FINAL REPORT ON THE START PROGRAMME

# Using Machine Learning for Particle Identification in MPD

Supervisor:
**Dr. Artem A. Korobitsin**

Student:
**Tolkachev Grigorii**, Russia
National Research Nuclear University MEPhI

**Participation period:**
August 1 - September 10
Dubna, 2022

# Abstract

Currently, various tasks in particle physics can be solved using Machine Learning techniques. For instance, particle reconstruction (clustering), identification (classification), and energy or direction measurement (regression) in calorimeters and tracking devices. Multivariate classification techniques are typically used to combine the all available response off all involved detectors into variables, called particle identification classifiers. Thus, study was made of the optimal MLP classifier selection for particle identification.

# Contents

# Introduction

The diversity of the data in the field of relativistic heavy-ion collisions, obtained by experiments at the SIS[1], AGS[2], SPS[3], RHIC[4] and LHC[5], is already quite large and impressive. In recent years, the STAR[6] program has produced a wealth of results to describe the bulk properties of the medium created in Au+Au reactions for $\sqrt{S_{NN}}$=7.7, 11.5, 14.6, 19.6, 27, 39, 62.4 and 200 GeV [7] by measuring several observables at mid rapidity. The Nuclotron-base Ion Collider fAcility (NICA)[8] is a major science project realised by the Joint Institute for Nuclear Research (JINR). Its aims to investigate phase diagram of QCD matter in the region of maximum baryonic density by studying (heavy-)ion collisions in the energy range from 2.5 to 11 GeV/nucleon.

The Multi-Purpose Detector (MPD)[9] program planned with the high intensity NICA beams promises to provide deeper knowledge of the dynamics of hadronic interactions and multiparticle production mechanisms at maximum baryon density, determine the nature of the phase transition between the deconfined and hadronic matter and search for the critical point. The new experimental program at the NICA-MPD will fill a niche in the energy scale, which is not yet fully explored, and the results will bring about a deeper insight into hadron dynamics and multiparticle production in the high baryon density domain.

Particle IDentification (PID) at the MPD experiment relies on several sub-detectors such as Time-Of-Flight (TOF) and Time Projection Chamber (TPC). There are many technique for PID such as, the n-Sigma or the Bayesian approach. In this paper, PID classifier based on Multilayer Perceptron (MLP) are presented.

# Chapter 1

# Introduction to the detector

## 1.1  The NICA accelerator complex

The NICA accelerator facility (see Fig. 1.1.1) consists of an injection, a booster, the upgraded Nuclotron accelerator and two storage rings of a collider [10]. The injector, which is designed and fabricated in cooperation with German
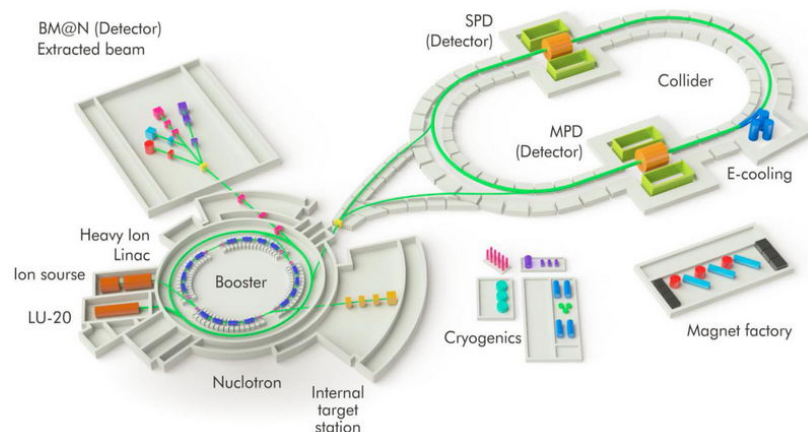


Figure 1.1.1 – A schematic view of the NICA accelerator complex.

and Russian research companies, consists of a heavy-ion source followed by a heavy-ion linac (HILAC) and a RFQ fore-injector. The booster synchrotron with its magnetic ring to be located inside of the existing yoke of the Dubna Synchrophasotron will accelerate ions up to 600 MeV/nucleon energy. The most challenging characteristics of the booster will be ultrahigh vacuum and electron cooling. The upgraded Nuclotron should provide p, d, and heavy ion beams with the maximum energy per nucleon of 5.8 GeV for A/Z=0.5 specie and 4.5 GeV for $^{197}Au$ nuclei. The initial luminosity is planned to be at least $10^{24}cm^{-2}s^{-1}$ with a relatively quick increase to at least $10^{25}cm^{-2}s^{-1}$. The

design luminosity goal for NICA with all components, such as an Electron Cooling System and the full set of RF cavities, is $10^{27} cm^{-2} s^{-1}$. Symmetric collisions of heavy ions will be performed in the initial stages of the NICA operation. Several types of ions are under consideration. These include $^{197}Au$ ions, which were used in previous and ongoing experiments at RHIC; $^{208}Pb$ ions, which were used for extensive data runs at SPS; and $^{209}Bi$ ions, which are very similar to $Pb$ ions, but provide more reliable operation of the NICA injection and acceleration complexes during the commissioning and first running phases. For heavy ions, such as $Au$ and $Bi$, the kinetic energy of the beam provided by the Nuclotron will be in the range from 2.5 to 3.8 GeV per nucleon. In the first year of operation, additional acceleration of the beams in the NICA collider is not foreseen. Therefore the initial collision energy $\sqrt{s_{NN}}$ may vary from 7 up to 9.46 GeV, with the collision energy of 9.2 GeV being preferred, so that results can be compared with those of RHIC-STAR that collected data at the same energy. Delivering $Au + Au$ collisions at $\sqrt{s_{NN}}$ up to 11 GeV remains the key goal of the NICA project that will be accomplished after the initial commissioning stage of operation. NICA will also provide the beams of polarised protons and deuterons up to centre of mass energy of 27 GeV with luminosity of $10^{32} cm^{-2} s^{-1}$ . Two collider rings of about 503 m circumference each are based on double-aperture superconductive (SC) magnets which are designed and manufactured in JINR Laboratory. The maximum field of SC dipole magnets is of 1.8 T. For luminosity preservation in the collider, both an electron and stochastic cooling systems will be constructed.

## 1.2   Multi-Purpose Detector (MPD)

The main physics experiment at NICA is the MPD, which will be operating at the collider. In 2018 an international scientific collaboration of MPD has been established. Currently, it is composed of 42 institutes from 12 countries, as well as JINR as a host institution. The collaboration will be operating the MPD apparatus, which is shown schematically in Fig. 1.2.1. Its designated position is the MPD Hall, which is located at the northern straight section of the NICA collider. Since 2020 the building has been available for MPD activities.

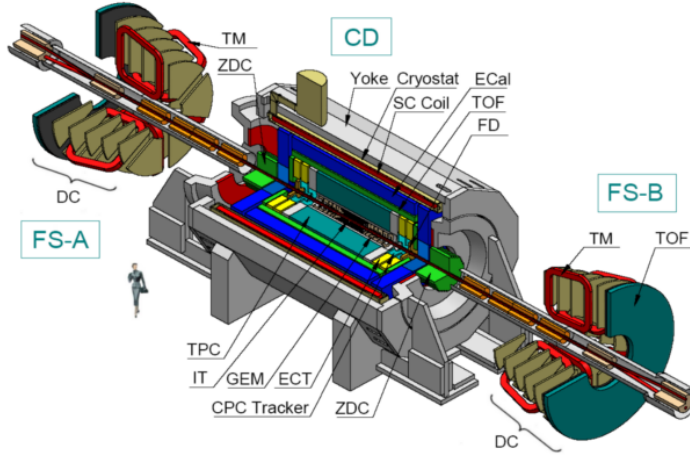The MPD is designed as a $4\pi$ spectrometer capable of detecting charged

Figure 1.2.1 – The overall schematic of the MPD subsystems.

hadrons, electrons and photons in heavy-ion collisions at high luminosity. The beam line is surrounded by the large gaseous Time Projection Chamber (TPC) which is enclosed by the Time-of- Flight (TOF) barrel. The TPC is the main tracker, and in conjunction with the TOF they will provide precise momentum measurements and particle identification.

The Electromagnetic Calorimeter (ECal) is placed in between the TOF and the MPD Magnet. It will be used for detection of electromagnetic showers, and will play the central role in photon and electron measurements. In the forward direction, the Fast Forward Detector (FFD) is located still within the TPC barrel. It will play the role of a wake-up trigger. The Forward Hadronic Calorimeter (FHCal) is located near the Magnet endcaps. It will serve for determination of the collision centrality and the orientation of the reaction plane for collective flow studies. The silicon-based Inner Tracker System (ITS) will be installed close to the interaction point in the second stage of the MPD construction. It will greatly enhance tracking and secondary vertex reconstruction capabilities. The miniBeBe detector, placed between the beam pipe and the TPC, close to the beam, is designed to aid in triggering and start time determination for the TOF. The MPD Cosmic Ray Detector (MCORD), installed on the outside of the MPD Magnet Yoke, will measure muons, also from the cosmic showers.

It is expected that the MPD will produce event-by-event information on charged particle tracks coming from the primary interaction vertices, together with identification of those particles, and information on the collision centrality.

The MPD identification power obtained for charged hadrons with combined mass-squared ($m^2$) from TOF and energy loss per distance ($dE/dx$) from TPC. In this paper we show results which was received using simulation from the TPC and TOF sub-detectors.

## 1.2.1   Time Projection Chamber (TPC)

The TPC is the main tracking detector of the MPD central barrel. It is designed to perform three-dimensional precise tracking of charged particles and momentum measurements for transverse momentum $p_T > 50$ MeV/c. The track reconstruction is based on the drift time and R-$\phi$ cylindrical coordinate measurement of the primary ionisation clusters created by a charged particle crossing the TPC. Charged particles traversing this volume ionise the gas mixture of $90\%Ar+10\%CH_4$ along helix-shaped trajectories. It covers momentum resolution for charged particles under 3% in the transverse momentum range $0.1 \leq p_T \leq 1$ GeV/c. The efficient tracking at pseudorapidities up to $|\eta| \leq 1.2$. Two-track resolution of about 1 cm. Hadron and lepton identification by $dE/dx$ measurements made with a resolution better than 8%.

Momentum measurements in combination with $dE/dx$ travelled by a particle measurements, $dE/dx$ are used for flow analysis of identified particles. $dE/dx$ traveled by a particle through a specific material is described by the Bethe-Bloch formula:

$$\frac{dE}{dx} = \frac{4\pi}{m_e c^2} \frac{n z^2}{\beta^2} \left( \frac{e^2}{4\pi\varepsilon_0} \right)^2 \left[ ln \left( \frac{2m_e c^2 \beta^2}{I(1-\beta^2)} \right) - \beta^2 \right], \qquad (1.1)$$

where $\beta$ equals $v/c$, $v$ is velocity of the particle, $c$ is the speed of light, $E$ is the energy of the particle, $x$ is the distance travelled by the particle, $m_e$ is the rest mass of the electron, $n$ is the electron density of the target, $z$ is the particle charge, $e$ is the charge of the electron, $\varepsilon_0$ is the vacuum permittivity an $I$ the mean excitation potential of the target. The TPC is calibrated in order to describe $dE/dx$ travelled with a parameterisation of the Bethe-Bloch formula. By combining total momentum ($p_{tot}$) and $dE/dx$ measurements the particles can be identified (see Fig. 1.2.2).
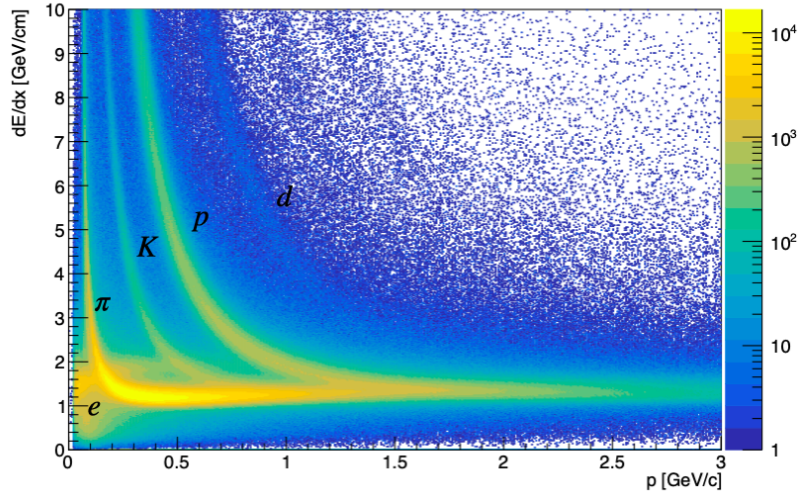
Figure 1.2.2 – Distribution of $dE/dx$ as a function of $p_{tot}$ for $e$ $\pi$, $K$, protons and deitrons. The units of dE/dx and $p_{tot}$ are GeV/cm and GeV/c, respectively.

## 1.2.2 Time-Of-Flight (TOF)

The Time-Of-Flight (TOF) system is intended to perform particle identification for momenta up to 2 GeV/c. The gas mixture is composed of 90% of $C_2H_2F_4$, 5% of $SF_6$ and 5% of $i - C_4H_{10}$. The system includes the barrel part and two endcaps and covers the pseudorapidity interval $|\eta| < 2$. The TOF is based on Multigap Resistive Plate Counters (mRPC) with satisfactory timing properties and efficiency in particle fluxes up to $103 cm^{-2}s^{-1}$. The 2.5 $m$ diameter barrel of TOF has a length of 500 $cm$ to cover the pseudorapidity region $|\eta| < 1.4$. The basic element of TOF is a 7 $cm \times 62$ $cm$ mRPC built of 12 glass plates separated by 220 $\mu m$ thick spacers forming 10 equal gas gaps. All the counters are assembled in 12 azimuthal modules providing an overall geometric efficiency of about 95%. Two options of the signal readout geometry are still being considered. One is a pad structure with lateral dimensions of 3 $\times$ 3.5 $cm$, the other makes use of 3 $cm$ wide strips with readout from both sides of the strips. The endcap TOF system consists of two discs situated at both sides of the TOF barrel. The inner diameter of the discs is 40 $cm$, the outer one is 250 $cm$ resulting in a pseudorapidity coverage of $1.5 < |\eta| < 2$. TOF measure the particle velocity relative to the speed of light in vacuum, $\beta = v/c$ which allowing the mass square $m^2$ of the particle to be determined:

$$m^2 = p^2(1/\beta^2 - 1), \tag{1.2}$$

By combining $p_{tot}$ and $m^2$ the particles can be identified (see Fig. 1.2.3).
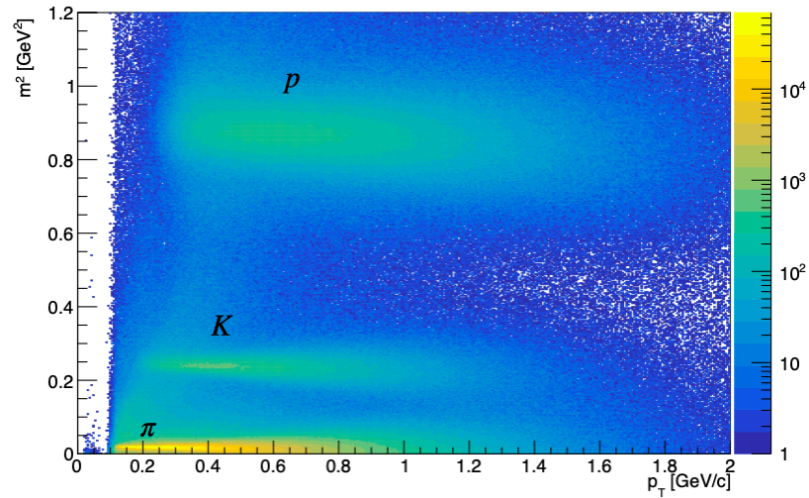


Figure 1.2.3 – Distribution of $m^2$ a function of $p_{tot}$ for $\pi^\pm$, $K^\pm$ and protons. The units of $m^2$ and $p_{tot}$ are GeV$^2$ and GeV/c, respectively.

In this analysis particle tracks are reconstructed using global tracks, which are tracks reconstructed using both TPC and TOF signals. Particles are also identified by using only the TPCs signal.

# Chapter 2

# Data

The simulated data used in the work were obtained by the Monte Carlo method using the generators UrQMDv3.4[11] and underwent the entire chain of reconstructions, on the condition of real Bi-Bi collisions of the MPD experiment with $\sqrt{s_{NN}}$=9.2 GeV.

## 2.1 Selection criteria

For a neural network research, as well as to compare particle identification results, it is necessary to conduct a preliminary selection. The selection criterion includes a restriction on the events and tracks that present in Tab. 2.1.1.

| $p_{tot}$ | $|\eta|$ | r | nHints | dca | Vz |
|---|---|---|---|---|---|
| >0.1 GeV | 1.5 | 1.25 | > 15 | < 5 | > 10 |

Table 2.1.1 – selection criteria

Where $p_{tot}$ is total momentum of particle, $\eta$ is pseudorapidity, r is sum of squares of primary x and y vertices, nHints is the number of points at which the tracks were restored, dca is distance of the closest approach. Selection on momentum and pseudorapidity are associated with the acceptance and incomplete coverage the pseudorapidity of range in TPC detector. Other criteria were used for quality improvements to the tracks.

# Chapter 3

# Neural Network research

Nowadays, various tasks in particle physics can be solved using Machine Learning (ML) techniques. For instance, particle reconstruction (clustering), identification (classification), and energy or direction measurement (regression) in calorimeters and tracking devices. In large high-energy physics experiments, particles leave the traces in several detectors, some of which are specialised in charged or neutral particle identification (PID). Usually a combination of techniques among Cherenkov and transition radiations, ionisation loss, time-of-flight measurements and calorimetry are simultaneously employed to guarantee redundancy and a wide kinematic coverage. Multivariate classification techniques are typically used to combine the all available response off all involved detectors into variables, called PID classifiers.

## 3.1   Neural Network setup

There are many Neural Network including Deep Neural Network. Each of them can be used for particle identification[12–14]. For first step of work the simplest Neural Network was chosen Multilayer perceptron (MLP). MLP is one of the standard method for multi-class and binary classification the evaluation of which for PID is shown in this paper. Neural Network research was carried out using the `scikit-learn` library [15]. During MLP research and efficiency comparison were used data of 6 particles species such as: charged pion $\pi^{\pm}$, charged kaon $K^{\pm}$, proton $p$ and anti-proton $\bar{p}$. Each of particle species dataset contain 200000 lines.

### 3.1.1 Feature selection

One of the initial part of the model preparation is feature selection. Feature selection is the process of reducing the number of input variables when developing a predictive model. It is desirable to reduce the number of input variables to both reduce the computational cost and, in some cases, to improve the performance of the model. In the beginning we have the assemble of variable which include: $p_{tot}$(momentum), $m^2$(m2), $dE/dx$(dedx), charge, nHints, dca, eta, Vz, Vy, Vz. The fig. 3.1.1 shows correlation matrix for all input features. It is not hard to notice that most features are not correlated or have very small
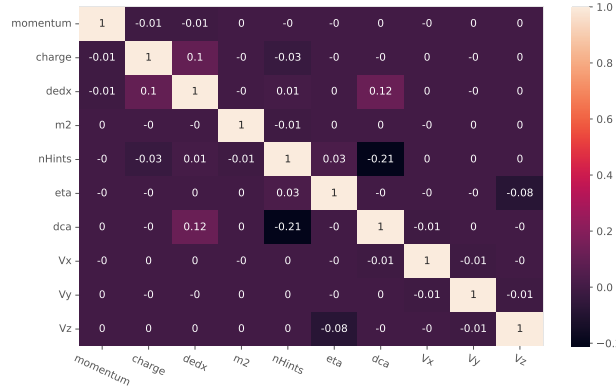


Figure 3.1.1 – Correlation matrix for input features: $p_{tot}$, $m^2$, $dE/dx$, charge, nHints, dca, eta, Vz, Vy.

correlation ($< 20\%$) with each other.

During feature selection was used MLP model with basic hyper-parameters: activation function - logistic, layer sizes - 50, number of iteration - 50, learning rate - 0.01.

On the first step in feature selection was evaluated weight matrix of MLP (see Fig. 3.1.2) As can be seen, not zero value of weight almost for each element
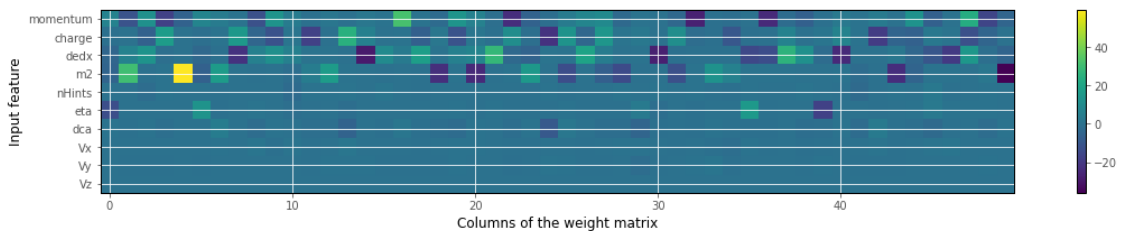


Figure 3.1.2 – Weight matrix

of MLP have four features: $p_{tot}$, charge, $dEdx$ and $m^2$. Nevertheless, such

11

features as nHints, eta and dca in some of column have a non-zero weight. That is why it is necessary to evaluate the quality of the classification with different combination of the features. For model evaluation in this section and following is used f1 score which is defined as:

$$f_1 = 2 * \frac{recall * precision}{recall + precision}, \tag{3.1}$$

F1 Score is the Harmonic Mean between precision and recall. The range for F1 Score is $[0, 1]$ and a perfect model has an F-score of 1.

$$precision = \frac{TP}{TP + FP}, \qquad recall = \frac{TP}{TP + FN} \tag{3.2}$$

Precision is the fraction of true positive examples among the examples that the model classified as positive (see Eq. 3.2 left). Recall, also known as sensitivity, is the fraction of examples classified as positive, among the total number of positive examples (see Eq. 3.2 right).

The fig. 3.1.3 demonstrates dependence of the f1-score for each class on the combination of the features. The basis of features is: $p_{tot}$, charge, $dEdx$ and $m^2$. As can be noted, the larges f1-score value have $\pi^{\pm}$ and $p$ species.
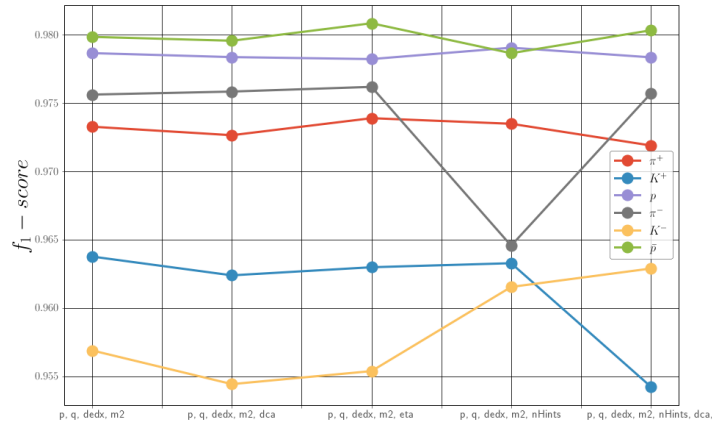


Figure 3.1.3 – Dependence of the f1-score on set of features

The reason $K^{\pm}$ species have the lowest f1-score is that, for instance, their $m^2$ distribution are between $p$ and $\pi$ and mixes with all of them (see Fig.1.2.3). Evaluation of the impact of the additional features show that each of features reduce or increase f1-score for different species. For example, usage the nHints feature increase f-score for $p, K^-$ species and decrease f-score for $\pi^-, \bar{p}$ species. Hence, in the following work were used features: $p_{tot}$, charge, $dE/dx$ and $m^2$.

## 3.1.2 Hyperparameters selection

The performance of a model can drastically depend on the choice of its hyperparameters. The choice of the optimal hyperparameters is more art than science, if we want to run it manually. Since the algorithms, the goals, the data types, and the data volumes change considerably from one project to another, there is no single best choice for hyperparameter values that fits all models and all problems. Instead, hyperparameters must be optimized within the context of each machine learning project. Even with in-depth domain knowledge by an expert, the task of manual optimisation of the model hyperparameters can be very time-consuming. An alternative approach is to set aside the expert and adopt an automatic approach. An automatic procedure to detect the optimal set of hyperparameters for a given model in a given project in terms of some performance metric is called an optimization strategy. Four commonly used optimization strategies: Grid search, Random search, Hill climbing and Bayesian optimization.

In this work was used Bayesian optimisation based on the `optuna`[16] package. The Bayesian optimisation strategy selects the next hyperparameter value based on the function outputs in the previous iterations. Unlike hill climbing, Bayesian optimisation looks at past iterations globally and not only at the last one. In Bayesian optimisation for MLP were used set of hyperparameters that are presented in Table 3.1.1.

| | |
|---|---|
| hidden_layer_sizes | 10 - 70 |
| max_iter | 10 - 100 |
| learning_rate_init | 0.0001 - 0.01 |
| activation | logistic, tanh, relu |
| learning_rate | constant, invscaling, adaptive |

Table 3.1.1 – Set of hyperparameters that were used in Bayesian optimisation

Where hidden_layer_sizes is the $i$th element represents the number of neurons in the $i$th hidden layer, max_iter is number of epochs, learning_rate_init is the initial learning rate used which controls the step-size in updating the weights, activation is the name of the output activation function and learning_rate is the learning rate schedule for weight updates.

For evaluation the quality of the model in Bayesian optimisation was used weighted f1-score which is defined as the weighted F1-score of each class by the number of samples from that class and divide by the total number of samples. For Bayesian optimisation was performed 50 trials. Almost all models with different parameters found during optimisation have f1-score > 0.97.

In Fig.3.1.4 the dependence of the weighted f1-score on the value of hyperparameters is presented. As can be noted, logistic activation function has
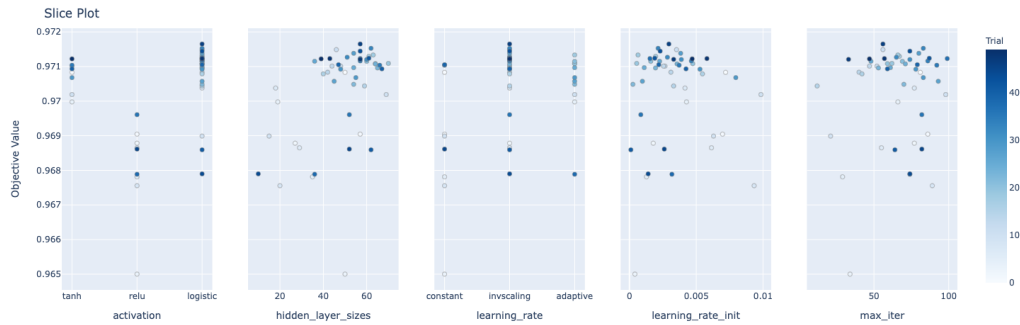


Figure 3.1.4 – Dependence of the weighted f1-score on the value of each of the hyperparameters

been selected for a larger number of models, models with relu and tanh activation function has been selected for the fewer models. All models with differnet learning_rate value type have good f1-score. Invscaling and adaptive learning rate has been selected for a larger number of models, constant learning rate has been selected for fewer number of models. The wast majority of models have learning_rate_init 0.0001-0.0050.In Fig.3.1.5 the map of hyperpartmetrs
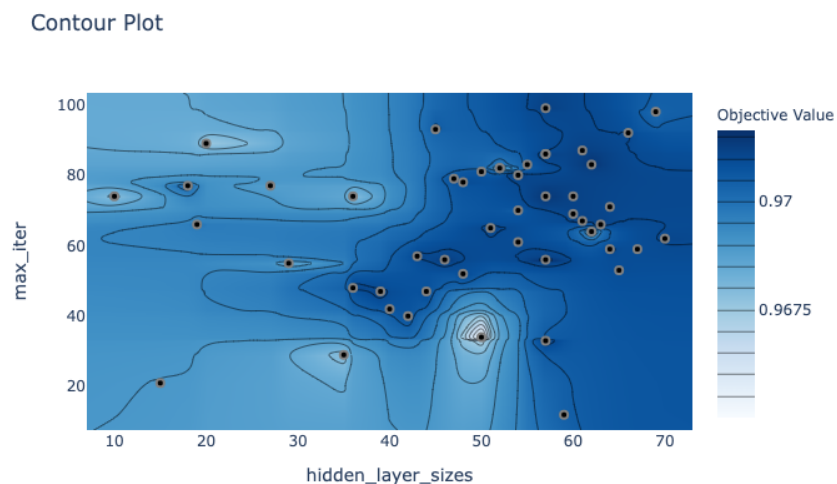


Figure 3.1.5 – Map of hyperparameters hidden_layer_sizes and max_iter

hidden_layer_sizes and max_iter. As mentioned above almost all models have f-score > 0.97, therefore, it is necessary to choose a model which have small number of the hidden_layer_sizes and max_iter. In this case, simplifying the model reduces the computational coast. The model with hidden_layer_sizes = 33 and max_iter = 64 was chosen. The full set of parameters that were chosen presented in Table 3.1.2.

| | |
|---|---|
| hidden_layer_sizes | 36 |
| max_iter | 48 |
| learning_rate_init | 0.006 |
| activation | logistic |
| learning_rate | constant |

Table 3.1.2 – Optimal set of hyperparemetrs

It is important to clarify that we can choose max_iter less than 48 when the loss or score is not improving by at least 0.0001 for 10 consecutive iterations. Different models combination with the bigger number of hidden layers were also considered. Nevertheless, adding additional hidden layers does not make significant contribution in the f1-score.

## 3.2 Training in different range of momentum

In this section are presented one of two addition approach that was researched for improve the correctly classification quality of the species. The main idea of this approach is to train MLP model using data from different range of $p_{tot}$. For this research the $p_{tot}$ ranges were selected:

$$[0.0, 0.5, 1.0, 1.5, 2.5]$$

The number of particles of a given species $i$ that are identified correctly $dN_{true}$ (aka TP) by model which was trained in full range of $p_{tot}$ and models that were trained in different range of $p_{tot}$ were compered.

$$Ratio = \frac{dN_{true}^{\text{Partial range model}}/dp}{dN_{true}^{\text{Full range model}}/dp}, \tag{3.3}$$

The models that was trained in different range of the $p_{tot}$ had the same set of hyperparameters as full momentum model. In Fig.3.2.1 and 3.2.2 the rations (see Eq. 3.3) for each species are demostraited. As visible, ratio has upper
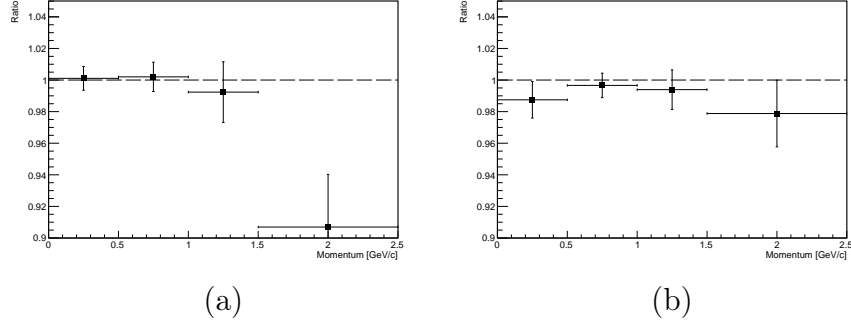


(a)

(b)

Figure 3.2.1 – Ratio for $\pi^-$, $K^-$

and lower deviation from 1.0. Thus, it can be concluded that train in different range of momentum does not make significant contribution in the particle identification. Nevertheless, the chosen binning is not optimal. That is why
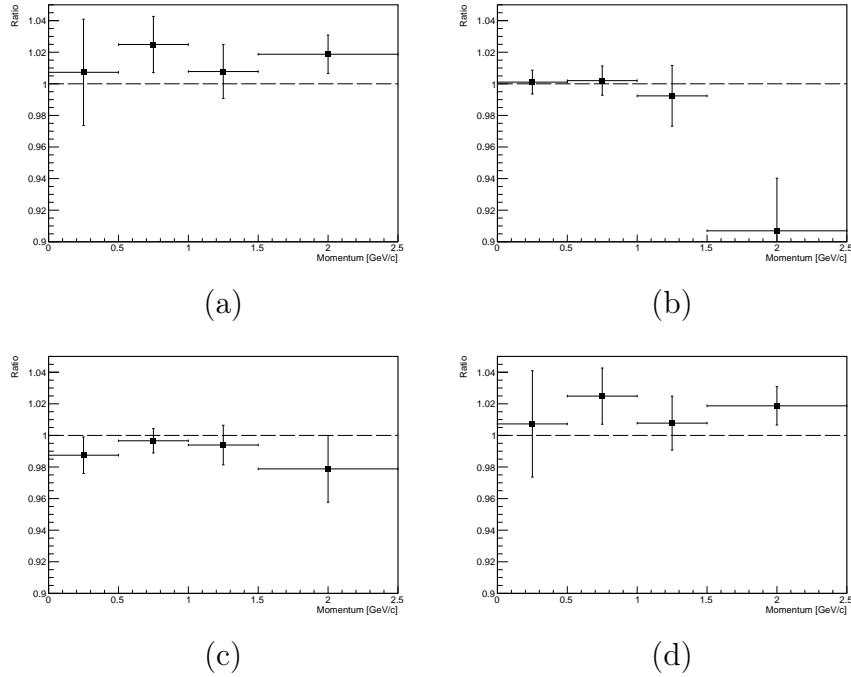


(a)

(b)

(c)

(d)

Figure 3.2.2 – Ratio for $p$, $\pi^+$, $K^+$, $\bar{p}$

studies with different range of momentum can be preform in future.

## 3.3 Binary classification for each particle

In this section are presented the second addition approach that was researched for improve the correctly classification quality of the species. The main idea of this approach is to do binary MLP classifier for each species. In other words, to train models on the relabled data for each species and after that combine the model output. For each classifier have been done Bayesian optimisation. By analogy, which is described in the Sec. 3.1.1 have been selected hyperparameters for each binary MLP classifier. The sets of hyperparametrs for each binary MLP classifier almost have coincided with each other. That is why for each species was used the same set of hyperparametrs. The selected hyperparameters shown in Table 3.3.1

| | |
|---|---|
| hidden_layer_sizes | 30 |
| max_iter | 40 |
| learning_rate_init | 0.006 |
| activation | logistic |
| learning_rate | constant |

Table 3.3.1 – Optimal set of hyperparemetrs for binary classifier

It is not difficult to see that hyperparametrs for binary classification almost coincides to the hyperparametrs for multi classification. Considering early stopping, which is described in the end of the section 3.1.2 the sets of hypertparemtres could be the same



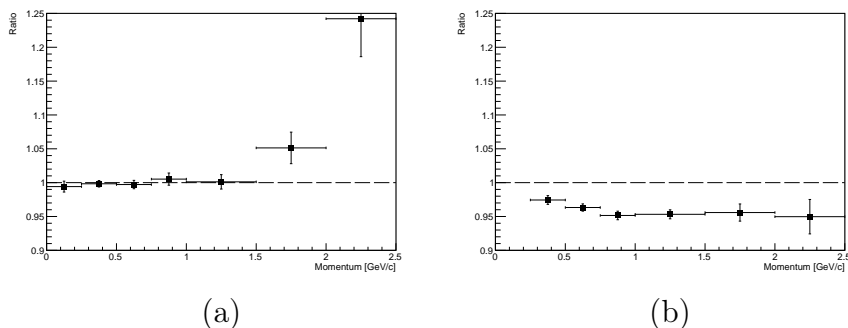(a)                                                                          (b)

Figure 3.3.1 – Ratio for $\pi^-$, $K^-$

Prediction of the binary classifiers combination have been compered with the multi-classification prediction. By analogy to the previous additional approach (see Sec. 3.2) ratio of the right answers was compered. The results of
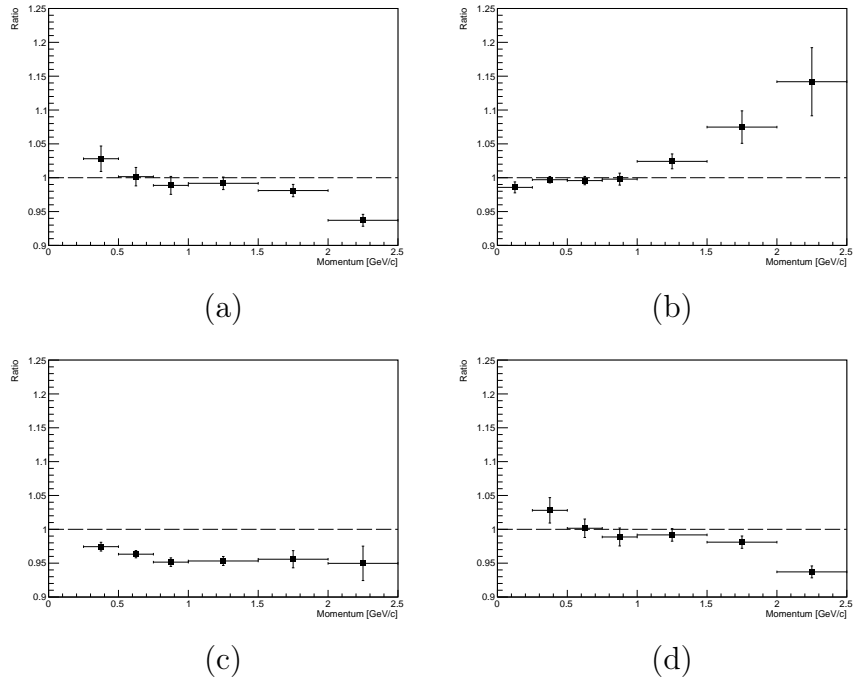
Figure 3.3.2 – Ratio for $p$, $\pi^+$, $K^+$, $\bar{p}$

models comparison shown in Fig 3.3.1 and 3.3.2. This plots show us that for wast majority of the particle species this approach decrease or does not make significant contribution. But for $\pi^+$ and $\pi^-$ species (see Fig. 3.3.1a, 3.3.2b) approach make significant contribution about $20 - 30\%$ in range of $p_{tot} > 1.5$ [GeV/c]. Despite the large contribution in the high $p_{tot}$ range for $\pi^{\pm}$ the approach was not used in the final research of efficiency (see Sec. 4.2).

# Chapter 4

# Particle Identification

Simple selections based on the individual Particle Identification (PID) signals of each sub-detector do not take full advantage of the PID capabilities of MPD. An example of this is illustrated in Fig. 4.0.1, which shows the separation of the expected TPC and TOF signals for $\pi^+$, $K^+$ and $p$ with $p_{tot}$ in the range $0.7 < p_{tot} < 1.0$ GeV/c. Clearly, the separation in the two-dimensional plane
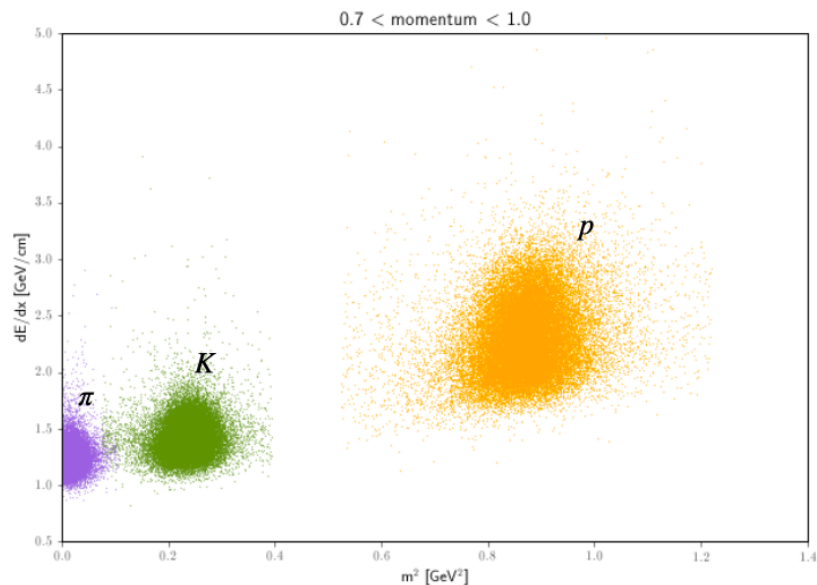


Figure 4.0.1 – dE/dx (from TPC) versus mass-squared (from TOF) distribution for $\pi^+$, $K^+$ and $p$ species.

(the peak-to-peak distance) of e.g. pions and kaons is larger than the separation of each individual one-dimensional projection. That is why, it is necessary to use different approaches for particle identification.

## 4.1   n-Sigma approach

One of the standard approach for combining the information from several sub-detectors is n-Sigma approach. Basic idea of the n-Sigma approach is to perform parameterisation for Bethe-Bloch energy loss formula (see Eq.1.1) using asymmetric Gaussian function. The most probability value $\langle dE/dx \rangle$ and standard deviation $\sigma_{dE/dx}$ are obtained from projection energy loss in TPC for each given momentum. The most commonly used discriminating variable for PID is the $N_\sigma$ variable, defined as the deviation of the measured signal from the most probability value for every species of particle $i$. For TPC $N_{\sigma_{TPC}}$ is defined as:

$$N^i_{\sigma_{TPC}} = \frac{dE/dx - \langle dE/dx \rangle^i}{\sigma^i_{TPC}}, \tag{4.1}$$

The same procedure is performed for $m^2$ from TOF:

$$N_{\sigma^i_{TOF}} = \frac{m^2 - \langle m^2 \rangle^i}{\sigma^i_{m^2}}, \tag{4.2}$$

where $\langle m^2 \rangle$ is most probability value and $\sigma_{m^2}$ standard deviation for $m^2$.

A certain species is assigned to a particle if this value lies within a certain range around the expectation $N_{\sigma_{TPC}} = 2$ and $N_{\sigma_{TOF}} = 2$.

$$N_\sigma \leq \sqrt{N^2_{\sigma^i_{TOF}} + N^2_{\sigma^i_{TPC}}}, \tag{4.3}$$

If the condition 4.3 is met for $N^i_{TPC}$ and $N^i_{TOF}$ the particle compatible to $i$-species. In case a particle can be compatible with more than one species n-Sigma approach corresponds to false decision. The application of the PID n-Sigma approach was based on `MpdRoot`[17] framework.

## 4.2 Comparison of MLP and n-Sigma approach

In this section show results of comparison of two approach for particle identification. In order to evaluate the quality of PID approach it is necessary to compute efficiency. The PID Efficiency of the species $i$ is defined as the proportion of particles of a given species $i$ that are identified correctly $dN^i_{\text{true}}$ in certain momentum range $dp$ divided by all generated $i$ particle $dN^i_{textallgen.}$ in the same momentum range $dp$.

$$Efficiency = \frac{dN^i_{\text{true}}/dp}{dN^i_{\text{all gen.}}/dp},\tag{4.4}$$

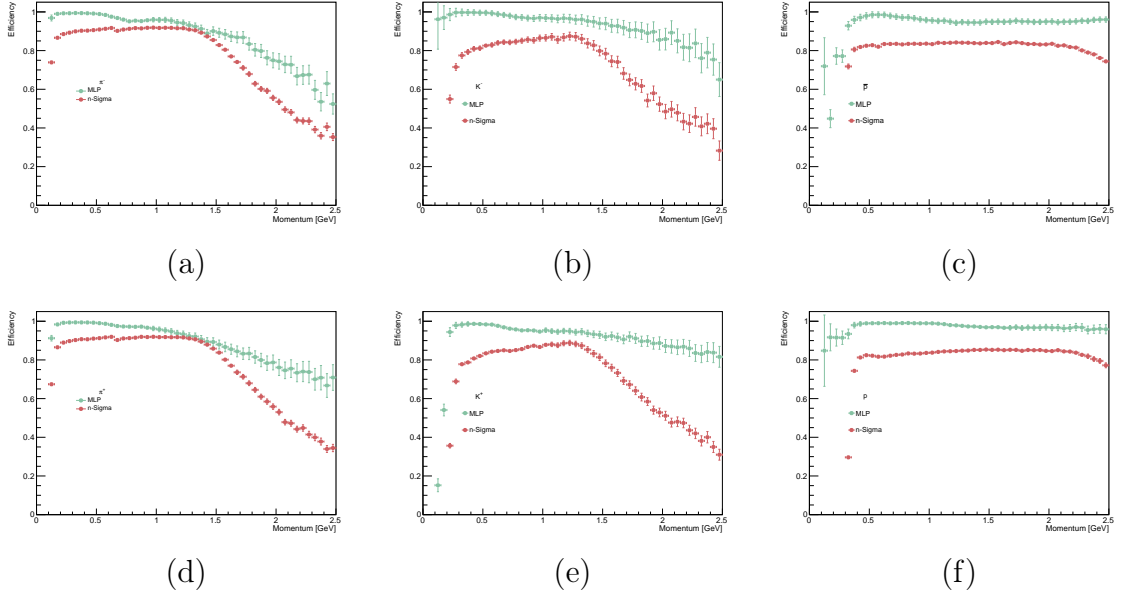Results of efficiency comparison is demonstrated in Fig. 4.2.1. As can be seen,



Figure 4.2.1 – Comparison of the MLP(green) and nSigma(red) approach efficiency for $\pi^-$, $K^-$, $p$, $\pi^+$, $K^+$, $\bar{p}$

for each particle species MLP approach have higher efficiency then n-Sigma approach for full range of momentum.

# Conclusion

In this paper, several studies have been performed to improve the quality of PID using MLP classifier. For MLP multi-classifier have been chosen the number of features that contribute the biggest contribution in identification. Using Bayesian optimisation have chosen the hyperparametrs that do not complicate MLP model and allow to get high $f_1$-score.

Additional approaches for improve the correctly classification quality of the species were researched. Each of them have not shown significant results for PID, however these approaches can be researched in future with another setting.

The n-Sigma approach was studied and compared with MLP approach for particle identification. It was shown usage MLP classifier for particle identification considerably improves efficiency for each particle species.

The improvement is shown only for the certain version of MC simulation data. In the future, it is planned to conduct research for a wide set of MC data.

# References

1. The SIS Heavy Ion Synchrotron Project / K. Blasche [et al.] // IEEE Transactions on Nuclear Science. — 1985. — Vol. 32. — P. 2657–2661.

2. *Beth R. A.*, *Lasky C.* The Brookhaven Alternating Gradient Synchrotron // Science. — 1958. — Vol. 128, no. 3336. — P. 1393–1401. — ISSN 00368075, 10959203.

3. *Evans L. R.* The SPS collider: status and outlook. — 1987.

4. *Ludlam T. W.*, *Samios N. P.* The relativistic heavy ion collider project: a status report // Quark Matter / ed. by H. Satz, H. J. Specht, R. Stock. — Berlin, Heidelberg : Springer Berlin Heidelberg, 1988. — P. 353–359.

5. *Green D.* The Status of the CERN Large Hadron Collider (LHC) // Frontiers in Optics 2010/Laser Science XXVI. — Optica Publishing Group, 2010. — STuB4.

6. *Marx J. N.* The STAR Experiment at RHIC // Advances in Nuclear Dynamics 2 / ed. by W. Bauer, G. D. Westfall. — Boston, MA : Springer US, 1996. — P. 233–237. — ISBN 978-1-4757-9086-3.

7. Bulk Properties of the Medium Produced in Relativistic Heavy-Ion Collisions from the Beam Energy Scan Program / L. Adamczyk [et al.] // Phys. Rev. C. — 2017. — Vol. 96, no. 4. — P. 044904. — arXiv: 1701.07065 [nucl-ex].

8. Status of NICA / V. Kekelidze [et al.] // EPJ Web of Conferences. — 2018. — Vol. 182. — P. 02063.

9. Status and initial physics performance studies of the MPD experiment at NICA / V. Abgaryan [et al.] // Eur. Phys. J. A. — 2022. — Vol. 58, no. 7. — P. 140. — arXiv: 2202.08970 [physics.ins-det].

10. Three Stages of The NICA Accelerator Complex Nuclotron-based Ion Collider fAcility / V. D. Kekelidze [et al.] // 3rd Large Hadron Collider Physics Conference. — Gatchina : Kurchatov Institute, 2016. — P. 565–569.

11. UrQMD (Ultra relativistic Quantum Molecular Dynamics) hadron-String Transport Model. — URL: http://urqmd.org/..

12. Using machine learning for particle identification in ALICE / Ł. K. Graczykowski [et al.] // JINST. — 2022. — Vol. 17, no. 07. — P. C07016. — arXiv: 2204.06900 [nucl-ex].

13. A Neural-Network-defined Gaussian Mixture Model for particle identification applied to the LHCb fixed-target programme / G. Graziani [et al.] // JINST. — 2022. — Vol. 17, no. 02. — P02018. — arXiv: 2110.10259 [hep-ex].

14. *Fanelli C.*, *Mahmood A.* Artificial Intelligence for imaging Cherenkov detectors at the EIC // JINST. — 2022. — Vol. 17, no. 07. — P. C07011. — arXiv: 2204.08645 [physics.ins-det].

15. scikit-learn: machine learning in Python. — URL: https://scikit-learn.org/.

16. A hyperparameter optimization framework. — URL: https://optuna.org/.

17. The MpdRoot Framework. — URL: https://git.jinr.ru/nica/mpdroot/-/blob/dev/README.md.