

# New approaches for boosting to uniformity

---

**Alex Rogozhnikov<sup>a,b,\*</sup>, Aleksandar Bukva<sup>c</sup>, Vladimir Gligorov<sup>d</sup>, Andrey Ustyuzhanin<sup>b,e,f</sup> and Mike Williams<sup>g</sup>**

<sup>a</sup>*Lomonosov Moscow State University, Moscow, Russia*

<sup>b</sup>*Yandex, Moscow, Russia*

<sup>c</sup>*Faculty of Physics, Belgrade, Serbia*

<sup>d</sup>*Organisation Européenne pour la Recherche Nucléaire (CERN), Geneva, Switzerland*

<sup>e</sup>*Moscow Institute of Physics and Technology, Moscow, Russia*

<sup>f</sup>*Imperial College, London, UK*

<sup>g</sup>*Massachusetts Institute of Technology, Cambridge, MA, United States*

*E-mail: alex.rogozhnikov@yandex.ru*

**ABSTRACT:** The use of multivariate classifiers has become commonplace in particle physics. To enhance the performance, a series of classifiers is typically trained; this is a technique known as boosting. This paper explores several novel boosting methods that have been designed to produce a uniform selection efficiency in a chosen multivariate space. Such algorithms have a wide range of applications in particle physics, from producing uniform signal selection efficiency across a Dalitz-plot to avoiding the creation of false signal peaks in an invariant mass distribution when searching for new particles.

---

\*Corresponding author.

---

## Contents

<b>1. Introduction</b>	<b>1</b>
<b>2. Uniformity Boosting Methods</b>	<b>2</b>
<b>3. Example Analysis</b>	<b>4</b>
<b>4. CPU Resources</b>	<b>7</b>
<b>5. Summary</b>	<b>8</b>
<b>6. Source code</b>	<b>8</b>
<b>A. Measures of uniformity</b>	<b>9</b>

---

## 1. Introduction

Methods of machine learning play an important role in modern particles physics. Multivariate classifiers, *e.g.*, boosted decision trees (BDTs) and artificial neural networks (ANNs), are now commonly used in analysis selection criteria. BDTs are now even used in software triggers [1, 2]. To enhance the performance, a series of classifiers is typically trained; this is a technique known as boosting. Boosting involves training many simple classifiers and then building a single composite classifier from their responses. The classifiers are trained in series with the inputs of each member being augmented based on the performance of its predecessors. This augmentation is designed such that each new classifier targets those events which were poorly classified by previous members of the series. The classifier obtained by combining all members of the series is typically much more powerful than any of the individual members.

In particle physics, the most common usage of BDTs is in classifying candidates as signal or background. The BDT is determined by optimizing some figure of merit (FOM), *e.g.*, the signal purity or approximate signal significance. This approach is optimal for a counting experiment; however, in many cases the BDT-based selection obtained in this way is not optimal. For example, in a Dalitz-plot (or any angular or amplitude analysis) analysis, obtaining a selection efficiency for signal candidates that is uniform across the Dalitz-plot is more important than any integrated FOM. Similarly, when measuring a mean particle lifetime, obtaining an efficiency that is uniform in lifetime is what is desired. In both cases, obtaining a uniform selection efficiency greatly reduces the systematic uncertainties involved in the measurement. When searching for a new particle, an analyst may want a uniform efficiency in mass for selecting background candidates so that the BDT-based selection does not generate a fake signal peak. Furthermore, the analyst may also desire a uniform selection efficiency of signal candidates in mass (or other variates) since the mass of the

new particle is not known. In such cases, the BDT is often trained on simulated data generated with several values of mass (lifetime, *etc.*). A uniform selection efficiency in mass ensures that the BDT is sensitive to the full range of masses involved in the search.

## 2. Uniformity Boosting Methods

The variates used in the BDT are denoted by  $\vec{x}$ , while the variates in which uniformity is desired are denoted by  $\vec{y}$ . Some (perhaps all) of the  $\vec{x}$  variates will be *biasing* in  $\vec{y}$ , *i.e.* they provide discriminating power between signal and background that varies in  $\vec{y}$ . A uniform BDT selection efficiency can be obtained by removing all such variates; however, this will also reduce the power of the BDT. The goal of boosting algorithms presented in this paper is to *balance* the biases to produce the optimal uniform selection.

One category of boosting works by assigning training events more weight based on classification errors made by previous members of the series. For example, the AdaBoost [3] algorithm updates the weight of event  $i$ ,  $w_i$ , according to

$$w'_i = w_i \times \exp[-\gamma p_i], \quad (2.1)$$

where  $\gamma = +1(-1)$  for signal(background) events and  $p$  is the prediction for each event produced by last classifier in the series. The uBoost technique, described in detail in Ref. [4], alters the event-weight updating procedure to achieve uniformity in the signal-selection efficiency.

Another approach to obtain uniformity, introduced in this paper, involves defining a more general expression of the AdaBoost criteria:

$$w'_i = w_i \times \exp \left[ -\gamma \sum_j a_{ij} p_j \right], \quad (2.2)$$

where  $a_{ij}$  are the elements of some square matrix  $A$ . For the case where  $A$  is the identity matrix, the AdaBoost weighting procedure is recovered. Other choices of  $A$  will induce non-local effects, *e.g.*, consider the sparse matrix  $A_{\text{knn}}$  given by

$$a_{ij}^{\text{knn}} = \begin{cases} \frac{1}{k}, & j \in \text{knn}(i), \text{ events } i \text{ and } j \text{ belong to the same class} \\ 0, & \text{otherwise,} \end{cases} \quad (2.3)$$

where  $\text{knn}(i)$  denotes the set of  $k$ -nearest-neighbor events to event  $i$ . This procedure for updating the event weights, which we refer to as kNNAdaBoost, accounts for the score of each event's  $k$  nearest neighbors and not just each event individually.

The gradient boosting [5] (GB) algorithm category requires the analyst to choose a differentiable loss function with the goal of building a classifier that minimizes the loss. A popular choice of loss function is the so-called AdaLoss function

$$L_{\text{ada}} = \sum_{i=1}^n \exp[-\gamma s_i]. \quad (2.4)$$

The *scores*  $s$  are obtained for each event as the sum of predictions of all elements in the series. At each stage in the gradient boosting process, a *regressor* (a decision tree in our case) is trained

whose purpose is to decrease the loss. This is accomplished using the gradient decent method and the *pseudo-residuals*

$$-\frac{\partial L_{\text{ada}}}{\partial s_i} = \gamma_i \exp[-\gamma_i s_i], \quad (2.5)$$

which are positive(negative) for signal(background) events and have larger moduli for poorly classified events.

The gradient-boosting algorithm is general in that it only requires the analyst specify a loss function and its gradient. The AdaLoss function considers each event individually, but can easily be modified to take into account non-local properties of the classifier as follows:

$$L_{\text{general}} = \sum_{i=1}^n \exp \left[ -\gamma_i \sum_j a_{ij} s_j \right]. \quad (2.6)$$

For example, the loss function obtained from Eq. 2.6 using<sup>1</sup>  $A_{\text{knn}}$ , which we refer to as kN-NAdaLoss and denote  $L_{\text{knn}}$ , accounts for the score of each event's  $k$  nearest neighbors and not just each event individually. The pseudo-residuals of  $L_{\text{knn}}$  are

$$-\frac{\partial L_{\text{knn}}}{\partial s_k} = \sum_i \gamma_i a_{ik}^{\text{knn}} \exp \left[ -\gamma_i \sum_j a_{ij}^{\text{knn}} s_j \right]. \quad (2.7)$$

One can see that the direction of the gradient will be influenced the most by events whose  $k$ -nearest-neighbor events are classified poorly. We generically refer to GB methods designed to achieve uniform selection efficiency as uniform GB (uGB). The specific algorithm that uses kNNAdaLoss will be called uGBkNN.

Another approach is to include some uniformity metric in the definition of the loss function. Consider first the case where the data have been binned in  $\vec{y}$ . If the distribution of classifier responses in each bin,  $f_b(s)$ , is the same as the global response distribution,  $f(s)$ , then any cut made on the response will produce a uniform selection efficiency in  $\vec{y}$ . Therefore, performing a one-dimensional goodness-of-fit test of the hypothesis  $f_b \equiv f$  in each bin provides an assessment of the selection uniformity. For example, one could perform the Kolmogorov-Smirnov test in each bin and define a loss function as follows:

$$L_{\text{flat(KS)}} = \sum_b w_b \max |F_b(s) - F(s)|, \quad (2.8)$$

where  $F_b(s)$  denotes the cumulative distribution of  $f_b(s)$  and  $w_b = \sum \delta(\text{bin}_i - b)/n_{\text{signal}}$ , i.e.  $w_b$  is the fraction of signal events in the bin<sup>2</sup>.

The gradient of the Kolmogorov-Smirnov loss function is zero for events with responses greater than the value of  $s$  at which  $\max |F_b(s) - F(s)|$  occurs. Thus, it is not suitable for gradient boosting due to its instability. Instead, we use the following *flatness* loss function:

$$L_{\text{flat}} = \sum_b w_b \int |F_b(s) - F(s)|^2 ds, \quad (2.9)$$

<sup>1</sup>A natural choice is a square  $n \times n$  matrix, but this is not required.

<sup>2</sup>If weighted events are used, then the fractional sum of weights should be used for  $w_b$ .

Table 1: Description of uniform boosting algorithms.

Name	Description
uBoost	algorithm introduced in Ref. [4]
kNNAda	AdaBoost modification using matrix $A_{\text{knn}}$
uGBkNN	gradient boost using kNNAdaLoss loss function
uGBFL(bin)	gradient boost using flatness loss $+\alpha$ AdaLoss as in Eq. 2.11 (data binned for FL)
uGBFL(kNN)	same as uGBFL(bin) except kNN events are used rather than bins

whose pseudo-residuals are ( $b$  is the bin containing the  $k$ th event)

$$-\frac{\partial L_{\text{flat}}}{\partial s_k} \approx -2w_b [F_b(s_k) - F(s_k)]. \quad (2.10)$$

This so-called flatness loss penalizes non-uniformity but does not consider the quality of the classification. Therefore, the full loss function used is

$$L_{\text{ada+flat}} = L_{\text{flat}} + \alpha L_{\text{ada}}, \quad (2.11)$$

where  $\alpha$  is a real-valued parameter that is typically chosen to be small. The first term in Eq. 2.11 penalizes non-uniformity, while the second term penalizes poor classification. We refer to this algorithm as uGB with flatness loss (uGBFL). In principle, many different flatness loss functions can be defined and could be substituted for our choice here. See Appendix A for a detailed discussion on this topic.

The loss function given in Eq. 2.11 can also be constructed without binning the data using k-nearest-neighbor events. The cumulative distribution  $F_{\text{knn}}(s)$  is easily obtained and the bin weight,  $w_b$ , is replaced by a k-nearest-neighbor weight,  $w_{\text{knn}}$ . First, each event is weighted by the inverse of the number of times it is included in the k-nearest-neighbor sample of another event. Then,  $w_{\text{knn}}$  is the sum of such weights in a k-nearest-neighbor sample divided by the total sum of such weights in the full sample. This procedure is followed to offset the fact that some events are found in more k-nearest-neighbor samples than other events. We study two versions of uGBFL below: uGBFL using bins denoted by uGBFL(bin) and uGBFL using kNN collections denoted by uGBFL(kNN). The algorithms are summarized in Table 1.

### 3. Example Analysis

The example analysis studied here involves a so-called Daltiz-plot analysis. In such analyses, the distribution of events in a 2-D space is typically fit to extract some information of physical interest. The regions of the Daltiz-plot that tend to have the highest sensitivity to the desired information are the *edges*. Unfortunately, the edge regions also typically have the most background contamination and the least discrimination against background. Therefore, traditional classifier-based selections tend to produce selections for Dalitz-plot analyses with lower efficiency near the edges.

This study uses simulated event samples produced using the official LHCb simulation framework. The software used for the generation of the events is described in LHCb publications as follows :

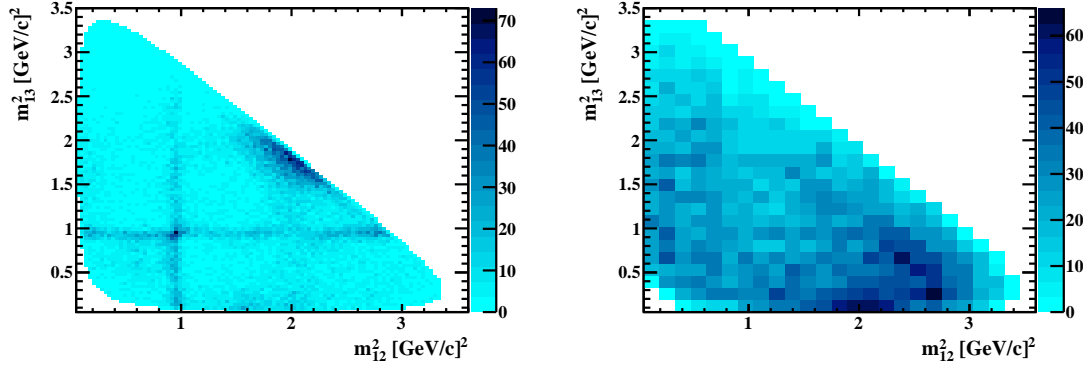


Figure 1: Dalitz-plot distributions for (left) signal and (right) background for the  $D_s^\pm \rightarrow \pi^+ \pi^- \pi^\pm$ . The three pions are labeled here as 1, 2 and 3 and ordered according to increases momentum.

In the simulation,  $pp$  collisions are generated using PYTHIA [6] with a specific LHCb configuration [7]. Decays of hadronic particles are described by EvtGen [8], in which final state radiation is generated using PHOTOS [9]. The interaction of the generated particles with the detector and its response are implemented using the GEANT toolkit [10, 11] as described in Ref. [12].

All simulated event samples are generated inside the LHCb detector acceptance. The signal used in this analysis consists of  $D_s^\pm \rightarrow \pi^+ \pi^- \pi^\pm$  decays, simulated using the D\_DALITZ model of EvtGen to simulate the intermediate resonances which contribute to the three pion final state. The background candidates are three pion combinations reconstructed in simulated samples of  $c\bar{c}$  and  $b\bar{b}$  events, where the charm and bottom quark decays are inclusively modelled by EvtGen. The simulated events contain “truth” information which identifies them as signal or background, and which identifies the physical origin of the three pion combinations reconstructed in the  $c\bar{c}$  and  $b\bar{b}$  simulated samples.

Figure 1 shows the Dalitz-plot distributions for signal and background events. These samples are split into training and testing samples and then various BDTs are trained. For the BDTs designed to produce uniform selections, the  $\vec{y}$  variates are the Dalitz masses with the choice of uniform selection efficiency on signal candidates in the Dalitz-plot. Figure 2 shows the ROC curves obtained for the various classifiers studied in this paper. For the uGBFL algorithms, there is a choice to be made for the value  $\alpha$  which defines the relative weight of the flatness loss vs AdaLoss. As expected, increasing  $\alpha$ , which increases the weight of AdaLoss, drives the ROC curve to be similar to AdaBoost. Analysts will need to choose how much ROC performance to sacrifice to gain uniformity in the selection efficiency. In general, the ROC curves for the uniform-driven BDTs are not too different from AdaBoost.

Figure 3 shows the efficiency obtained for each classifier vs distance from the a corner of the Dalitz-plot<sup>3</sup>. The AdaBoost algorithm, as expected, produces a much lower efficiency in the interesting corner regions. The kNNAdaBoost algorithm does not improve upon the AdaBoost result

<sup>3</sup>The Dalitz-plot is essentially a triangle with three corners. Our definition of this distance is  $\text{MIN} \left[ (m(D_s) - m(\pi))^2 - m_{ij}^2 \right]$ , where  $ij$  is 12, 13 and 23.

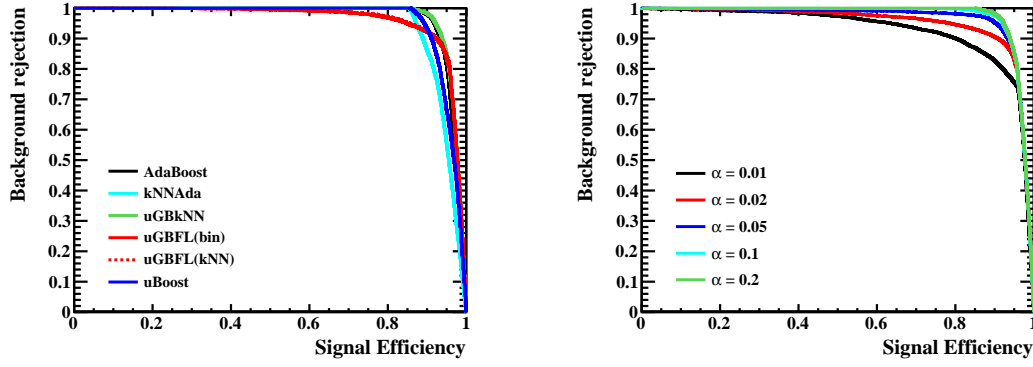


Figure 2: (left) ROC curves for classifier algorithms studied in this paper. For the uGBFL algorithms  $\alpha = 0.02$  is shown. (right) ROC curves for uGBFL(bin) for different values of  $\alpha$ .

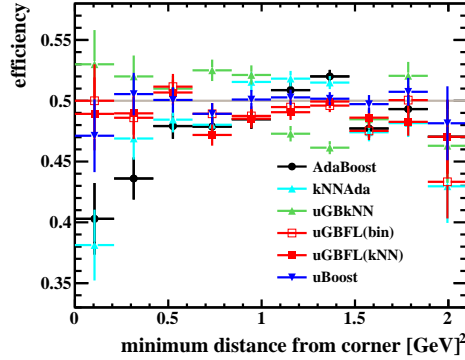


Figure 3: Efficiency vs distance to a corner of the Dalitz-plot. An arbitrary working point of 50% integrated efficiency is displayed. For the uGBFL algorithms  $\alpha = 0.02$  is shown.

much. This is likely due to the fact that while kNNAdaBoost uses non-local kNN information, it does not utilize global information. The uGBkNN algorithm overcompensates and drives the efficiency higher at the corners. This suggests that if this algorithm is to be used some stopping criteria or *throttle* of the event-weighting updating should be implemented. The uGBFL (binned and unbinned kNN) and uBoost algorithms each produce an efficiency which is statistically consistent with uniform across the Dalitz plot. As stated above, the analyst is free to optimize the choice of  $\alpha$  for uGBFL by defining a metric that involves signal efficiency, background rejection and uniformity, *e.g.*, using uniformity metrics discussed in detail in Appendix A.

As a separate study using the same data samples, consider the case where one has simulated signal events and uses data from a *nearby* region, a so-called sideband, for background. This is a common situation in particle-physics analyses. Figure 4 shows the training samples used. A major problem can arise in these situations as typically input variables to the BDT are correlated with the parent particle mass. Therefore, the BDT may learn to reject the background in the training using the fact that the mass of the background and signal candidates is different. This is just an artifact of how the background sample is obtained and will not be true for background candidates *under* the signal peak. Figure 4 shows the background mis-identification rate vs  $D$  candidate mass. AdaBoost

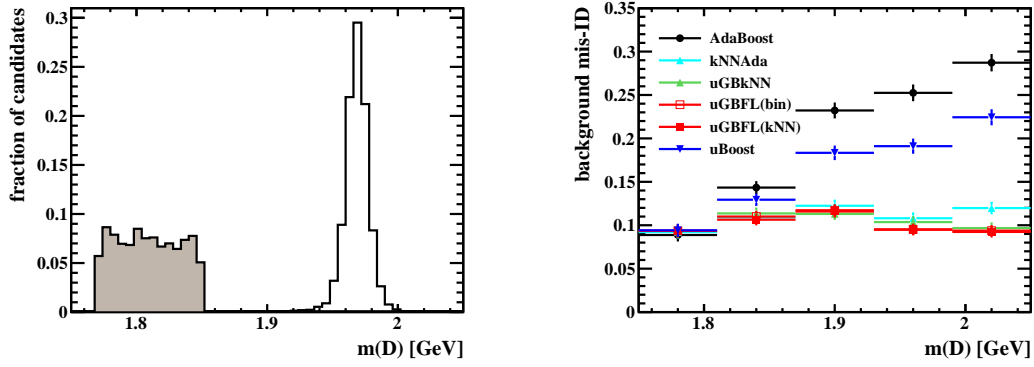


Figure 4: (left) Signal and (filled) background samples used in training. (right) Background mis-identification *vs*  $D$  candidate mass for an arbitrary working point of 10% background mis-identification in the training region  $1.75 < m(D) < 1.85$  GeV is displayed.

has clearly learned to use this mis-match in signal and background candidate masses in the training. The background in the region of the signal is about three times higher than one would expect from looking only at the sideband data.

Figure 4 also shows the background mis-identification rate *vs*  $D$  candidate mass for the various uniform classifiers where  $y = m(D)$  and the choice is for uniformity in the background efficiency<sup>4</sup>. The uBoost algorithm does better than AdaBoost here but is still not optimal. The way that uBoost achieves uniformity is not such that it can be trusted to work outside the region of training. The algorithms presented in this paper each does well in achieving similar performance in the training and signal regions. Consider, *e.g.*, the uGBFL approach to achieving uniform selection efficiency. In this case the training drives the BDT response itself to have the same PDF everywhere in the region  $1.75 < m(D) < 1.85$  GeV (the training region). This does not guarantee that the BDT response is truly independent of  $m(D)$  outside the training region, but does strongly suppress learning to use  $m(D)$  and in this example results in the desired behavior. Finally, if both high and low  $m(D)$  sidebands had been used, it is possible for a BDT to create a fake peak near the signal peak location. The use of uGBFL greatly reduces the chances and possible size of such an effect.

#### 4. CPU Resources

One drawback of the uBoost technique is that it has a high degree of computational complexity: while AdaBoost trains  $M$  trees (a user-defined number), uBoost builds  $100 \times M$  trees. The algorithms presented in this paper only build  $M$  trees; however, the boosting involves some more complicated algorithms. Training each of the  $M$  trees scales as follows for  $N$  training events:

- kNNAdaBoost:  $O(k \times N)$ ;
- uGBkNNknn:  $O(k \times N)$  for  $A_{\text{knn}}$ , and  $O(\text{\#nonzero elements in the matrix})$  for arbitrary matrix  $A$ ;

<sup>4</sup>The algorithms in this paper can easily be made uniform on signal, background or both.



- uGBFL(bin):  $O(N \ln N)$ ;
- uGBFL(kNN):  $O(N \ln N + Nk \ln k)$ .

In the example analysis studied in this paper, we find that the training time for these new algorithms is within a factor of two the same as AdaBoost. The CPU-resource usage of these new algorithms is not prohibitive.

## 5. Summary

A number of novel boosting algorithms have been presented that consider uniformity of selection efficiency in a multivariate space in addition to mis-classification errors. Of these, the uGBFL algorithm has the best performance on the example analyses studied in this paper. This algorithm is expected to be useful in a wide-variety of analyses performed in particle physics.

## 6. Source code

The code for classifiers proposed in this article as well as for metrics of uniformity is publicly available at repository [https://github.com/anaderi/lhcb\\_trigger\\_ml](https://github.com/anaderi/lhcb_trigger_ml).

## Acknowledgments

These results were obtained using events generated with the official LHCb simulation, and we are grateful to the LHCb collaboration for this privilege. We particularly acknowledge the work of the LHCb Simulation and Core Computing teams in tuning the simulation software and managing the productions of simulated events. MW is supported by NSF grant PHY-1306550.

## References

- [1] R. Aaij *et al.* [LHCb Trigger Group], *The LHCb trigger and its performance*, JINST **8**, P04022 (2013). [arXiv:1211.3055]
- [2] V. V. Gligorov and M. Williams, *Efficient, reliable and fast high-level triggering using a bonsai boosted decision tree*, JINST **8**, P02013 (2013). [arXiv:1210.6861]
- [3] Y. Freund, and R. Schapire, *A Decision-Theoretic Generalization of On-Line Learning and an Application to Boosting*, *Journal of Computer and System Sciences* **55**(1) (August 1997) 119-139.
- [4] J. Stevens and M. Williams, *uBoost: A boosting method for producing uniform selection efficiencies from multivariate classifiers*, JINST **8**, P12013 (2013). [arXiv:1305.7248]
- [5] Jerome H. Friedman, *Greedy Function Approximation: A Gradient Boosting Machine*, *Annals of Statistics* **29** (2000) 1189-1232.
- [6] T. Sjöstrand, S. Mrenna and P. Skands, *PYTHIA 6.4 physics and manual*, JHEP **05** (2006) 026. [hep-ph/0603175]
- [7] I. Belyaev, I. *et al.*, *Handling of the generation of primary events in GAUSS, the LHCb simulation framework*, NSS/MIC IEEE (2010) 1155.

- [8] D.J. Lange, *The EvtGen particle decay simulation package*, *NIM* **A462** (2001) 152-155.
- [9] P. Golonka and Z. Was, *PHOTOS Monte Carlo: a precision tool for QED corrections in Z and W decays*, *EPJ* **C45** (2006) 97-107. [[hep-ph/0506026](#)]
- [10] J. Allison *et al.*, *Geant4 developments and applications*, *IEEE TNS* **53** (2006) 270.
- [11] S. Agostinelli *et al.*, *Geant4: a simulation toolkit*, *NIM* **A506** (2003) 250.
- [12] M Clemencic *et al.*, *The LHCb simulation application, GAUSS : design, evolution and experience*, *J. Phys. Conf. Ser.* **331** (2011) 032023.

## A. Measures of uniformity

In this section we discuss different methods for measuring the uniformity of prediction. One typical way of ‘checking’ uniformity of prediction used by physicists is fitting the distribution of the events that were classified as signal (or background) over the feature for which you wish to check uniformity. This approach requires assumptions about the shape of the distribution, which makes quantitative comparisons of different classifiers difficult. Our aim here is to explore uniformity figures of merit which make comparing classifiers easier, analogously to how the area under the ROC curve can be used to compare absolute classifier performance.

The output of event classification is the probability of each event being signal or background, and it is only after we apply a cut on this probability that events are classified. An ideal uniformity of signal prediction can then be defined for a given “uniform feature” of interest. It means that whichever cut we select, the efficiency for a signal event to pass the cut doesn’t depend on the uniform feature. Uniformity for background can be defined in the same manner, but for simplicity, in what follows we will only discuss the uniformity of efficiency for signal events.

A trivial example of a classifier that has ideal uniformity is a classifier which returns a random classification probability, but such a classifier is of course not very useful. One can try to design a uniform classifier with respect to a given feature by not using this feature, or any correlated features, in the classification; in practice, however, this approach also tends to lead to poorly performing classifiers. The approach which we take in this paper is to explicitly let the classifier learn how to balance non-uniformities coming from different features in such a way as to generate a classification which is uniform on average. It is then important to be able to accurately measure the uniformity of classification.

Before proceeding, it is useful to define some desirable properties of uniformity metrics

1. The metric shouldn’t depend strongly on the number of events used to test uniformity;
2. The metric shouldn’t depend on the normalization of the event weights: if we multiply all the weights by some arbitrary number, it shouldn’t change at all;
3. The metric should depend only on the order of predictions, not the exact values of probabilities. This is because we care about which events pass the cut and which don’t, not about the exact values of predictions. For example: correlation of prediction and mass doesn’t satisfy this restriction.

4. The metric should be stable against any of its own free parameters: if it uses bins, changing the number of bins shouldn't affect the result, if it uses  $k$ -nearest neighbors, it should be stable against different values of  $k$ .

In what follows we will consider different metrics which satisfy these criteria, and then compare their performance in some test cases.

### Standard Deviation of Efficiency on Bins (SDE)

If the space of uniform features is split into bins, it is possible to define the global efficiency

$$\text{eff} = \frac{\text{total weight of signal events that passed the cut}}{\text{total weight of signal events}},$$

as well as the efficiency in every bin,

$$\text{eff}_{\text{bin}} = \frac{\text{weight of signal events in bin that passed the cut}}{\text{weight of signal events in this bin}}.$$

One measure of non-uniformity is the standard deviation of bin efficiencies from the global efficiency:

$$\sqrt{\sum_{\text{bin}} (\text{eff}_{\text{bin}} - \text{eff})^2}.$$

To make the metric more stable against fluctuations in bins which contain very few events, we add weights to the bins (note that  $\sum_{\text{bin}} \text{weight}_{\text{bin}} = 1$ ):

$$\text{weight}_{\text{bin}} = \frac{\text{total weight of signal events in bin}}{\text{total weight of signal events}},$$

giving the weighted standard deviation (SDE) formula

$$\text{SDE}(\text{eff}) = \sqrt{\sum_{\text{bin}} \text{weight}_{\text{bin}} \times (\text{eff}_{\text{bin}} - \text{eff})^2}.$$

This formula is valid for any given cut value. To measure the overall non-flatness of the selection, we take several global efficiencies and use

$$\text{SDE}^2 = \frac{1}{k} \sum_{\text{eff} \in [\text{eff}_1 \dots \text{eff}_k]} \text{SDE}^2(\text{eff})$$

Another power  $p \neq 2$  can be used as well, but  $p = 2$  is considered as the default value.

### Theil Index of Efficiency

The Theil Index is frequently used to measure economic inequality:

$$\text{Theil} = \frac{1}{N} \sum_i \frac{x_i}{\langle x \rangle} \ln \frac{x_i}{\langle x \rangle}, \quad \langle x \rangle = \frac{1}{N} \sum_i x_i$$

In our case we have to alter formula a bit to take into account that different bins have different impact, thus the formula turns into

$$\text{Theil}(\text{eff}) = \sum_{\text{bin}} \text{weight}_{\text{bin}} \frac{\text{eff}_{\text{bin}}}{\text{eff}} \ln \frac{\text{eff}_{\text{bin}}}{\text{eff}}.$$

To measure the overall non-flatness, we average values for several global efficiencies:

$$\text{Theil} = \frac{1}{k} \sum_{\text{eff} \in [\text{eff}_1 \dots \text{eff}_k]} \text{Theil}(\text{eff})$$

## Distribution Similarity Approach

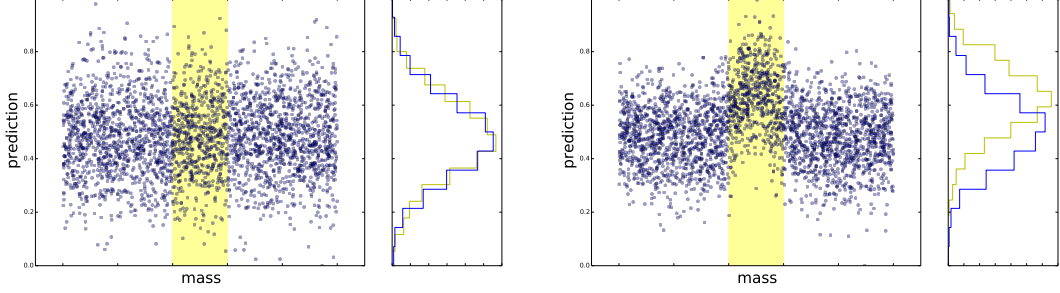


Figure 5: Demonstration of the distribution similarity approach. (left) Predictions are uniform in mass, the distribution of predictions in the bin (yellow) is close to the global (blue). (right) Distribution with peak in the middle, the distribution in the bin is quite different from the global distribution. In both cases the yellow rectangle shows the events in the bin over mass.

Instead of measuring uniformity in terms of binned efficiencies, it is possible to consider the distribution of the binned classifier predictions,  $F_{\text{bin}}$ , directly. Ideal uniformity means that all the distributions  $F_{\text{bin}}$  are equal and hence equal to the global distribution  $F(x)$ . This is demonstrated on figure 5. To 'measure' non-flatness we can use some distribution distance, like Kolmogorov-Smirnov:

$$\sum_{\text{bin}} \text{weight}_{\text{bin}} \max_x |F_{\text{bin}}(x) - F(x)|,$$

but Cramér–von Mises similarity is more informative (usually  $p = 2$  is used):

$$\sum_{\text{bin}} \text{weight}_{\text{bin}} \int |F_{\text{bin}}(x) - F(x)|^p dF(x),$$

in particular because Kolmogorov-Smirnov measures are too sensitive to local non-uniformities. The advantage of this method is that we don't need to select some global efficiencies like in the previous metrics.

## Knn-based modifications

Though operating with bins is usually both simple and very efficient, in many cases it is hard to find the optimal size of bins in the space of uniform features (specifically in the case of more than two dimensions). As mentioned earlier, problems can also arise due to bins with very low populations.

In these cases we can switch to  $k$ -nearest neighbors: for each signal event we find  $k$  nearest signal events (including the event itself) in the space of uniform features. Now we can compute the efficiency  $\text{eff}_{\text{knn}(i)}$ , from the empirical distribution  $F_{\text{knn}(i)}$  of nearest neighbors. The weights for  $\text{knn}(i)$  are proportional to the total weight of events in  $\text{knn}(i)$ :

$$\text{weight}_{\text{knn}(i)} = \alpha \sum_{j \in \text{knn}(i)} w_j, \quad \alpha^{-1} = \sum_i \sum_{j \in \text{knn}(i)} w_j,$$

so again weights are normed to 1:  $\sum_i \text{weight}_{\text{knn}(i)} = 1$ .

It is then possible to write the knn versions of SDE

$$\text{knnSDE}^2(\text{eff}) = \sum_{i \in \text{events}} \text{weight}_{\text{knn}(i)} |\text{eff}_{\text{knn}(i)} - \text{eff}|^2,$$

$$\text{knnSDE}^2 = \frac{1}{k} \sum_{\text{eff} \in [\text{eff}_1, \dots, \text{eff}_k]} \text{knnSDE}^2(\text{eff}),$$

the Theil index of efficiency

$$\text{knnTheil}(\text{eff}) = \sum_{i \in \text{events}} \text{weight}_{\text{knn}(i)} \frac{\text{eff}_{\text{knn}(i)}}{\text{eff}} \ln \frac{\text{eff}_{\text{knn}(i)}}{\text{eff}},$$

$$\text{knnTheil} = \frac{1}{k} \sum_{\text{eff} \in [\text{eff}_1, \dots, \text{eff}_k]} \text{knnTheil}(\text{eff}),$$

and the similarity-based measure:

$$\sum_{i \in \text{events}} \text{weight}_{\text{knn}(i)} \int |F_{\text{knn}(i)}(x) - F(x)|^p dF(x).$$

The knn approach suffers from a drawback: the impact of different events has very little connection with the weights, because some events are selected as nearest neighbours much more frequently than others. This effect can be suppressed by dividing the initial weight of the event by the number of times it is selected as a nearest neighbour.

## Advantages and Disadvantages of Different Metrics

### Theil and SDE

Let us compare two metrics that have proven most appropriate for our problem, SDE and Theil. We have some masses distributed uniformly in  $[0, 1]$ , some constant  $\alpha$  from interval  $[0, 1]$ . The predictions are correlated with mass via beta distribution. First distribution with a peak is obtained by generating a prediction for each event according to its mass:

$$p_{\text{peak}} \sim \text{Beta}(1 + \alpha f(m), 1), \quad f(m) = 5 \times e^{-100|m - \bar{m}|^2}$$

The second distribution is obtained by flipping:

$$p_{\text{pit}} \sim 1 - \text{Beta}(1 + \alpha f(m), 1)$$

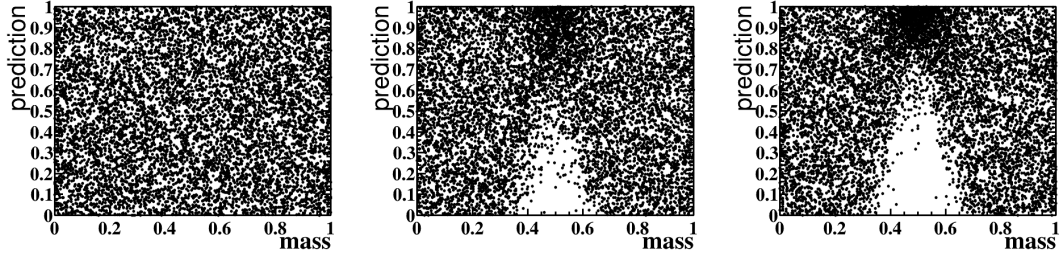
SDE should show no changes if we flip the distribution, while the Theil should make difference between pit and peak.

Figures 6 and 7 show the distribution of the predictions and classifier efficiency as a function of mass for the peak and pit distributions, respectively.

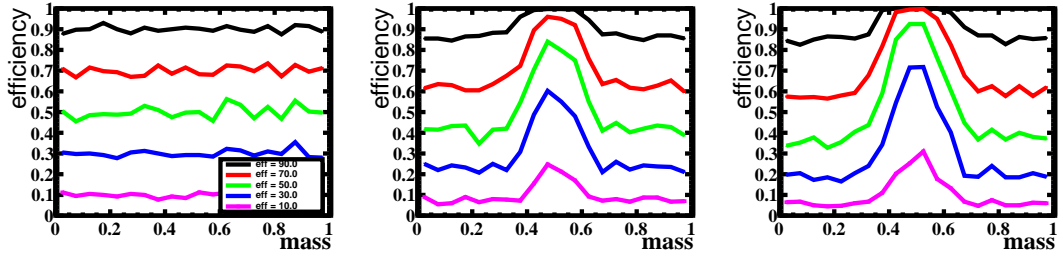
From figure 8 we can see that SDE doesn't make any difference between these distributions, while Theil has lower value in the second case which indicates that distribution is flatter. This example demonstrates that Theil has larger penalty for distributions with narrow peaks rather than with narrow pits in the distribution of efficiencies.

### $D \rightarrow hhh$

Finally, we compare the SDE and Theil uniformity measures for the classifiers applied to the  $D \rightarrow hhh$  data set, as shown in Fig. 9. Both measures show similar results, so there is no significant difference between them, nor the CvM metric, on this dataset.

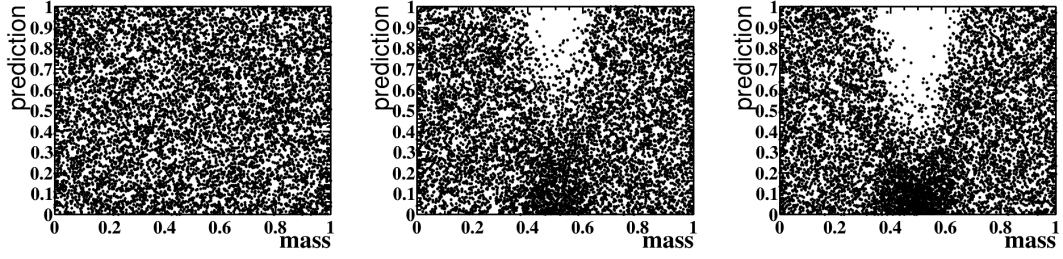


(a) Mass vs prediction.

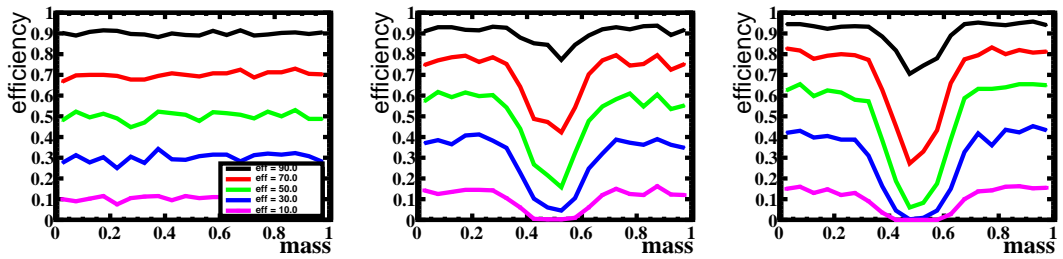


(b) Mass vs efficiency. The different coloured lines correspond to different global classifier efficiencies, as explained in the legend on the leftmost subplot.

Figure 6: Peak response distribution and efficiencies as a function of mass.



(a) Mass vs prediction.



(b) Mass vs efficiency. The different coloured lines correspond to different global classifier efficiencies, as explained in the legend on the leftmost subplot.

Figure 7: Pit response distribution and efficiencies as a function of mass.

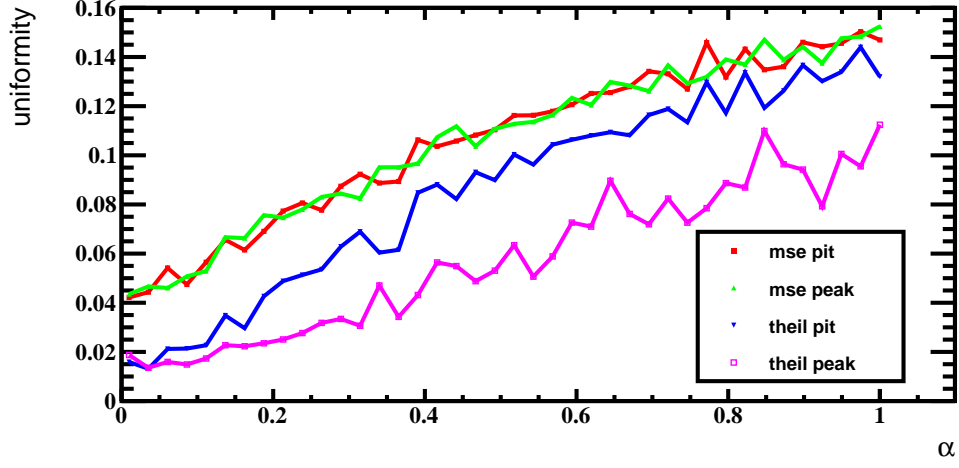


Figure 8: Comparison of SDE and Theil uniformity measures for different values of the constant  $\alpha$  for the peak and pit distributions. The Theil can clearly distinguish between these two while the SDE measure cannot.

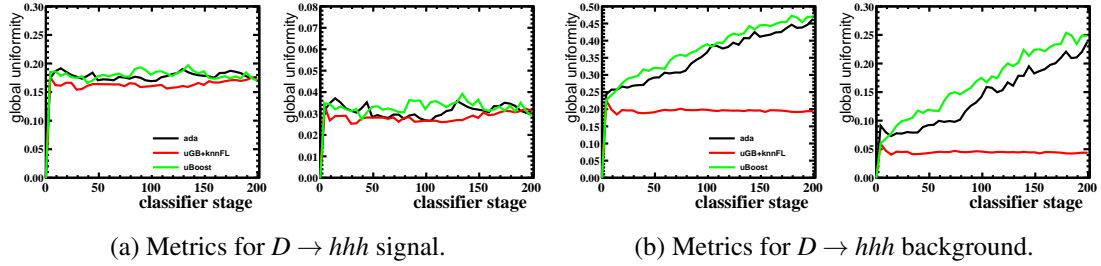


Figure 9: Uniformity metrics for the  $D \rightarrow hhh$  data set. The (top) signal and (bottom) background uniformities are plotted as a function of the training stage of a given classifier, listed in the legends on the leftmost plots. Within each figure the leftmost plot is the SDE and the rightmost plot is the Theil uniformity metric.