

Higgs Boson Machine Learning Project

Ettore Fincato, Hannah Sansford, Harry Tata

January 11, 2022

Abstract

Write abstract here.

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 H 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 (LHC), 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,

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

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

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

Now, let $g : \mathbb{R}^d \rightarrow \{b, s\}$ be a classifier. Let the *selection region* $\mathcal{G} = \{\mathbf{x} : g(\mathbf{x}) = s\}$ 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.,

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

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.

3 Maximising the AMS

Having decided a statistical model for classification, its parameters will be tuned by maximising the AMS. This will 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 s.t. it minimises a logistic loss), and then, on a second stage, choose the classification threshold based on AMS. This approach is presented in [3]

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. This section is dedicated to the explanation of the two stage procedure 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 and b from (3) and (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. The paper 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

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

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

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

and

$$f^* := \underset{f}{\operatorname{argmax}} L(f)$$

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 $h(\theta) := \operatorname{sgn}(f - \theta)$, let $\hat{\theta} := \underset{\theta}{\operatorname{argmax}} \text{AMS}(h)$. Then,*

$$R_{AMS}(h(\hat{\theta})) \leq \frac{s(h^*)}{b(h^*)} \sqrt{\frac{1}{2} R_{log}(f)} \quad (9)$$

Proof. See [3]-Theorem 2. □

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 \mathbf{x}_0 can be done by setting $\hat{\mathbf{w}}' \mathbf{x}_0 > \theta$ for some suitable $\theta \in \mathbb{R}$. In a “proper” logistic regression, θ is directly set to 0. Here we choose it during a second stage via maximising AMS.*

3.2 Logistic regression

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

$$p(y = 1|\mathbf{x}) = \frac{p(\mathbf{x}|y = 1)p(y = 1)}{p(\mathbf{x}|y = 1)p(y = 1) + p(\mathbf{x}|y = -1)p(y = -1)} = \frac{1}{1 + \frac{p(\mathbf{x}|y=-1)p(y=-1)}{p(\mathbf{x}|y=1)p(y=1)}} \quad (10)$$

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

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

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

Specifically, it models the log density ratio as

$$\log \left(\frac{p(\mathbf{x}|y = -1)p(y = -1)}{p(\mathbf{x}|y = 1)p(y = 1)} \right) \stackrel{\text{model}}{=} f(\mathbf{x}, \mathbf{w}) := \langle \mathbf{w}, \mathbf{x} \rangle + w_0$$

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

$$p(y|\mathbf{x}; \mathbf{w}) = \sigma(f(\mathbf{x}, \mathbf{w})y) \quad (11)$$

for a so-called "activation/logistic function" $\sigma(t) := \frac{1}{1+\exp(-t)}$.

ML estimates of \mathbf{w} are computed from the log-likelihood as

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

3.3 Generalising logistic regression's decision function

After ML estimates $\hat{\mathbf{w}}$ have been computed, classification of an observation \mathbf{x}_0 can be based on whether $p(y = 1|\mathbf{x}; \mathbf{w})$ is bigger than $p(y = -1|\mathbf{x}; \mathbf{w})$. That is,

$$\begin{aligned} \hat{y}_{\mathbf{x}_0} &= 1 \text{ if } \sigma(f(\mathbf{x}_0, \hat{\mathbf{w}})) > \sigma(-f(\mathbf{x}_0, \hat{\mathbf{w}})) \\ \text{or, equivalently, } \hat{y}_{\mathbf{x}_0} &= 1 \text{ if } \log \left(\frac{\sigma(f(\mathbf{x}_0, \hat{\mathbf{w}}))}{\sigma(-f(\mathbf{x}_0, \hat{\mathbf{w}}))} \right) > 0 \end{aligned}$$

We know that

$$\log \left(\frac{\sigma(f(\mathbf{x}_0, \hat{\mathbf{w}}))}{\sigma(-f(\mathbf{x}_0, \hat{\mathbf{w}}))} \right) \propto \langle \hat{\mathbf{w}}, \mathbf{x}_0 \rangle + \hat{w}_0$$

Therefore, classification from standard logistic regression can also be made as follows

$$\hat{y}_{\mathbf{x}_0} = 1 \text{ (resp. -1) if } \langle \hat{\mathbf{w}}, \mathbf{x}_0 \rangle + \hat{w}_0 > 0 \text{ (resp. } < 0).$$

Remark 3.4. *The two-stage procedure "generalises" this classification method by admitting that there is a $\theta \in \mathbb{R}$, possibly different from 0, such that the classification based on*

$$\hat{y}_{\mathbf{x}_0} = 1 \text{ (resp. -1) if } \langle \hat{\mathbf{w}}, \mathbf{x}_0 \rangle + \hat{w}_0 > \theta \text{ (resp. } < \theta).$$

could be "better" (for some purposes) than the standard logistic regression one.

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.