

## Thoughts on ATLAS MVA Challenge

This note records and extends some thoughts from the discussion at LAL on 29/3/13 (GDC, David Rousseau, Balazs Kegl, Cecile Germain) concerning a Data Challenge, in which participants would be given  $n$ -tuples of signal and background processes and try to construct an analysis with a maximum discovery sensitivity for the signal.

A possible set of rules for the analysis would be that the participant should define a single search region of the space of input variables in which one would count  $n$  events. These can be modeled as following a Poisson distribution with mean  $\mu s + b$ , where  $s$  and  $b$  are the expected numbers of events from the signal and background processes and  $\mu$  is a strength parameter. In particular one is interested in testing the hypothesis  $\mu = 0$ .

A complication that one would like to include is that the background MC events will have weights, and the distribution of weights will not be known a priori. (In principle the signal events could also have weights, but in general it is much easier to generate a large sample of signal MC so let us assume for the moment that this can be produced without weights.)

Suppose the Challenge is to design an analysis for a certain integrated luminosity,  $L_{\text{data}}$ , and that the MC samples correspond to effective integrated luminosities  $L_s$  and  $L_b$ . These are related to the data luminosity by

$$\begin{aligned} L_s &= \tau_s L_{\text{data}} , \\ L_b &= \tau_b L_{\text{data}} , \end{aligned}$$

where  $\tau_s$  and  $\tau_b$  are known scale factors.

The analyst produces in the multivariable space a selection region and this results in  $n_s$  events found in the signal sample and  $n_b$  events in the background sample, with weights  $w_1, w_2, \dots, w_{n_b}$ .

From these values one can produce estimates<sup>1</sup> for  $s$ ,

$$\tilde{s} = n_s / \tau_s , \tag{1}$$

and for  $b$

$$\tilde{b} = \frac{1}{\tau_b} \sum_{i=1}^{n_b} w_i . \tag{2}$$

Given large enough MC samples,  $s$  and  $b$  could be estimated with arbitrary precision. Here we will assume  $\tilde{s}$  is determined with negligible uncertainty but that this is not true for  $b$ .

---

<sup>1</sup>Here we use tildes for estimators rather than hats to distinguish them from other estimators for the same quantities that will enter later in the analysis.

Rather we will take  $\tilde{b}$  as Gaussian distributed with a mean equal to the true  $b$  and a variance  $\sigma_{\tilde{b}}^2$ , which can be estimated as

$$\tilde{\sigma}_{\tilde{b}}^2 = \frac{1}{\tau_{\tilde{b}}^2} \sum_{i=1}^{n_b} w_i^2 . \quad (3)$$

In the “real” experiment, one finds therefore  $n$  and  $\tilde{b}$ , which are modeled as

$$n \sim \text{Poisson}(\mu s + b) , \quad (4)$$

$$\tilde{b} \sim \text{Gauss}(b, \sigma_{\tilde{b}}) , \quad (5)$$

and for  $\sigma_{\tilde{b}}$  we use the estimate from Eq. (3). The likelihood function is thus

$$L(\mu, b) = \frac{(\mu s + b)^n}{n!} e^{-(\mu s + b)} \frac{1}{\sqrt{2\pi\sigma_{\tilde{b}}}} e^{-(\tilde{b}-b)/2\sigma_{\tilde{b}}^2} . \quad (6)$$

The log-likelihood function is therefore found to be (up to an additive constant)

$$\ln L(\mu, b) = n \ln(\mu s + b) - (\mu s + b) - \frac{1}{2} \frac{(b - \hat{b})^2}{\hat{\sigma}_{\tilde{b}}^2} . \quad (7)$$

To test a hypothetical value of  $\mu$  we can use the profile likelihood ratio (see, e.g., [1]),

$$\lambda(\mu) = \frac{L(\mu, \hat{\hat{b}}(\mu))}{L(\hat{\mu}, \hat{b})} , \quad (8)$$

where the double hat notation indicates the value of  $b$  that maximizes  $L$  for the given value of  $\mu$ , and single hats denote the (unconditional) maximum-likelihood estimators. In particular we are interested in testing  $\mu = 0$ , so we need  $\lambda(0)$ , and for this we require  $\hat{\hat{b}}(0)$ . The ingredients are found to be

$$\hat{\mu} = \frac{n - \tilde{b}}{s} , \quad (9)$$

$$\hat{b} = \tilde{b} , \quad (10)$$

$$\hat{\hat{b}}(0) = \frac{\tilde{b} - \sigma_{\tilde{b}}^2}{2} + \frac{1}{2} \sqrt{(\tilde{b} - \sigma_{\tilde{b}}^2)^2 + 4\sigma_{\tilde{b}}^2 n} \quad (11)$$

Using these quantities one can then evaluate the statistic

$$q_0 = \begin{cases} -2 \ln \lambda(0) & \hat{\mu} > 0 , \\ 0 & \text{otherwise} , \end{cases} \quad (12)$$

and from this one can obtain the discovery significance

$$Z = \sqrt{q_0} . \tag{13}$$

This is equivalent to the  $p$ -value of the  $\mu = 0$  hypothesis; the two quantities are related by

$$Z = \Phi^{-1}(1 - p) , \tag{14}$$

where  $\Phi^{-1}$  is the quantile of standard Gaussian.

For purposes of the Challenge one would like to maximize the sensitivity for discovery of the signal process, which can be quantified as the median discovery significance under assumption of  $\mu = 1$  (i.e., assuming that the nominal signal model is correct). To obtain an estimate of this one can simply take the significance from Eq. (13) and evaluate it replacing the measured values  $n$  and  $\tilde{b}$  by their expectation values,  $\mu s + b$  and  $b$  (the so-called Asimov data set). For  $s$  and  $b$  one then uses the best available estimates, which here are simply  $\tilde{s}$  and  $\tilde{b}$ , and in addition one can use  $\sigma_{\tilde{b}}$  from Eq. (3). The task of computing the sensitivity (median significance) can thus be carried out with a few relatively simple formulas.

The procedure can be modified in a straightforward way to include more than one background component. One would in addition like to allow for a more realistic distribution of weights. That is, instead of assuming that  $\tilde{b}$  follows a Gaussian distribution, a more realistic model would take, say, a log-normal distribution for the weights. Unfortunately this does not allow for a simple calculation of the ingredients required to obtain the significance (the profiled parameters must be determined numerically) and so this approach is perhaps less appropriate for the data challenge.

## References

- [1] Glen Cowan, Kyle Cranmer, Eilam Gross and Ofer Vitells, Eur. Phys. J. C 71 (2011) 1-19; arXiv:1007.1727.