

# Punzi-loss

## A non-differentiable metric approximation for sensitivity optimisation in the search for new particles

F. Abudinén<sup>14</sup>, M. Bertemes<sup>16</sup>, S. Bilokin<sup>18</sup>, M. Campajola<sup>4,9</sup>, G. Casarosa<sup>3,11</sup>, S. Cunliffe<sup>2</sup>, L. Corona<sup>3,11</sup>, M. De Nuccio<sup>2</sup>, G. De Pietro<sup>12</sup>, S. Dey<sup>20</sup>, M. Eliachevitch<sup>21</sup>, P. Feichtinger<sup>a,16</sup>, T. Ferber<sup>15</sup>, J. Gemmler<sup>15</sup>, P. Goldenzweig<sup>15</sup>, A. Gottmann<sup>15</sup>, E. Graziani<sup>12</sup>, H. Haigh<sup>16</sup>, M. Hohmann<sup>22</sup>, T. Humair<sup>19</sup>, G. Inguglia<sup>16</sup>, J. Kahn<sup>7</sup>, T. Keck<sup>15</sup>, I. Komarov<sup>2</sup>, J.-F. Krohn<sup>22</sup>, T. Kuhr<sup>18</sup>, S. Lacaprara<sup>10</sup>, K. Lieret<sup>18</sup>, R. Maiti<sup>16</sup>, A. Martini<sup>2</sup>, F. Meier<sup>5</sup>, F. Metzner<sup>15</sup>, M. Milesi<sup>22</sup>, S.-H. Park<sup>8</sup>, M. Prim<sup>21</sup>, C. Pulvermacher<sup>15</sup>, M. Ritter<sup>18</sup>, Y. Sato<sup>8</sup>, C. Schwanda<sup>16</sup>, W. Sutcliffe<sup>21</sup>, U. Tamponi<sup>13</sup>, F. Tenchini<sup>11</sup>, P. Urquijo<sup>22</sup>, L. Zani<sup>1</sup>, R. Žlebčík<sup>6</sup>, A. Zupanc<sup>17</sup>

<sup>1</sup>Aix Marseille Université, CNRS/IN2P3, CPPM, 13288 Marseille, France

<sup>2</sup>Deutsches Elektronen-Synchrotron, Hamburg, Germany

<sup>3</sup>Dipartimento di Fisica, Università di Pisa, I-56127 Pisa, Italy

<sup>4</sup>Dipartimento di Scienze Fisiche, Università di Napoli Federico II, I-80126 Napoli, Italy

<sup>5</sup>Duke University, Durham, USA

<sup>6</sup>Faculty of Mathematics and Physics, Charles University, Prague, Czech Republic

<sup>7</sup>Helmholtz AI, Karlsruhe Institute of Technology, 76131, Karlsruhe, Germany

<sup>8</sup>High Energy Accelerator Research Organization (KEK), Tsukuba, Japan

<sup>9</sup>INFN - Sezione di Napoli, I-80126 Napoli, Italy

<sup>10</sup>INFN - Sezione di Padova, Padova, Italy

<sup>11</sup>INFN - Sezione di Pisa, I-56127 Pisa, Italy

<sup>12</sup>INFN - Sezione di Roma Tre, Roma, Italy

<sup>13</sup>INFN - Sezione di Torino, Torino, Italy

<sup>14</sup>INFN - Sezione di Trieste, Trieste, Italy

<sup>15</sup>Institut für Experimentelle Teilchenphysik, Karlsruher Institut für Technologie, Karlsruhe, Germany

<sup>16</sup>Institute of High Energy Physics, 1050, Vienna, Austria

<sup>17</sup>Jožef Stefan Institute, Ljubljana, Slovenia

<sup>18</sup>Ludwig Maximilians University, Munich, Germany

<sup>19</sup>Max-Planck-Institut für Physik, München, Germany

<sup>20</sup>Tel Aviv University, Tel Aviv, Israel

<sup>21</sup>University of Bonn, Bonn, Germany

<sup>22</sup>University of Melbourne, Melbourne, Australia

Received: date / Accepted: date

**Abstract** We present the novel implementation of a non-differentiable metric approximation and a corresponding loss-scheduling aimed at the search for new particles of unknown mass in high energy physics experiments. We call the loss-scheduling, based on the minimisation of a figure-of-merit related function typical of particle physics, a Punzi-loss function, and the neural network that utilises this loss function a Punzi-net. We show that the Punzi-net outperforms standard multivariate analysis techniques and generalises well to mass hypotheses for which it was not trained. This is achieved by training a single classifier that provides a coherent and optimal classification of all signal hypotheses over the whole search space. Our result constitutes a complementary approach to fully differentiable analyses in par-

ticle physics. We implemented this work using PyTorch and provide users full access to a public repository containing all the codes and a training example.

### 1 Introduction

The standard model (SM) of particle physics is the theoretical framework that describes fundamental interactions and the fundamental constituents of matter. Although successful in predicting phenomena, there is a general consensus that this framework is not a complete description of nature, and new physics (NP) has to exist. Searches for NP beyond the SM can be grouped into two main categories: searches for direct production and decays of new, unknown particles; and searches for deviations from the theoretical predic-

<sup>a</sup>e-mail: paul.feichtinger@oeaw.ac.at (corresponding author)

tions in precision measurements. When searching for new particles, for example, in a collider experiment, one of the main challenges is correctly reconstructing and identifying the new particles (the signal) and rejecting any (or most) contributions from potential background sources. This is a common problem referred to as event classification. A common approach to correctly classify a signal with respect to background uses Monte Carlo (MC) simulation to generate signal- and background-like event distributions. MC simulation can help find underlying features or patterns in the signal and the background distributions that allow one to disentangle the two (possibly) unambiguously. In the last decade, advanced data analysis methodologies, such as multivariate analysis (MVA) methods, have often improved analysis signal selection power, allowing for more precise analyses, usually performed in a shorter time. Typical MVA methods in use in the field of particle physics include, but are not limited to, decision trees, boosted decision trees (BDTs) [1], or shallow and deep neural networks (NNs) [2]. This paper focuses on the implementation of NNs. We propose and describe how to implement a new loss function, called Punzi-loss, based on the so-called Punzi figure-of-merit (FOM) [3]. We henceforth refer to a neural network trained with the Punzi-loss function as a Punzi-net. As a benchmark study to test the performance of the Punzi-loss and compare it to other techniques, we consider the search for invisible decays of the hypothetical  $Z'$  boson produced in the reaction  $e^+e^- \rightarrow \mu^+\mu^-Z'$  at the Belle II experiment [4, 5] at the SuperKEKB collider [6], based on MC simulations.

## 2 Neural networks

There exist many implementations of neural networks (e.g. convolutional neural networks (CNNs), transformers, etc.) that are used in various applications ranging from image classification in the case of CNNs to natural language processing with transformers. In this work, we focus on a fully connected feed-forward neural network for our experiments. We nonetheless emphasise that the concepts outlined in this work apply to all neural network implementations that use backpropagation.

A neural network comprises a collection of connected neurons. In a fully connected neural network, these constitute a series of layers in which each neuron is connected to all those in both the previous and subsequent layers. Each neuron describes a mathematical function that produces an output dependent on those input connections and a unique bias, defined as

$$a_j^l = \sigma \left( \sum_k w_{jk}^l a_k^{l-1} + b_j^l \right). \quad (1)$$

Here  $w_{jk}^l$  is the weighting of the connection to the  $k^{\text{th}}$  neuron in the previous  $(l-1)$  layer,  $b_j^l$  is the bias and  $\sigma$  is the *activation function*. A variety of different activation functions can be applied here, and most have specific traits that may be desirable depending on the application. Commonly used examples include sigmoid, rectified linear activation (ReLU), or hyperbolic tangent functions. The key requirements are that they are non-linear and have a derivative defined everywhere.

Using Eq. 1, a network of individual neurons is able to map input variables to some desired output. For this to be possible, however, the weight and bias parameters must be optimised. In the implementation we present here, this is done via *supervised training*, whereby training data,  $x$ , is passed to the network along with the set of corresponding labels,  $y$ . The actual output of the network,  $f(x) = \hat{y}$ , can then be compared with this desired output to measure how well it maps input data. This comparison is quantified by way of a loss function, a commonly used example of which is the Binary Cross Entropy loss,

$$L = -y \ln \hat{y} - (1-y) \ln(1-\hat{y}), \quad (2)$$

where  $y \in \{0, 1\}$  and  $\hat{y} \in [0, 1]$ . With this measure of the error, the training process becomes a minimisation problem: what weights and biases will minimise the loss and therefore provide the most effective network? This is solved by employing a method such as *gradient descent*, by which the parameters can be iteratively adjusted in the direction opposite that of the loss function's gradient,

$$w_{n+1} = w_n - \eta \frac{\delta L}{\delta w_n} \quad \text{and} \quad (3)$$

$$b_{n+1} = b_n - \eta \frac{\delta L}{\delta b_n}, \quad (4)$$

where  $\eta$  is the *learning rate*, the step size by which the parameters are adjusted at each iteration of the learning process. Each of these iteration steps constitutes a complete pass through a randomly sampled *batch* from the full training data set, a pass through the entirety of which is referred to as an *epoch*. The derivatives  $\frac{\delta L}{\delta w_j}$  and  $\frac{\delta L}{\delta b_j}$  are calculated through use of the *backpropagation* algorithm. This starts from the final layer and utilises the chain rule to incrementally calculate all derivatives through one full backward pass to the first layer. The key is carefully selecting a loss function whose minimum solves the given task while remaining differentiable across all possible neural network outputs.

## 3 Figure of merit

As highlighted in Section 1, one of the main challenges when performing a precision test of the SM, or in the search for

NP, is the fact that some background processes may mimic the signal and therefore contaminate the results. In the search for a new particle, for example, one is often performing a counting experiment which is described by the Poisson distribution. As discussed in [3], the number of events  $n$  in a counting experiment in the case of a background ( $B$ ) only hypothesis ( $H_B$ ), and in the case of a signal ( $S$ ) in the presence of the same background ( $H_{S+B}$ ) follows the Poisson distributions

$$p(n | H_B) = \frac{B^n e^{-B}}{n!} \quad (5)$$

and

$$p(n | H_{S+B}) = \frac{(S+B)^n e^{-(S+B)}}{n!}. \quad (6)$$

When MC simulations for both the signal and the background are available, it is possible to identify quantities or features in the data to separate and classify them correctly by applying specific selection criteria. This would eventually enable one to choose between the (null) background only and the signal plus background hypotheses. In general, however, applying some selection criteria to reduce the background contamination will also remove some of the signal. It is, therefore, fundamental to define some additional criteria that would indicate the best balance between reducing the background without compromising the signal. This is done via the implementation of a FOM. One can define  $S(t)$  and  $B(t)$  as the number of signal and background events that pass some selection criteria (e.g. particles having a momentum or energy larger than a specified threshold  $t$ ). In that case, standard FOMs used in particle physics are:

$$FOM = \frac{S(t)}{\sqrt{B(t)}} \text{ and} \quad (7)$$

$$FOM = \frac{S(t)}{\sqrt{S(t)+B(t)}}. \quad (8)$$

Neither of the above is usable in the search for new particles since the number of expected signal events depends on the cross-section of the process, and this is not known *a priori*. An alternative FOM for this specific case was proposed in [3], often referred to as the Punzi FOM after the author, and is now in widespread use. The Punzi FOM to maximise is the inverse of the minimum detectable cross-section  $\sigma_{\min}$ , which defines a sensitivity region for which the experiment will certainly give conclusive results: it either will be excluded, or a discovery will be claimed. An analytic formula for  $\sigma_{\min}$  is given by

$$\sigma_{\min}(t) = \frac{\frac{b^2}{2} + a\sqrt{B(t)} + \frac{b}{2}\sqrt{b^2 + 4a\sqrt{B(t)} + 4B(t)}}{\varepsilon(t) \cdot L}, \quad (9)$$

where  $L$  is the target luminosity,  $\varepsilon(t)$  is the signal efficiency and  $B(t)$  is the number of background events after the selection defined by  $t$ . The constants  $a$  and  $b$  are the number of sigmas corresponding to one-sided Gaussian tests at some predefined significance level,  $\alpha$  and  $\beta$ . Here  $\alpha$  is the probability of rejecting  $H_B$  when it is true (type I error), and  $\beta$  is the probability of not rejecting  $H_B$  when instead  $H_{S+B}$  is true (type II error). Since we are interested in cases where the signal hypothesis depends on a free parameter,  $\beta$  will change with this parameter. The sensitivity region for a given experiment is obtained for the parameter space that fulfils  $1 - \beta > CL$ , where  $CL$  is the confidence level for setting limits in case of no discovery. So for example when choosing  $\alpha$  to correspond to a significance of  $5\sigma$  and a desired confidence level of 90 %,  $a$  and  $b$  would be set to 5 and 1.28.

#### 4 Punzi-loss

We propose here a quantity approximating the Punzi FOM, appropriate for optimising neural networks for physics selections.

This loss function is based on the equation for the Punzi sensitivity region (Eq. 9). However, Eq. 9 can not be used directly because the number of background events  $B$  and the signal efficiency  $\varepsilon$  are discrete functions of the network parameters for any given fixed cut on the classifier output, whereas the loss function must be differentiable. We can build a differentiable function by replacing the fixed cut on the output with a sum over all events, weighted with the respective value of the output. If events classified as signal cluster around an output of 1 and events classified as background at 0, this quantity will closely approximate the original function. In Eq. 9 this weighting can be captured by performing the replacements

$$\varepsilon(t) \rightarrow \varepsilon(\mathbf{w}, \mathbf{b}) = \sum_{\mathbf{x}} \frac{y_i \cdot \hat{y}_i(\mathbf{w}, \mathbf{b}) \cdot s_{\text{sig}}}{N_{\text{gen}}} \quad \text{and} \quad (10)$$

$$B(t) \rightarrow B(\mathbf{w}, \mathbf{b}) = \sum_{\mathbf{x}} (1 - y_i) \cdot \hat{y}_i(\mathbf{w}, \mathbf{b}) \cdot s_{\text{bkg}}^i, \quad (11)$$

where the sum is over all training inputs  $\mathbf{x}$  and the index  $i$  denotes the  $i^{\text{th}}$  training event. The collection of weights and biases that constitute the free parameters of the network are denoted as  $\mathbf{w}$  and  $\mathbf{b}$ .  $N_{\text{gen}}$  is the total number of generated signal events,  $s_{\text{sig}}$  is a scale factor for the signal and  $s_{\text{bkg}}^i$  is a scale factor for the background, which can include a weight factor to scale the luminosity for the individual simulated background samples to the target luminosity. The scale factors can also include correction factors such as trigger efficiencies and should account for the sample size when only a subset of the generated data is used to compute the loss.

A similar approach of building a differentiable metric based on a FOM was taken by Elwood and Krücker [7], with a loss function based on the discovery significance.

The Punzi-loss function is given by the arithmetic mean of this continuous Punzi sensitivity calculated for all signal hypotheses ( $m_{Z'}$ ) that are used in training,

$$C_{\text{Punzi}} = \frac{1}{N_{Z'}} \sum_{m_{Z'}} \sigma_{\min}(\mathbf{w}, \mathbf{b}), \quad (12)$$

with  $N_{Z'}$  being the total number of hypotheses that were considered. Here we present an implementation in which all mass hypotheses are treated equally since they have equal weights in the calculation. However, one could introduce some weightings in the case of an analysis where the hypotheses do not have flat priors. Note that this loss function can no longer be calculated using single training events but is instead based on a set of training data.

To test the Punzi-loss function, we implemented a simple fully-connected network in PyTorch [8] with four input neurons, one output neuron, and two hidden layers with 8 and 4 neurons, respectively. The size of the net was determined empirically to give good results while keeping the network relatively small.<sup>1</sup>

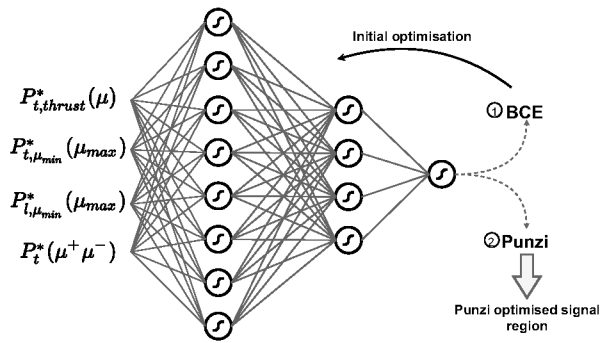
## 5 Training strategy

For the Punzi-loss training to converge, we found that the parameters of the network should already be initialised in a way that defines some separation between signal and background (similar to the loss scheduling scheme described in [9]). This can be achieved by pretraining the network using a conventional loss function and subsequently fine-tuning this through the use of the Punzi-loss function.

For the activation function of the neurons in the hidden layers, a hyperbolic tangent is used while the output neuron uses a sigmoid function. Before training, the input variables were scaled to lie between 0 and 1, and the network parameters were randomly initialised. A weighted binary cross-entropy (BCE) loss function was used for the pretraining. A weighting was attributed to the signal events such that their weighted sum was equal to the weighted sum of all background events. An outline of the network architecture is given in Figure 1.

Initially, using the BCE loss function, the network was trained with a batch size of 2048 and a learning rate (LR) of 1. When the loss did not decrease for 10 epochs, the LR was reduced by a factor of 0.5. The pretraining was stopped after 200 epochs. Training is then continued using the Punzi-loss function with  $a = 3$  and  $b = 1.28$ . Here we used a learning rate of 0.0001 and again reduced it upon plateauing. This

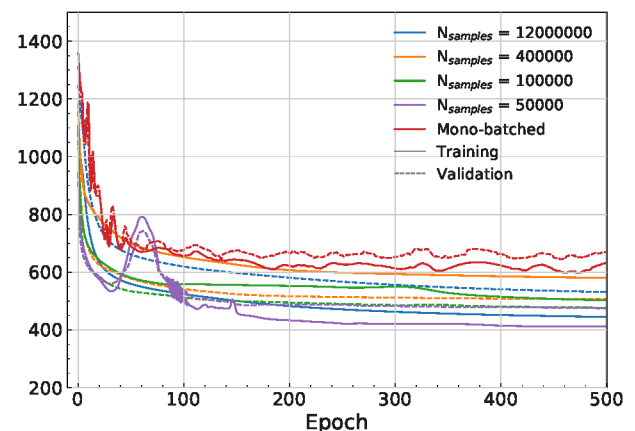
<sup>1</sup>The network size and architecture is not relevant for our approach.



**Fig. 1** An outline of the network architecture. The first training with the BCE loss function was used to set the weights and biases of the net for the second training with the custom loss function based on the Punzi FOM.

training was stopped after 1000 epochs. The gradient descent algorithm was used for optimisation for both of these trainings. All hyperparameters were optimised to give the best results for the training methods. One particularly important hyperparameter is the batch size, the variation of which presents some unique aspects of the Punzi-loss function that must be considered.

Due to the nature of the Punzi-loss function concerning the optimisation for a desired luminosity, utilising training data in excess of this requires the addition of weightings in the loss calculation. The background data used for training contained  $1000 \text{ fb}^{-1}$ ,  $450 \text{ fb}^{-1}$  and  $3000 \text{ fb}^{-1}$  worth of events from three main background processes; however, in this study, we wish to optimise the classifier for just  $50 \text{ fb}^{-1}$  of real-world data. Naturally, it is preferred that all background data is utilised, and thus we introduce a background scaling factors of 0.05, 0.111 and 0.0167 respectively. Additionally, dividing the training data into batches brings about the additional requirement of multiplying both  $\epsilon(\mathbf{w}, \mathbf{b})$  and  $B(\mathbf{w}, \mathbf{b})$  by the number of batches used.



**Fig. 2** Evolution of Punzi-loss during training with batch sizes of  $5 \times 10^4$ ,  $1 \times 10^5$ ,  $4 \times 10^5$  and  $1.2 \times 10^6$  and mono-batched.

Fig. 2 shows the evolution of loss during training with batch sizes of  $5 \times 10^4$ ,  $1 \times 10^5$ ,  $4 \times 10^5$  and  $1.2 \times 10^6$ , and also mono-batched (where the whole data set is passed as a single batch). These correspond to batch sizes of between roughly 0.5% and 13% of the total dataset. A batch size of  $1 \times 10^5$  was chosen for the following experiment. We note that small batch sizes bring a degree of instability to the loss, as can be seen in the line representing a batch size of  $5 \times 10^4$  in fig. 2. It was found that batches smaller than those shown in Fig. 2 led to increasing loss values over the training, with a batch size of  $1 \times 10^4$  leading to training regularly failing with the Punzi-loss increasing and plateauing at a value above the initial loss. This can be understood as a result of the limited number of signal events present in any given batch of small size, leading to large statistical fluctuations in the calculated loss values. This lower limit is, of course, study dependent. Similarly, we note large instability in the mono-batched case. The batching introduces an additional stochastic component during training, making the training more robust and helping the algorithm escape local minima.

## 6 Results

In this section, we present the results of utilising a Punzi-net in a search for  $e^+e^- \rightarrow \mu^+\mu^-Z'$  signals amongst various common backgrounds found in  $e^+e^-$  collider experiments. At the Belle II experiment, this search was performed with the commissioning data for the specific case of invisible decays of the  $Z'$  boson [10], a final state in which only the two muons produced by the electron-positron annihilation can be reconstructed. Therefore, all information about the production and decay of the  $Z'$  boson is to be inferred by the two-muon system. The signal events are generated with MadGraph 5 [11] for a range of candidate  $Z'$  masses, spanning  $0.1 \text{ GeV}/c^2$  to  $8.9 \text{ GeV}/c^2$  in steps of  $0.1 \text{ GeV}/c^2$  with 20000 events produced at each. Additionally, MC samples for the background process  $e^+e^- \rightarrow e^+e^-\mu^+\mu^-$ ,  $e^+e^- \rightarrow \tau^+\tau^-$  and  $e^+e^- \rightarrow \mu^+\mu^-(\gamma)$  corresponding to  $1000 \text{ fb}^{-1}$ ,  $450 \text{ fb}^{-1}$  and  $3000 \text{ fb}^{-1}$  respectively were used, since these can mimic the signal. The simulation and reconstruction of the events were done using GEANT4 [12], and the Belle II Analysis Software Framework [13]. The analysis is carried out via the search for a peak in the distribution of the squared mass recoiling against the two-muon system. An excess of entries beyond that of the expected background at a given mass would indicate the presence of such a  $Z'$  particle of that mass. This distribution is divided into (potentially overlapping) bins with bin widths corresponding to  $\pm 2\sigma$  of the fitted  $Z'$  signal distributions.

During both the initial BCE and subsequent Punzi-loss training, only every second generated  $Z'$  mass was used. For the calculation of  $\sigma_{\min}$  in Eq. 12 only signal and background

variable	description
$p_{t,\text{thrust}}^*(\mu)$	The transverse momentum component of the muons with respect to the thrust axis.
$p_{t,\mu_{\min}}^*(\mu_{\max})$	The transverse momentum component of the higher energetic muon with respect to the lower energetic muon.
$p_{l,\mu_{\min}}^*(\mu_{\max})$	The longitudinal momentum component of the higher energetic muon with respect to the lower energetic muon.
$p_t^*(\mu^+\mu^-)$	The transverse momentum of the dimuon system.

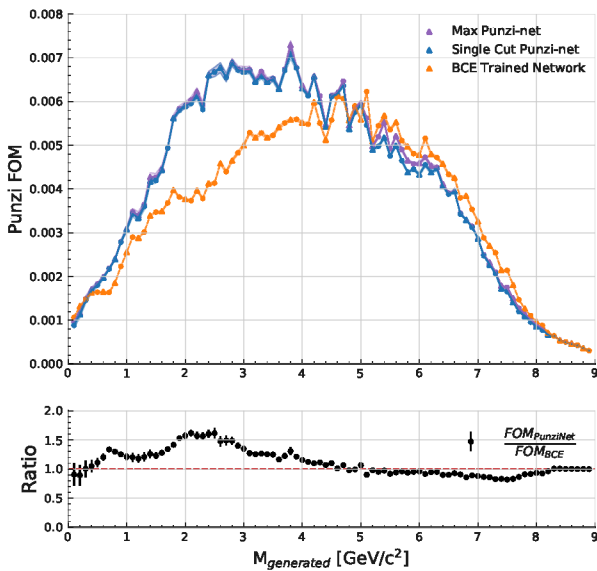
**Table 1** The most important features found after training BDTs with many variables. All variables are computed in the centre-of-mass system of the  $e^+e^-$  collisions. These features are used for training the NN.

events that lie within the respective  $\pm 2\sigma$  mass windows are considered, using only signal events that were generated for the corresponding mass. Thus, events that are not contained in any of the mass windows of the used signal samples are not taken for the training. This results in a data set of approximately 9 million total events, of which  $\sim 2.5\%$  are signal and the rest background. This is then split by using 80% of the events for training and the remaining 20% for validation. The unused signal hypotheses are utilised for validation and to check the trained networks ability to generalise to signals unseen in training. The network was trained with four carefully selected features related to the event kinematics that showed a good discrimination power when using a boosted decision tree classifier. These features are described in Tab. 1. A more detailed description of these features and the analysis can be found in [14].

The resulting maximum achievable Punzi FOM spanning the range of generated  $Z'$  signals is shown in Fig. 3. Included in this figure are the Punzi-net along with the BCE pretrained network. These values are calculated using the background data contained within the  $\pm 2\sigma$  bin around each generated mass point. The maximum achievable Punzi FOM in each bin is found using the cut to the network output that provides the highest FOM for that respective bin. In addition to this, the resulting Punzi FOM after applying a single cut value to the output of the Punzi-net across the full recoil mass spectrum is shown. The plot shows the average result found over ten independently trained networks, along with the associated standard error. This serves to demonstrate that not only can the Punzi-loss function produce better FOMs, but can do so consistently. The Punzi-loss function shows greater effectiveness through the lower half of the recoil mass spectrum, providing clear improvements to the FOM below approximately  $5 \text{ GeV}/c^2$ . For mass hypotheses above this point, there is some slight degradation of the maximum achievable FOM with the Punzi-net. It is important to note that the single cut applied to the Punzi-net output can still

provide a FOM near to that of the maximum achievable with the BCE trained network in this region.

This means that even when compared to an optimal varied cut applied to the BCE network output, interpolated over the recoil mass spectrum, the Punzi-net provides comparable or even improved results. As discussed previously, this cut interpolation can lead to discontinuities in the final recoil mass distribution. So the ability to achieve comparable results with a single cut to the Punzi-net output is much preferable, meaning that even in the higher recoil mass region where the BCE network appears to outperform the Punzi-net, it may not be a preferable method due to the need for cut interpolation.

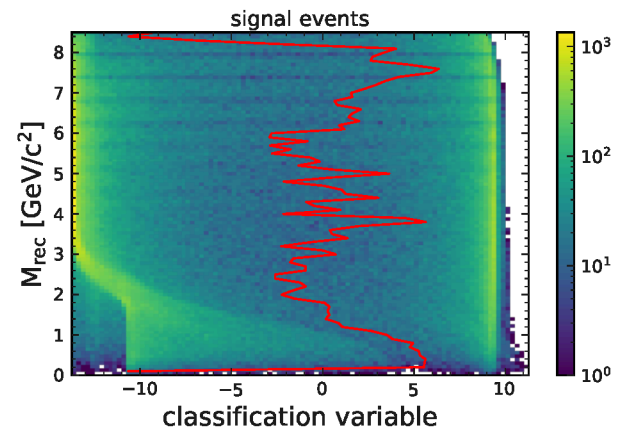


**Fig. 3** The average maximum Punzi FOM achievable in each bin across range of generated  $Z'$  signals, with standard error spread taken from 10 independently trained networks. Triangles indicates those masses which were left out of training while circles indicates those used.

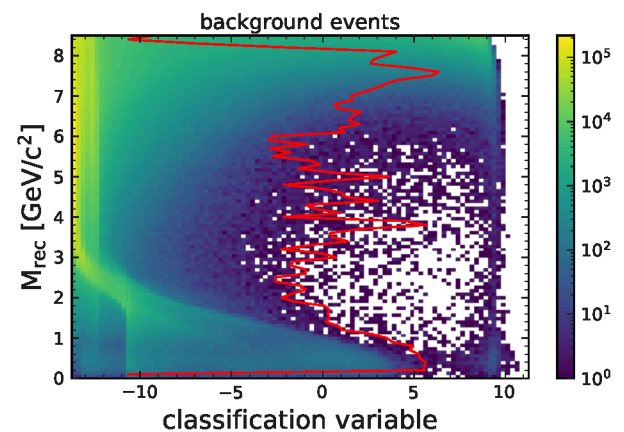
The generated  $Z'$  masses used for training the network are shown with circles, and those not used in training are shown with triangles. The figures show little to no difference in the network's ability between these training and validation masses, indicating that the model generalises well to unseen signals. In the region between approximately  $4.5 \text{ GeV}/c^2$  and  $5.5 \text{ GeV}/c^2$  some dependence on whether or not a mass was used in training does appear. This could be combated by generating a larger set of  $Z'$  signals covering more mass points in that region.

Fig. 4 shows the output of the Punzi-trained network for all signal events and Fig. 5 shows the same for all background events. Here the variable on the x-axis shows the

NN output before applying the last sigmoid activation function to resolve the distribution of events better. The y-axis corresponds to the reconstructed recoil mass ( $M_{\text{rec}}$ ), which discriminates between the different signal hypotheses. The classified signal and background events are separated into two clusters, corresponding to an output of 0 and 1. The overlaid line shows the cut value that would give the maximum achievable Punzi FOM for each  $Z'$  mass. The line separates the two clusters, showing that the training using the approximations in Eq. 10 and 11 worked as expected.



**Fig. 4** The output distribution of all signal events using the Punzi-loss trained NN, overlaid with the optimal decision threshold for each signal hypothesis. The classification variable shows the NN output before applying the last sigmoid function in order to better see the separation. The optimal cut value can be replaced by a uniform cut without any significant difference in the resulting selection.

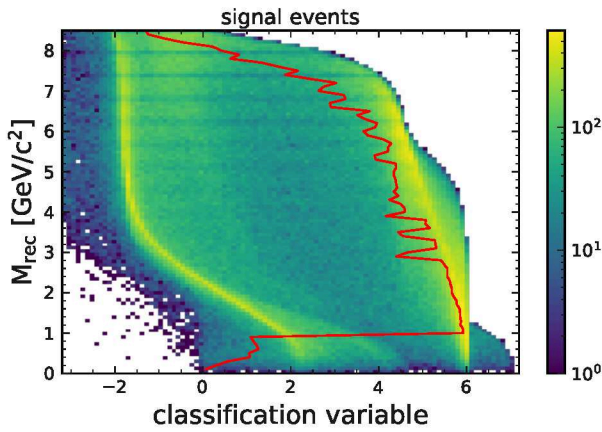


**Fig. 5** The output distribution of all background events using the Punzi-net, overlaid with the optimal decision threshold for each signal hypothesis.

The events are separated so that when only the events classified as signal are selected (for example, by applying a

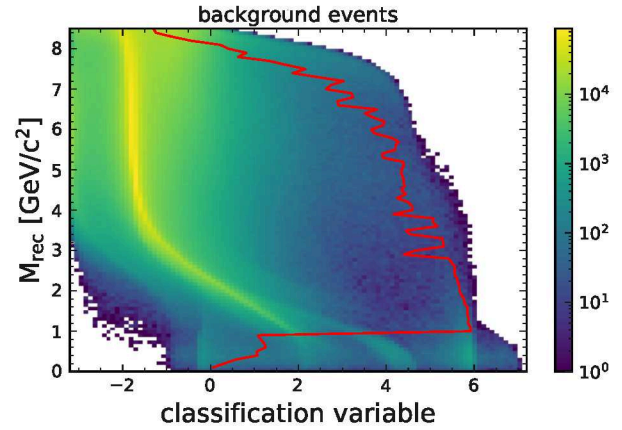
cut at a NN output of 0.5), this gives the optimal Punzi FOM for the whole mass range. This is a significant advantage for an analysis since no additional interpolation between output values is required, which can introduce discontinuities in the final recoil mass distribution. Additionally, since the selection generalises to all signal hypotheses, it gives also the best possible FOM for a signal in-between trained masses, which would otherwise have non-optimal results.

For comparison we also show the output distribution of the BCE pretrained network in Fig. 6 for the signal and Fig. 7 for the background. Again, the optimal cut that gives the highest Punzi FOM at each mass hypothesis is shown. While a separation between signal and background is also achieved here, the division is not as pronounced as with the Punzi-net and the best cut value varies significantly with the mass.

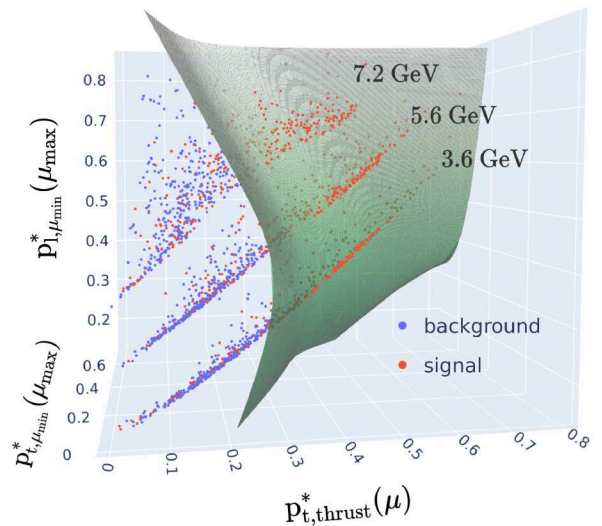


**Fig. 6** The output distribution of all signal events after the BCE pretraining, overlaid with the optimal decision threshold for each signal hypothesis. The optimal cut value varies significantly across the mass spectrum.

An understanding of why the model successfully generalises, and one network can be utilised for the full squared recoil mass spectrum, can be inferred from Fig. 8, which shows a 3D scatter plot of the  $p_{t,\text{thrust}}^*(\mu)$ ,  $p_{t,\mu_{\min}}^*(\mu_{\max})$  and  $p_{t,\mu_{\min}}^*(\mu_{\max})$  variables (after being normalised to values between 0 and 1) for three of the mass bins at a region of  $p_t^*(\mu^+\mu^-) = (2.2 \pm 0.5) \text{ GeV}/c$ . The green plane is the chosen signal/background classification boundary obtained with a single cut. One can see the masses describing three respective planes in the parameter space which occupy distinct regions. This partitioning allows the network to adapt between the different mass regions and so negates any need for multiple classifiers for different regions.



**Fig. 7** The output distribution of all background events after the BCE pretraining, overlaid with the optimal decision threshold for each signal hypothesis.



**Fig. 8** A 3D scatter plot showing the input space of the NN with  $p_t^*(\mu^+\mu^-)$  fixed around  $2.2 \text{ GeV}/c$ . The separation boundary defined by the final selection (green sheet) separates the planes corresponding to different recoil masses in a way that optimises the selection for all signal hypotheses.

e

## 7 Conclusions

In this work, we have demonstrated that it is possible to implement a non-differentiable metric approximation and a corresponding loss-scheduling, combining the approach of particle physics and that of machine learning. Our proposed method applies to the search for new particles with unknown parameters in high energy physics experiments.

We designed a new loss function directly related to the Punzi figure-of-merit, intended to be calculated on a set of training events at once. Training instabilities could be solved by a batched training that helps in the algorithm's conver-

gence. We showed that this loss function can be used to achieve an optimal selection for all signal hypotheses with a single cut on the classifier output, also achieving overall better performance than standard methods. The main advantage of this method is that it simplifies the analysis since it does not require any further optimisation of the selection or training of multiple classifiers for subsets of signal hypotheses. We implemented the Punzi-loss in the training of a simple neural network and made the code publicly available [15]. However, the method is general and not restricted to the use of the presented architecture.

A universal approach to this problem would be to construct a fully differentiable analysis pipeline that can optimise any utility function, which is an active area of research [16]. Such analysis frameworks can also take into account systematic effects during the optimisation of the signal selection [17, 18]. Another interesting approach to incorporate systematic effects is to introduce an adversarial discriminator in addition to a classifier. This provides a handle for robust inference by learning a pivotal quantity - a predictive function that is insensitive against the unknown values of the nuisance parameters that model the systematic effects. [19]

**Acknowledgements** The authors would like to thank the Belle II Collaboration and Belle II software group for useful discussions and suggestions on how to improve this work. P. Feichtinger, H. Haigh and G. Inguglia would like to acknowledge funding received under the Horizon 2020 framework of the European Research Council, ERC StG Nr. 947006 *InterLeptons*, and under the FWF standalone framework with grant Nr. P31361 *Searches for dark matter and dark forces at Belle II*. J. Kahn's work is supported by the Helmholtz Association Initiative and Networking Fund under the Helmholtz AI platform grant. F. Meier work is supported by the US Department of Energy.

## References

1. T. Keck, *Comput. Softw. Big Sci.* **1**(1), 2 (2017). doi:10.1007/s41781-017-0002-8
2. K. Albertsson, et al. Machine learning in high energy physics community white paper. arXiv:1807.02876 (2019)
3. G. Punzi, in *Statistical problems in particle physics, astrophysics and cosmology. Proceedings, Conference, PHYSTAT 2003, Stanford, USA*, vol. C030908, ed. by L. Lyons, R.P. Mount, R. Reitmeyer (2003)
4. T. Abe, et al. Belle II technical design report. arXiv:1011.0352 (2010)
5. E. Kou, et al., *Prog. Theor. Exp. Phys.* **2019**(12) (2019). doi:10.1093/ptep/ptz106
6. Y. Ohnishi, et al., *Prog. Theor. Exp. Phys.* **2013**(3) (2013). doi:10.1093/ptep/pts083
7. A. Elwood, D. Krücker. Direct optimisation of the discovery significance when training neural networks to search for new physics in particle colliders. arXiv:1806.00322 (2018)
8. A. Paszke, et al., in *Advances in Neural Information Processing Systems 32*, ed. by H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, R. Garnett (Curran Associates, Inc., 2019), pp. 8024–8035
9. O. Taubert, M. Götz, A. Schug, A. Streit, in *19th IEEE International Conference on Machine Learning and Applications (ICMLA)* (2020), pp. 426–431. doi:10.1109/ICMLA51294.2020.00073
10. I. Adachi, et al., *Phys. Rev. Lett.* **124**, 141801 (2020). doi:10.1103/PhysRevLett.124.141801
11. J. Alwall, et al., *J. High Energy Phys.* **2014**(7), 79 (2014). doi:10.1007/JHEP07(2014)079
12. S. Agostinelli, et al., *Nucl. Instrum. Methods Phys. Res. A: Accel. Spectrom. Detect. Assoc. Equip.* **506**(3), 250 (2003). doi:10.1016/S0168-9002(03)01368-8
13. T. Kuhr, et al., *Comput. Softw. Big Sci.* **3**(1), 1 (2018). doi:10.1007/s41781-018-0017-9
14. P. Feichtinger, Search for an invisibly decaying  $Z'$  boson and study of particle identification at the Belle II experiment. Master's thesis, TU Wien (2021). doi:10.34726/hss.2021.84843
15. Punzi-loss repository. [github.com/feichtip/punzinet](https://github.com/feichtip/punzinet)
16. A.G. Baydin, et al., *Nucl. Phys. News* **31**(1), 25 (2021). doi:10.1080/10619127.2021.1881364
17. P. de Castro, T. Dorigo, *Comput. Phys. Commun.* **244**, 170 (2019). doi:10.1016/j.cpc.2019.06.007
18. S. Wunsch, S. Jörger, R. Wolf, G. Quast, *Comput. Softw. Big Sci.* **5**(1), 4 (2021). doi:10.1007/s41781-020-00049-5
19. G. Louppe, M. Kagan, K. Cranmer. Learning to pivot with adversarial networks. arXiv:1611.01046 (2017)