y_i = \text{s}. \\ p_b(\mathbf{x}_i)/q_b(\mathbf{x}_i), & \text{if } y_i = \text{b}. \end{cases} \quad (2)$$

Now, let $g : \mathbb{R}^d \rightarrow \{\text{b}, \text{s}\}$ be a classifier. Let the *selection region* $\mathcal{G} = \{\mathbf{x} : g(\mathbf{x}) = \text{s}\}$ be the set of points classified as signal, and let $\hat{\mathcal{G}}$ denote the *index set* of points that g classifies as signal, i.e.

$$\hat{\mathcal{G}} = \{i : f(\mathbf{x}_i) \in \mathcal{G}\} = \{i : g(\mathbf{x}_i) = \text{s}\}.$$

Then we have that

$$s = \sum_{i \in \mathcal{S} \cap \hat{\mathcal{G}}} w_i \quad (3)$$

is an unbiased estimator of the expected number of signal events selected by g and similarly,

$$b = \sum_{i \in \mathcal{B} \cap \hat{\mathcal{G}}} w_i \quad (4)$$

is an unbiased estimator of the expected number of background events selected by g . Alternatively, s and b are the *unnormalised* true and false positive rates, respectively.

High-energy physicists suggest the use of the *approximate median significance* (AMS) objective function defined by

$$\text{AMS} = \sqrt{2 \left[(s + b + b_{\text{reg}}) \ln \left(1 + \frac{s}{b + b_{\text{reg}}} \right) - s \right]} \quad (5)$$

to optimise the selection region for discovery significance, where b_{reg} is a regularisation term suggested to be set to $b_{\text{reg}} = 10$. Hence, our aim is simply to train a classifier g based on the training data \mathcal{D} with the goal of maximising the AMS on some unseen test data.

The ATLAS experiment dataset (SHOULD WE CITE THE ACTUAL DATASET?) provided for the Higgs Boson Machine Learning Challenge comprises 250000 training samples and 550000 testing samples.

3 Maximising the AMS

Having specified a statistical model for classification, its parameters will be tuned by maximising the AMS. This can be done in two ways.

- **Direct maximisation:** One can consider the AMS as an error/objective measure in a cross validation procedure. So the model will be tuned by a cross validation (cv) procedure by looking at the coefficients which directly maximise the AMS.

- **Two-stage maximisation:** Given a real-valued function f (e.g. a linear $f(\mathbf{w}, \mathbf{x}) = \mathbf{w}'\mathbf{x}$), one can always get a classification model by setting a threshold $\theta \in A \subset \mathbb{R}$ and a classifier

$$h(\theta) := \text{sgn}(f - \theta) \in \{-1, 1\}. \quad (6)$$

One can therefore train and tune f “independently” of the AMS, and then find the threshold θ for which the classifier $h(\theta)$ maximises the AMS. For example, one could train f in a logistic regression (i.e. train f such that it minimises a logistic loss), and then, on a second stage, choose the classification threshold based on AMS. This approach is presented in [3]

We will see that the two-stage maximisation is much more performant (in terms of AMS) than the direct one.

3.1 Two-stage maximisation

While the direct maximisation with cv can be seen as a usual tuning procedure by cross validation, the two-stage maximisation should be justified. In the following, we explain the two stage procedure in more detail for a classifier $h : \mathcal{X} \rightarrow \{-1, 1\}$, as defined in (6), where \mathcal{X} denotes the features space, and f is trained by minimising a logistic loss.

By incorporating the weights to the probabilities, without loss of generality one may write s (3) and b (4) as

$$s = s(h) := \mathbb{P}(h(\mathbf{X}) = 1, Y = 1), \quad b = b(h) := \mathbb{P}(h(\mathbf{X}) = -1, Y = 1)$$

where \mathbf{X} denotes a random input vector and Y a random binary response. [3] justifies the two stage procedure explained above for logistic loss transformation of f , as follows. Given a classifier h obtained from f as in (6), one can define the *AMS regret* as

$$R_{AMS}(h) := \text{AMS}^2(h^*) - \text{AMS}^2(h) \quad (7)$$

where

$$h^* := \underset{h}{\text{argmax}} \text{AMS}^2(h).$$

Similarly, the *logistic regret* of f is defined as

$$R_{log}(f) := L(f) - L(f^*) \quad (8)$$

for the expected logistic loss transformation of f ,

$$L(f) := \mathbb{E} \left[\log \left(1 + e^{-Yf(\mathbf{X})} \right) \right], \quad (9)$$

and

$$f^* := \underset{f}{\text{argmin}} L(f).$$

The formal justification of the idea of

- training a logistic loss transformation of f , and
- finding the threshold θ in (6) by maximising the AMS on a second stage

comes from the following theorem.

Theorem 3.1. *Given a real-valued function f and the related classifier $h(\theta) := \text{sgn}(f - \theta)$, let $\hat{\theta} := \underset{\theta}{\text{argmax}} \text{AMS}(h)$. Then,*

$$R_{AMS}(h(\hat{\theta})) \leq \frac{s(h^*)}{b(h^*)} \sqrt{\frac{1}{2} R_{log}(f)} \quad (10)$$

Proof. See Theorem 2 of [3]. □

That is, the AMS regret of classifier $h(\theta)$ is upper bounded by the logistic regret of the underlying f . It is possible to prove that it is sufficient to optimize θ on the empirical counterpart of AMS calculated on a separate validation sample.

Remark 3.2. *Training f by minimising a logistic loss is particularly suited for a classification task. For example, it allows us to train (on the first stage) models such as a logistic regression.*

Remark 3.3. *As far as logistic regression is concerned, recall that after model's coefficients \mathbf{w} have been estimated, classification for a new observation y_0 can be done by checking whether $\sigma(\hat{\mathbf{w}}' \mathbf{x}_0) > \theta$ for some suitable $\theta \in \mathbb{R}$. In a 'proper' logistic regression, θ is directly set to 0.5. Here we choose it during a second stage via maximising AMS.*

Our aim is to tune a logistic regression. The theoretical set-up of logistic regression can be found below.

3.2 Logistic regression

Suppose $y \in \{-1, 1\}$. Note that

$$p(y = 1|\mathbf{x}) = \frac{p(\mathbf{x}|y = 1)p(y = 1)}{p(\mathbf{x}|y = 1)p(y = 1) + p(\mathbf{x}|y = -1)p(y = -1)} = \frac{1}{1 + \frac{p(\mathbf{x}|y=-1)p(y=-1)}{p(\mathbf{x}|y=1)p(y=1)}} \quad (11)$$

Logistic regression is a discriminative learning model which aims at minimising an expected loss by predicting

$$\hat{y} = \underset{y_0}{\operatorname{argmin}} \mathbb{E}_{p(y|\mathbf{x})}(L(y, y_0)|\mathbf{x})$$

and by *inferring* $p(y|\mathbf{x})$ *directly*.

Specifically, it models the log density ratio as

$$\log \left(\frac{p(\mathbf{x}|y = -1)p(y = -1)}{p(\mathbf{x}|y = 1)p(y = 1)} \right) \stackrel{\text{model}}{=} f(\mathbf{x}, \boldsymbol{\beta}) := \langle \boldsymbol{\beta}, \mathbf{x} \rangle + \beta_0$$

i.e. models the *log-odds* by a linear model (the term on the left of the proportionality symbol is added for intuition). In this way we can write

$$p(y|\mathbf{x}; \boldsymbol{\beta}) = \sigma(f(\mathbf{x}, \boldsymbol{\beta})y) \quad (12)$$

for a so-called 'activation/logistic function' $\sigma(t) := \frac{1}{1 + \exp(-t)}$.

ML estimates of $\boldsymbol{\beta}$ are computed from the log-likelihood as

$$\boldsymbol{\beta}_{ML} = \underset{\boldsymbol{\beta}}{\operatorname{argmax}} \sum_{i \in D} \log \sigma(f(\mathbf{x}_i, \boldsymbol{\beta})y_i)$$

Remark 3.4. *Notice that*

$$\begin{aligned} \underset{\boldsymbol{\beta}}{\operatorname{argmax}} \sum_{i \in D} \log \sigma(f(\mathbf{x}_i, \boldsymbol{\beta})y_i) &= \underset{\boldsymbol{\beta}}{\operatorname{argmax}} \sum_{i \in D} -\log(1 + \exp\{-y_i \boldsymbol{\beta}' \mathbf{x}_i\}) \\ &= \underset{\boldsymbol{\beta}}{\operatorname{argmin}} \sum_{i \in D} \log(1 + \exp\{-y_i \boldsymbol{\beta}' \mathbf{x}_i\}) \end{aligned}$$

so that we are consistent with [3] by minimising a log loss as in equation (9) (of course, we are minimising its sample counterpart).

3.2.1 Weighted logistic regression

Remark 3.5. *In this dataset weights have an essential role in rebalancing the likelihood of signals and background events, as data have been artificially simulated. We will therefore fit the model with a weighted logistic regression.*

While unweighted logistic regression maximises

$$\hat{\beta} = \operatorname{argmax}_{\beta} \sum_{i \in D} -\log(1 + \exp\{-y_i \beta' \mathbf{x}_i\}),$$

weighted logistic regression takes into account weights w_i by weighting the corresponding observations, finding

$$\hat{\beta} = \operatorname{argmax}_{\beta} \sum_{i \in D} -\log(1 + \exp\{-w_i y_i \beta' \mathbf{x}_i\}). \quad (13)$$

Remark 3.6. *Elastic net (weighted) logistic regression, where*

$$\hat{\beta} = \operatorname{argmax}_{\beta} \left[\left(\sum_{i \in D} -\log(1 + \exp\{-w_i y_i \beta' \mathbf{x}_i\}) \right) + \lambda \left[\frac{1-\alpha}{2} \|\beta\|_2^2 + \alpha \|\beta\|_1 \right] \right],$$

penalises the L1 and L2 norms of the coefficients. This was implemented for both the direct and two-stage approaches with α, λ selected by cross validation. In both cases $\lambda = 0$ – the unpenalised model – performed best. Here the training dataset (250000 observations) was very large compared to the number of features (30); with higher-dimensional models an elastic net regression may be beneficial.

3.3 Generalising the logistic regression decision function

After ML estimates $\hat{\beta}$ have been computed, classification of an observation \mathbf{x}_0 can be based on which probability is larger out of $p(y = 1|\mathbf{x}_0; \beta)$ and $p(y = -1|\mathbf{x}_0; \beta)$. That is,

$$\begin{aligned} \hat{y}_{\mathbf{x}_0} &= 1 \text{ if } \sigma(f(\mathbf{x}_0, \hat{\beta})) > \sigma(-f(\mathbf{x}_0, \hat{\beta})) \\ \text{or, equivalently, } \hat{y}_{\mathbf{x}_0} &= 1 \text{ if } \sigma(f(\mathbf{x}_0, \hat{\beta})) > 0.5 \end{aligned}$$

On the other hand, we can classify \mathbf{x}_0 to class -1 if $\sigma(f(\mathbf{x}_0, \hat{\beta})) < 0.5$.

Hence, we classify observations based on whether the fitted probability that the class $y = 1$, falls either side of the *threshold* $\theta = 0.5$.

Remark 3.7. *The two-stage procedure ‘generalises’ this classification method by admitting that there is a $\theta \in (0, 1)$, possibly different from 0.5, such that the classification based on*

$$\hat{y}_{\mathbf{x}_0} = 1 \text{ (resp. -1) if } \sigma(f(\mathbf{x}_0, \hat{\beta})) > \theta \text{ (resp. } < \theta).$$

could be ‘better’ (for different evaluation metrics) than the standard logistic regression one.

4 Model fitting

4.1 AMS as direct measure - Cross Validation

As first attempt, a standard (i.e. $\theta = 0.5$) weighted logistic regression (WLR) has been fitted on the train set. Cross validation suggested that we can expect to get an AMS in the range [0.15, 0.41] on a test set. This is a relatively small value compared to results in following sections. Using the entire set of training data to fit the WLR model, we get an AMS value of 0.18 on the test dataset.

4.2 Two-stage maximisation

4.2.1 First stage

The training set $\mathcal{D} = \{(\mathbf{x}_1, y_1, w_1), \dots, (\mathbf{x}_n, y_n, w_n)\}$ has been split into D_t (containing 80% of observations) and a validation set D_v , containing 20% of the observations. A weighted logistic regression has been fitted on D_t , thus obtaining its estimated coefficients $\hat{\beta}$.

Remark 4.1. *Every time we subset the original training set, we must re-calibrate the weights such that the sum of selected background weights be equal to N_b and the sum of selected signal weights be equal to N_s . A specific function has been created for it.*

4.2.2 Tuning of threshold theta

The validation set D_v has been used to select the value of the threshold θ for which the AMS is maximised (on D_v). Uncertainty on AMS and on the best θ have been obtained by a cross validation procedure.

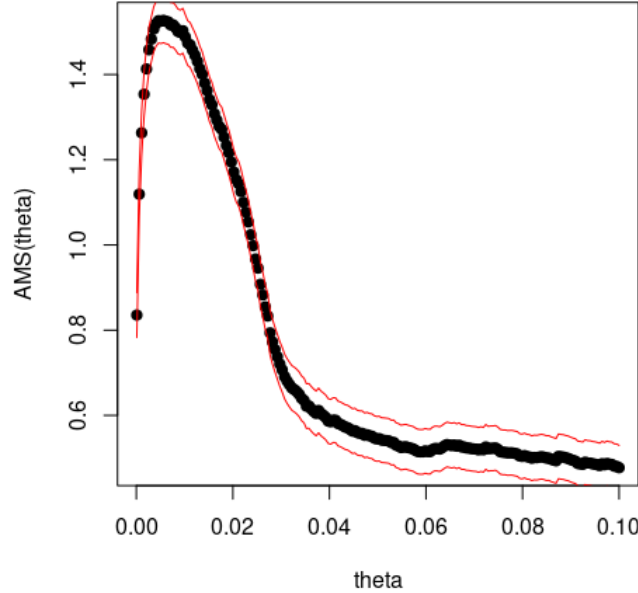


Figure 1: AMS values from CV procedure for different thresholds

Figure 1 shows that AMS is maximised for small values of θ (smaller than 0.02). The maximum value of AMS is around 1.45 - much higher than the one found with a standard logistic regression. A simple Monte Carlo maximisation of $\text{AMS}(\theta)$ finds the optimal threshold $\hat{\theta}$, which is generally found to be between around 0.004. The red lines in Figure 1 represent ± 1 standard deviation from the mean (using cross validation).

The same procedure has been tried after PCA. Analysis of variance explained by PC's suggested to try to use a small number of PC's (see Figure 2).

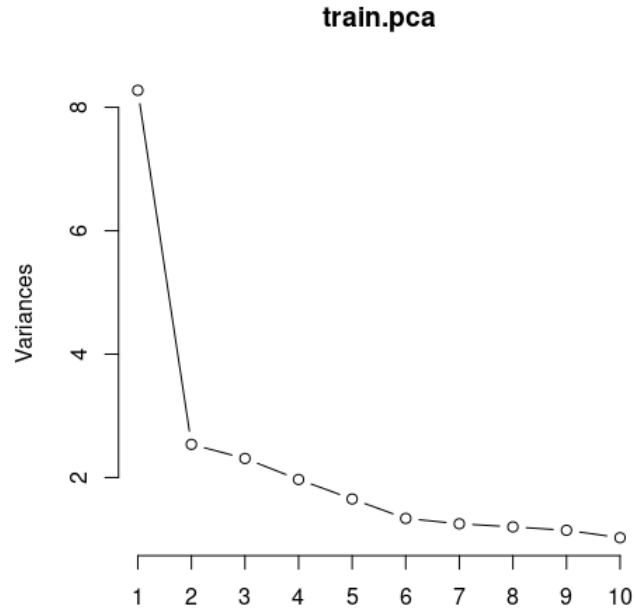


Figure 2: Screeplot

We decided to use the first 3 PCs as a new dataset, and repeat the procedure below. Cross-validated AMS values are plotted in Figure 3.

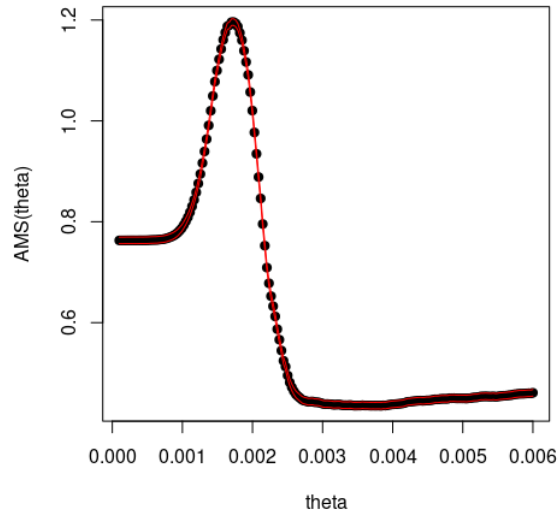


Figure 3: AMS values from CV procedure for different thresholds - on the first three PCAs

One can see that the AMS on validation set on the PCs is much more ‘peaked’. In particular, AMS variance is lower than before, and the range of suitable thetas is much shorter too. This might be a sign of underfitting due to the use of only three PCs. $\hat{\theta}_{pc}$ is almost always 0.0017.

4.3 Results on test set

It has already been said that a standard logistic regression on test set reaches an AMS of about 0.18, which is a very low performance.

On the contrary, the two stage procedure has a much better performance on the test set. With the threshold θ suggested from the second stage procedure on all the features, AMS on the test set is about 1.5, which is a big improvement. With the threshold suggested after PCA, we get an AMS of about 1.37. This is still much better than using a standard logistic regression classification rule, but smaller than for the threshold suggested by the original training set. This, as said before, may be connected to the model being too simplified - the range of suggested theta's was much shorter than before dimensional reduction.

There is further evidence suggesting that the two-stage procedure allows us to pick a threshold which is suited for the test set too. In Figure 4, AMS values from predictions on test set are plotted for different thetas. It is clear that few values of theta lead to a high value of AMS, and among them we find the ones suggested from the second stage of our procedure.

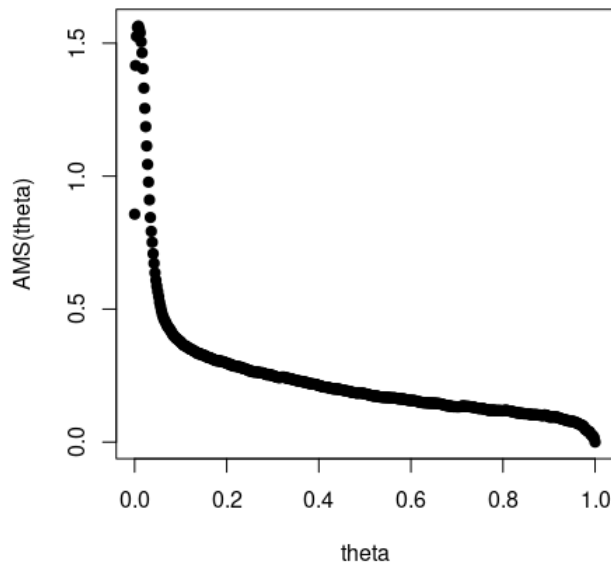


Figure 4: AMS values after predictions on test set against the decision rule threshold theta- for different thetas.

4.4 Discussion

The developments above go in redefining the use of a logistic regression for the Higgs boson decay classification problem for AMS maximisation. The winning solution of the challenge leads, using a complicated black box model with many neural networks, to an AMS of about 3.6, which is higher than our value (1.5). Surely, that winning model is more performant in prediction than ours. However, the scope of our investigation was the optimal tuning of logistic regression. We hope that since the resulting AMS (1.5) after the two-stage procedure was much higher than for a standard classification (0.18), this strategy might then be applied to more advanced and complex models, producing a sharp increase in AMS. An advantage of logistic regression is that the model's coefficients have a clear interpretation in terms of odds-ratio. Moreover our set up is theoretically justified in [3], Theorem 2 (see 3.1). Another advantage is that our model seems not to overfit: AMS values on train and test set are very close.

References

- [1] Evidence for Higgs Boson Decays to the $\tau^+\tau^-$ Final State with the ATLAS Detector. Technical report, CERN, Geneva, Nov 2013. All figures including auxiliary figures are available at <https://atlas.web.cern.ch/Atlas/GROUPS/PHYSICS/CONFNOTES/ATLAS-CONF-2013-108>.
- [2] Claire Adam-Bourdarios, Glen Cowan, Cecile Germain, Isabelle Guyon, Balazs Kegl, and David Rousseau. Learning to discover: the higgs boson machine learning challenge.
- [3] Wojciech Kotlowski. Consistent optimization of ams by logistic loss minimization, 2014.