

Applications of Machine Learning at BESIII

Beijiang Liu^{1,*}, Xian Xiong^{1,**}, Guoyi Hou¹, Shiming Song², and Lin Shen²

¹Institute of High Energy Physics, Chinese Academy of Sciences

²Sichuan University

Abstract. BESIII is an experiment at the high precision frontier of hadron physics in τ -charm region. Machine learning techniques have been used to improve the performance of BESIII software. In this proceeding, we present novel approaches with XGBoost for multi-dimensional distribution reweighting, muon identification and cluster reconstruction for CGEM (Cylindrical Gas Electron Multiplier) inner tracker.

1 Introduction

The BESIII detector is a magnetic spectrometer [1] located at the Beijing Electron Positron Collider (BEPCII) [2] which is a double ring e^+e^- collider running at the center of mass energies between 2.0 and 4.6 GeV and has reached a peak luminosity of $1 \times 10^{33} \text{cm}^{-2}\text{s}^{-1}$ at $\sqrt{s} = 3770$ MeV. The BESIII experiment has collected the world's largest data samples of J/ψ , $\psi(3686)$ and $\psi(3770)$ decays as well as data in the energy region above 4 GeV. These data samples with unrepresented precision are being used to make a variety of important and unique studies [3]. Machine learning (ML) techniques have been employed to improve the performance of BESIII software. Novel approaches with XGBoost (eXtreme Gradient Boosting) [4] for muon identification, multi-dimensional distribution reweighting and cluster reconstruction for the Cylindrical Gas Electron Multiplier inner tracker (CGEM-IT) are presented.

2 A new approach for muon identification

The BESIII detector has a geometrical acceptance of 93% of the full solid angle. The cylindrical core of the BESIII detector consists of a helium-based multilayer drift chamber (MDC), a plastic scintillator time-of-flight system (TOF), a CsI(Tl) electromagnetic calorimeter (EMC) and a muon chamber system (MUC) with layers of resistive plate chambers in the iron return yoke of a 1 T superconducting solenoid. Particle identification (PID) for charged tracks combines the measurements of the energy loss in the MDC (dE/dx), the time-of-flight information from the TOF and the information from EMC and MUC and forms a likelihood $\mathcal{L}(h)$ for each particle ($h = e, \mu, \pi, K, p$) hypothesis using $\mathcal{L} = \mathcal{L}_{dE/dx} \cdot \mathcal{L}_{TOF} \cdot \mathcal{L}_{EMC} \cdot \mathcal{L}_{MUC}$, where $\mathcal{L}_{dE/dx(TOF)}$ is calculated from χ^2 of particle hypothesis and $\mathcal{L}_{EMC(MUC)}$ is a normalized output of a shallow neural network of EMC(MUC). The discrimination of μ / π is crucial for many

*e-mail: liubj@ihep.ac.cn

**e-mail: xiongx@ihep.ac.cn

of the analyses. However the identification of μ is very challenging – μ and π are difficult to be discriminated with dE/dx and TOF because their masses are very close.

A nesting architecture with XGBoost classifiers for μ identification is proposed as shown in figure 1. We use two classifiers with all the reconstructed information of EMC and MUC as inputs, respectively. The outputs of the two classifiers together with $\chi^2_{dE/dx}$ and χ^2_{TOF} are submitted to another classifier for combination. Since PID of hadrons only uses dE/dx and TOF information and PID for electrons used dE/dx, TOF and EMC, user can easily choose which subdetectors to use with the architecture shown in figure 1. In this proceeding, the classifier is trained on a full simulation of μ and π sample uniformly distributed with momentum from 0.1 to 1.4 GeV, $\cos\theta$ (polar angle) from -0.8 to 0.8. Figure 2 (a) shows the result of ROC curve and figure 2 (b) shows the AUC values varying with particle momentum (one of the input variables of the classifier). The comparisons indicate the new approach with ML has a better performance than the default muon identification used in BESIII. The performance drops around 0.4 GeV because there is a cut-off of the MUC for those low momentum particles cannot reach the muon counter.

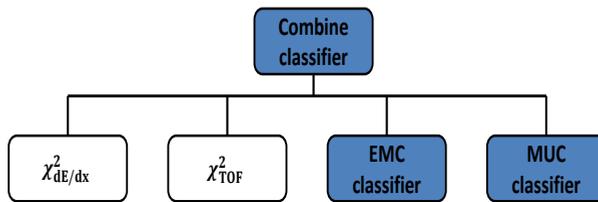


Figure 1. A nesting architecture with XGBoost classifiers for μ identification.

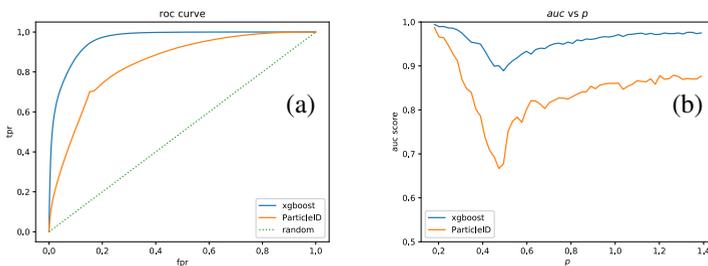


Figure 2. Comparison between new approach with XGBoost (blue curve) and the default PID (yellow curve) (a) ROC curve (b) AUC values varying with particle momentum.

3 Multi-dimensional reweighting with XGBoost

It is critical for physics analysis to model data with Monte Carlo (MC) simulation, e.g., for efficiency calculation and background estimation. In the energy regime of non-perturbative

Table 1. Detection efficiency calculated with different samples.

Sample	Efficiency
Pseudo data	(68.4±0.2)%
PHSP MC	(75.7±0.1)%
Reweighting	(67.9±1.1)%

QCD, resonances are plentiful in experimental data of BESIII leading to intricate interference patterns. Generic MC models can not describe the data in detail. Amplitude analysis [5] usually needs to extract the properties of resonances and model the data. For the channel of which the results of amplitude analysis are not available yet, multi-dimensional reweighting is an easy way to create a “data-like” MC . Two methods for reweighting with ML techniques [6, 7] have been proposed. Utilize the approach [6], we trained a XGBoost classifier to discriminate data and MC, which can provide probabilities $p_{data}(x)$ and $p_{MC}(x)$. The probabilities of an event x belongs to data or MC can be used to estimate the reweighting factor $p_{data}(x)/p_{MC}(x)$.

A MC sample (as “pseudo data”) is generated of $J/\psi \rightarrow N^*\bar{n} + c.c. \rightarrow p\pi^-\bar{n} + c.c.$ including a set of intermediate $p\pi$ resonances. Reweighting is applied to a phase-space-distributed MC (PHSP) with a 10-fold validation. Figure 3 shows the comparison of invariant mass distributions of $\pi^-\bar{n}$, $p\bar{n}$ and π^-p of pseudo data and MC before and after reweighting. χ^2/bin for $M(\pi^-\bar{n})$, $M(p\bar{n})$ and $M(\pi^-p)$ are reduced from 57.7, 15.3 and 18.2 to 2.8, 1.9 and 2.1 ,respectively. Tab. 1 shows the detection efficiencies calculated with different samples. The result of reweighed MC is more consistent with that of pseudo data.

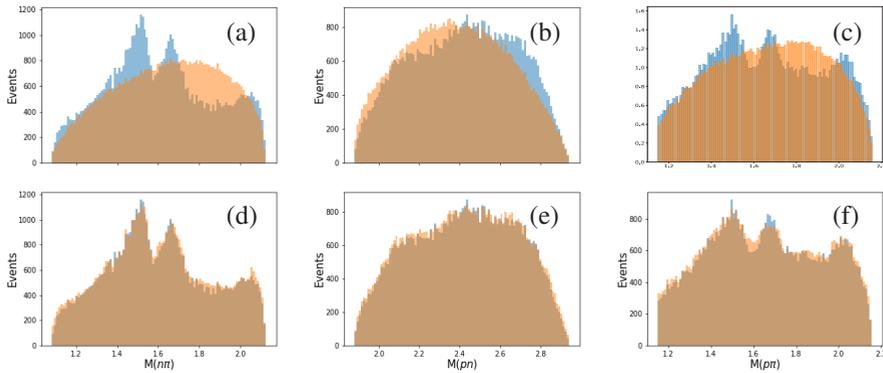


Figure 3. Comparison of pseudo data (blue) and MC(yellow) distributions before(a-c) and after (d-f) using the XGBoost reweighter.

4 Cluster reconstruction of cylindrical GEM inner tracker

BESIII will upgrade its inner tracker with 3 layers of cylindrical triple-GEMs in 2019 due to the aging effects of inner drift chamber. Cluster reconstruction is to measure the position of the ionizing particle in the drift cathode layer with the readouts from the anode strips which is the first step of track reconstruction for CGEM. There are two methods for cluster reconstruction of CGEM inner tracker [8]. The charge centroid method calculates the weighted average position of the anode strips with their charge (Q). The time-based method is based on the time

measurement (T) using the drift gap as a “micro time projection chamber” (micro-TPC) [9]. To improve the position resolution, the results of the two methods can be further combined according to their resolutions and correlations. However, the correlations between resolution and incident angle are quite complicated and difficult to handle.

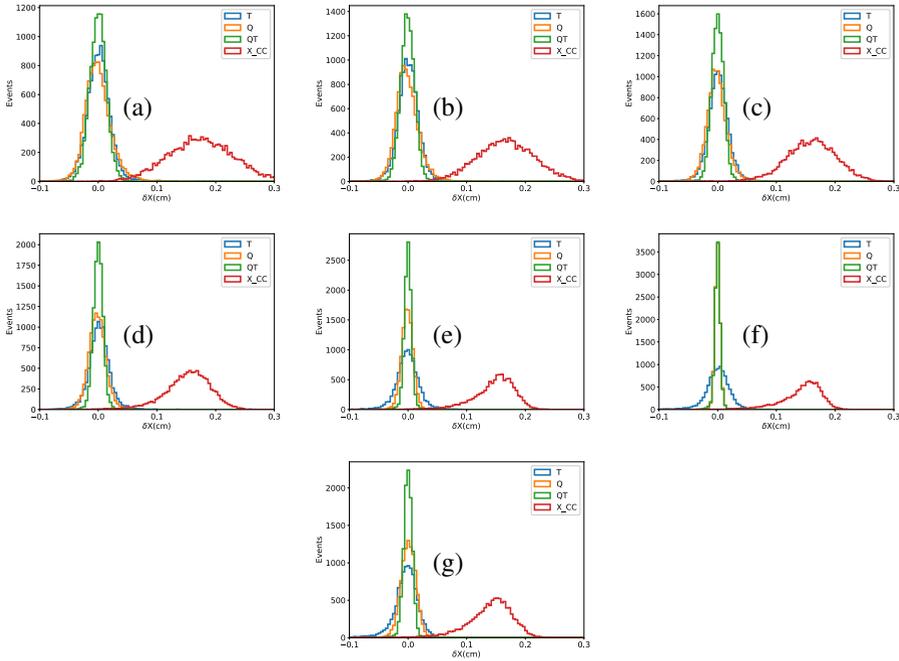


Figure 4. $\delta X(X_{reconstructed} - X_{truth})$ distributions under different circumstances, blue curve for T input only, yellow curve for Q input only, green curve for Q, T input together, red curve for charge centroid results, with different incident angle (a) -30° , (b) -20° , (c) -10° , (d) 0° , (e) 10° , (f) 20° , (g) 30° . The results of charge centroid method are on the anode readout plane, where a shift is induced by the Lorentz angle.

We propose a ML method based on XGBoost regressor to reconstruct the initial ionizing particle position with the readouts of Q and T from the fired strips. A simulation with a standalone digitization code, based on GARFIELD [10], is used to generate the event with 1 T magnetic field, incident angle between -30° to 30° and one layer of planar Triple-GEM. The results compared with the charge centroid method are shown in figure 4 and figure 5. The results show the dependency between incident angle and resolution is properly reflected with the Q or T input alone. The resolution of ML Q or T is significantly better than that of the charge centroid method. The information of Q and T can be combined by ML and gives a further improved resolution.

5 Summary

In this proceeding, we present three applications of ML techniques at BESIII and the results are promising. In the future, we plan to investigate the application of ML to further improve the performance of BESIII software, e.g., the tracking of low momentum charged particles, the tracking with high background rates, *etc.*

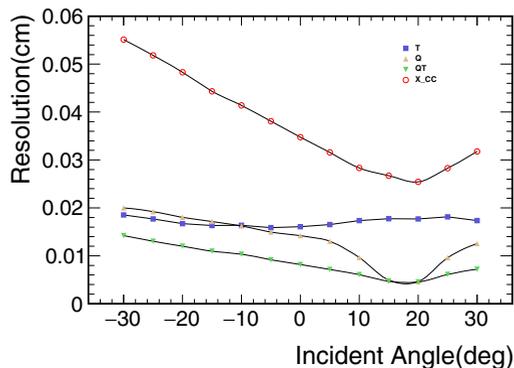


Figure 5. Resolution curve along with incident angle, blue curve for T input only, yellow curve for Q input only, green curve for Q, T input together, red curve for charge centroid results.

This work is supported in part by the CAS Large-Scale Scientific Facility Program; CAS Key Research Program of Frontier Sciences under Contract No. QYZDJ-SSW-SLH040; Joint Large-Scale Scientific Facility Funds of the NSFC and CAS under Contract No. U1732103.

References

- [1] M. Ablikim *et al.* (BESIII Collaboration), Nucl. Instrum. Meth. A **614**, 345 (2010).
- [2] C. H. Yu *et al.*, Proceedings of IPAC2016, Busan, Korea, 2016, doi:10.18429/JACoW-IPAC2016-TUYA01.
- [3] D. M. Asner *et al.*, Int. J. Mod. Phys. A **24**, S1-794 (2009)
- [4] T. Chen , C. Guestrin, Proceedings of the 22Nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (785-794) (2016)
- [5] M. Ablikim *et al.* (BESIII Collaboration), Phys. Rev. D **88** 112007, (2013)
- [6] D. Martschei *et al.*, Journal of Physics: Conference Series, 368(1), 012028 (2012)
- [7] A. Rogozhnikov, Journal of Physics: Conference Series, 762(1), 012036 (2016)
- [8] R. Farinelli *et al.*, arXiv:1807.00500 (2018)
- [9] T. Alexopoulos *et al.*, Nucl. Instrum. Meth. A **617**, 161 (2010).
- [10] R. Farinelli, L. Lavezzi *et al.*, arXiv:1807.01210 (2018)