

Machine learning in ALICE: status and perspectives



Francesco Mazzaschi ^{1, 2} Workshop di CCR, 24/05/21-28/05/21 Università degli Studi di Torino ¹, INFN Torino ²



1

Machine learning in ALICE





• INFN group of ALICE strongly involved in many of these activities, focusing mostly on light and heavy flavour analyses



Physics analyses in Run 2

Analysis workflow





- Training and testing on ROOT TTrees of particle candidates satisfying loose selections
 - signal candidates from Monte Carlo simulations
 - background candidates both from Monte Carlo and data (like-sign, sidebands)
- Application on local machines:
 - feasible if the TTree size is lower than few hundreds of GBs (pp and pPb collisions)
 - personal laptop for TTree size up to ~3-4 GBs, otherwise dedicated machine (e.g. Turin ML workstation: 256 GBs RAM, 128 threads)

Analysis workflow





- Training and testing on ROOT TTrees of particle candidates satisfying loose selections
 - signal candidates from Monte Carlo simulations
 - background candidates both from Monte Carlo and data (like-sign, sidebands)
- Application on grid:
 - used when the size of the TTree exceeds the TB
 - if training is done with a python package, require an interface with the ALICE environment (treelite, onnx-runtime)

Ο

two different approaches used: TMVA (ROOT) and python Ο



- Limited selection of models
- Usually BDT based on AdaBoost algorithm

Boosted Decision Trees

Widely used algorithms in HEP analyses



- few comparisons done between BDTs an DNNs, no significant differences for the physics analyses



- Widely used outside HEP
 - more models and analysis tools Ο available
- Need interface with ALICE ecosystem
 - treelite, onnx-runtime Ο



Light and heavy flavour analyses with ML

- From the ML point, both heavy and light flavour analyses are very similar
 - BDTs trained employing high-level physical variables (decay length, PID, ...)
- ML employed for reconstructing rare signals
 - low S/B ratio, especially in the very high multiplicity environment of Pb-Pb collisions
 - small systems (p-p, p-Pb): no specific triggers
 -> not negligible statistical uncertainties
- Next slides: examples of physics analyses in LF and HF sectors
 - not-comprehensive overview (e.g. applications to jet analyses not covered), focused on activities with large participation of INFN members



ML Analyses - Hypertriton

- Lightest hypernucleus: p, n, Λ
- Hypertriton lifetime and binding energy are powerful tools for investigating Λ-Nucleon interactions
- Reconstructed via its two-body decay channel

$$\circ \quad {}^3_{\Lambda}{
m H}
ightarrow {
m He} + \pi$$

- Pb-Pb @ √s = 5.02 TeV
- Selection using XGBoost BDT
 - signal extraction with high significance in 9 ct bins

Counts / (2.25 MeV/*c*²)

120

100

80

60

4(

20

2.96

ALT-PERF-335127

2 97

.....





ML Analyses - Hypertriton

ALICE

- Selection using XGBoost BDT
 - Most precise measurements of the lifetime and binding energy



ALI-DER-358857

ML Analyses - D mesons

- Measurements of prompt and non-prompt D^0 , D^+ , D^+_{s} mesons
 - prompt: from charm-quark hadronisation or excited charm hadron decay
 - non-prompt: from beauty hadron decays
- p-p @ √s = 5.02 TeV <u>arXiv:2102.13601</u>
- Multi-classification BDT (XGBoost) to separate prompt, non-prompt and bkg (D⁰ with 2-step BDTs, details in backup)





Francesco Mazzaschi

0

Non-prompt fraction determined by sampling the raw yields for different BDT outputs Most precise charm and beauty production measurements at midrapidity at the LHC

ML Analyses - D mesons





ML Analyses - Λ_{c}





Towards LHC Run 3

Francesco Mazzaschi

TPC space-charge distortion corrections with CNN

ALICE

- ALICE Time Projection Chamber (TPC):
 - identification of charged particles via reconstruction of the trajectory in the magnetic field and the specific energy loss
- Upgrade to allow for continuous readout up to 50kHz
 - significant amount of positive ions (space charge) from the readout chamber cause large space-charge distortions
 - average corrections used, but fluctuations of several mm are expected (TPC resolution: ~200 μm)

Space-charge density fluctuations

• TPC pad signals integrated over 1 ms intervals (IDCs) as an estimator





Train a ML algorithm for fast data-driven prediction

Space-charge distortion fluctuations

- Solve Poisson equation of the space-charge density to obtain the electrical field deviation
- Solve Langevin equations to obtain the deflection of the ionization electrons along their drift path



See vCHEP2021 talk of Ernst Hellbär

TPC space-charge distortion corrections with CNN



- Convolutional Neural Network (CNN) originally developed for biomedical image used
 - o <u>Ú-Net</u>
- Implementation in python
 - Keras, Tensorflow
 - ROOT and RootInteractive for evaluation of trained models in multi-dimensional TPC phase space
- Input variables
 - space-charge density fluctuations
 - mean space-charge density
- Output
 - \circ distorsion fluctuations in radial direction



See vCHEP2021 talk of Ernst Hellbär

Francesco Mazzaschi

TPC space-charge distortion corrections with CNN



- mean (μ) and RMSE of the difference between predicted and true distortion fluctuations
- Variation of
 - \circ grid size: 90 x 17 x 17, 180 x 33 , x 33
 - number of training samples: 5k, 10k, 18k
- Promising preliminary results:

See vCHEP2021 talk of Ernst Hellbär

- for 10k training events the fluctuations are always lower than the intrinsic TPC resolution
- systematic dependence of the predicted results on the TPC radius -> investigations on-going
- \circ $% \phi$ aim to extend the correction also to the ϕ and z coordinates





What's next



- In last years, the use of Machine Learning has spread to many areas of the ALICE experiment: this is a non-reversible process
- Run 3 is coming, we have many projects based on ML ongoing
 - online data quality monitoring/data anomaly detection
 - ➤ GANs, BDTs
 - \circ centrality determination in proton-nucleus and nucleus-nucleus collisions

➤ BDTs

- n dimensional reweighting of the Monte Carlo samples
 - > DNNs
- \circ centralized framework in the ALICE O² software for ML application on Grid
 - > onnx-runtime, treelite
- \circ training of ML models on large datasets (\approx 1 TB)
 - > new ALICE Analysis Facility (see <u>talk</u> of M.Concas on Friday)
 - HPC resources available in our local sites
 - > onnx-runtime

Thanks for your attention!



Backup



- ³ A production extremely sensitive to the nucleosynthesis model at low charged-particle multiplicity density
- Reconstructed via its two-body decay channel: $\circ {}^3_{\Lambda}{
 m H}
 ightarrow {}^3{
 m He} + \pi$
- p-Pb @ √s = 5.02 TeV
- Selection using XGBoost BDT



First measurement of hypertriton production in p-A collisions: significance of 4.6σ



- ³ A production extremely sensitive to the nucleosynthesis model at low charged-particle multiplicity density
- Reconstructed via its two-body decay channel: $\circ {}^{3}_{\Lambda} H \rightarrow {}^{3} He + \pi$
- p-Pb @ √s = 5.02 TeV
- Selection using XGBoost BDT



Comparison with available theoretical models: ${}^{3}_{\Lambda}$ H / Λ compatible within 1.2 σ with 2-body coalescence

ML Analyses - D mesons

- Measurements of non-prompt D⁰ meson
 - prompt: from charm-quark hadronisation or excited charm hadron decay
 - non-prompt: from beauty hadron decays
- Pb-Pb @ √s = 5.02 TeV
- 2-step binary classification:
 - 2 BDTs (XGBoost) to separate prompt, non-prompt and bkg







ML Analyses - D mesons



- Measurements of non-prompt D⁰ meson
 - prompt: from charm-quark hadronisation or excited charm hadron decay
 - non-prompt: from beauty hadron decays
- Pb-Pb @ √s = 5.02 TeV
- Probe beauty R_{AA} in central Pb-Pb collisions down to very low p_T
- Ratio of R_{AA} of non-prompt and prompt D^0
 - close agreement with most of the available theoretical models



ML Analyses - Λ_c





• Precise production measurement of Λ_c vs average charged-particle multiplicity