Energy reconstruction with machine learning techniques in JUNO: aggregated features approach

Arsenii Gavrikov^{1,2}, Yury Malyshkin², Fedor Ratnikov¹ on behalf of the JUNO collaboration

¹HSE University, Moscow, Russia

²Joint Institute for Nuclear Research, Dubna, Russia

The 6th International Workshop on Deep Learning in Computational Physics



NATIONAL RESEARCH UNIVERSITY A. Gavrikov (HSE+JINR)



Introduction to the JUNO experiment

Jiangmen Underground Neutrino Observatory:

- multipurpose experiment
- 53 km away from 8 reactor cores in China
- data taking expected in ${\sim}2023$
- JUNO Collaboration:
 - 77 institutions
 - 697 collaborators
- 2 The main goals of JUNO:
 - neutrino mass ordering (3σ in 6 years)
 - precise measure of oscillation parameters $\sin^2 \theta_{12}, \Delta m^2_{21}, \Delta m^2_{31}$

The Central Detector:

- detection channel: $\overline{\nu}_e + p \rightarrow e^+ + n$;
- deposited energy converts to optical light
- the largest liquid scintillator detector: 20 kt
- 77.9% photo-coverage: 18k 20", 26k 3" photo-multiplier tubes (PMTs)



A. Gavrikov (HSE+JINR)

Machine Learning (ML) in HEP

- ML methods are used at all levels of data processing in many HEP experiments:
 - signal/background discrimination
 - event selection in the trigger
 - event simulation
 - anomaly detection
 - identification, etc.
- Why is ML useful for HEP?
 - Faster. More precisely, with proper training
 - **Adequate** for many purposes simultaneously: event simulation, analysis, reconstruction, identification, etc.
 - GPU friendly by construction, which is important for big data processing
- Machine-learning algorithms use statistics to find patterns in massive amounts of data
- Our task is a supervised learning problem (regression)

通とくほとくほと

э.

Problem statement



A. Gavrikov (HSE+JINR)

DLCP 2022

2022-07-06

Datasets

- Two datasets: for training and for testing
- generated by the Monte Carlo method

Data description:

- ositron events
- uniformly spread in the volume of the central detector
- **3** $E_{kin} \in [0, 10]$ MeV. $E_{dep} = E_{kin} + 1.022$ MeV
- Training dataset:
 - 5 million events
 - S uniformly distributed in kinetic energy E_{kin}

- full detector and electronics simulation
- using the official JUNO software

- Testing dataset:
 - subsets with discrete kinetic energies:
 - **(a)** 0, 0.1, 0.3, 0.6, 1, 2, ..., 10 [MeV]
 - $\sum = 1.4$ million events: each subset contains 100k

Aggregated features

We use aggregated information from the whole array of PMTs as features for models:

- AccumCharge the accumulated charge on fired PMTs
- IPMTs the total number of fired PMTs
- Ocordinates of the center of charge:

$$(x_{\rm cc}, y_{\rm cc}, z_{\rm cc}) = \vec{r}_{\rm cc} = \frac{\sum_{i=1}^{N_{\rm PMTs}} \vec{r}_{\rm PMT_i} \cdot n_{{\rm p.e.},i}}{\sum_{i=1}^{N_{\rm PMTs}} n_{{\rm p.e.},i}}$$

and its radial component: $R_{
m cc} = |ec{r}_{
m cc}|$

Oordinates of the center of FHT:

$$(x_{\text{cht}}, y_{\text{cht}}, z_{\text{cht}}) = \vec{r}_{\text{cht}} = \frac{1}{\sum_{i=1}^{N_{\text{PMTs}}} \frac{1}{t_{\text{ht},i}+c}} \sum_{i=1}^{N_{\text{PMTs}}} \frac{\vec{r}_{\text{PMT}_i}}{t_{\text{ht},i}+c},$$

and its radial component: $R_{\mathrm{cht}} = |ec{r}_{\mathrm{cht}}|$

$$\begin{array}{l} \bullet \quad \gamma_z^{\text{cc}} = \frac{z_{\text{cc}}}{\sqrt{x_{\text{cc}}^2 + y_{\text{cc}}^2}} \\ \bullet \quad \gamma_y^{\text{cc}} = \frac{y_{\text{cc}}}{\sqrt{x_{\text{cc}}^2 + z_{\text{cc}}^2}} \\ \bullet \quad \gamma_x^{\text{cc}} = \frac{x_{\text{cc}}}{\sqrt{z_{\text{cc}}^2 + y_{\text{cc}}^2}} \\ \bullet \quad \theta_{\text{cc}} = \arctan \frac{\sqrt{x_{\text{cc}}^2 + y_{\text{cc}}^2}}{z_{\text{cc}}} \\ \bullet \quad \phi_{\text{cc}} = \arctan \frac{y_{\text{cc}}}{x_{\text{cc}}} \\ \bullet \quad J_{\text{cc}} = R_{\text{cc}}^2 \cdot \sin \theta_{\text{cc}} \\ \bullet \quad \rho_{\text{cc}} = \sqrt{x_{\text{cc}}^2 + y_{\text{cc}}^2} \\ \bullet \quad \phi_{\text{cc}} = \sqrt{x_{\text{cc}}^2 + y_{\text{cc}}^2} \\ \bullet \quad \phi_$$

with 7 similar features for the components of the center of FHT

御 とう ヨ とう アン

6/16

3

Aggregated features

- Percentiles of FHT and charge distributions:
 - { $ht_{2\%}$, $ht_{5\%}$, $ht_{10\%}$, $ht_{15\%}$, ..., $ht_{90\%}$, $ht_{95\%}$ }
 - $\bullet \ \{pe_{2\%}, pe_{5\%}, pe_{10\%}, pe_{15\%}, ..., pe_{90\%}, pe_{95\%}\}$

- E = 1.02 MeV - E = 5.02 MeV - E = 9.02 MeV

- Differences between percentiles for FHT:
 - { $ht_{5\%-2\%}$, $ht_{10\%-5\%}$, ..., $ht_{95\%-90\%}$ }
- Moments for FHT and charge distributions:
 - $\bullet \ \{ht_{mean}, ht_{std}, ht_{skew}, ht_{kurtosis}\}$
 - $\{pe_{mean}, pe_{std}, pe_{skew}, pe_{kurtosis}\}$

0.75 0.75 1000 30 F(nPE) 800 F(t) 0.5 0.5 200 600 400 0.25 0.25 100 200 Ω ſ 0 400 600 800 1000 200 1200 0 200 600 800 1000 1200 0 400 0 nPE t. ns

CDFs and PDFs for FHT (left) and charge (right) distributions. $R \simeq 0$ m, E_{kin} varied. Dashes lines show mean values, $q \approx 0$

A. Gavrikov (HSE+JINR)

- E = 1.02 MeV - E = 5.02 MeV - E = 9.02 MeV

Models description: BDT

A Decision Tree (DT) takes a set of input features and splits input data recursively based on those features.

Boosted Decision Trees (BDT):

- Ensemble model
- DT as base algorithm
- DTs in BDT are trained sequentially
- Each subsequent DT is trained to correct errors of previous DTs in the ensemble



Figure: BDT demonstration. Source: https://arogozhnikov.github.io/

Main tunable hyperparameters:

- **Max. depth**: The maximum depth of a tree (usually <12)
- Learning rate: This determines the impact of each tree on the final outcome (usually ≈ 0.1)
- Number of trees: How many trees in ensemble

Benefits:

- Fast for training and prediction
- Easier to tune
- Minimalistic

BDT: optimized set of features

BDT from XGBoost:

• Optimized **set of features** (sorted by *importance*):







- Optimized **hyperparameters** (using Grid Search):
 - The maximum depth of the tree: 10
 - 2 Number of trees in the ensemble: $\simeq 300$

-

э

Learning rate: 0.08

```
A. Gavrikov (HSE+JINR)
```

Models description: FCDNN



Input laver Hidden lavers

- A *fully connected neural network* consists of layers with sets of units called **neurons**
- Neuron computes a **linear combination** of its inputs and passes it to a non-linear activation function *h*:
 f(**x**) = h (W**x** + **b**)
- Each neuron in a layer is connected with each neuron in the next layer
- Many layers **deep** neural network

-

э

(日)

Models description: FCDNN

Fully-connected deep neural network (FCDNN):



- The search for hyperparameters was performed using *BayesianOptimizer*
- Training with *early stopping*
- Validation dataset: 400k events
- *Selected features* provided the same ۰ performance as full set:



 $\{h_{12\%}, h_{15\%}, h_{10\%}, h_{15\%}, \dots, h_{190\%}, h_{195\%}\}$

-

э

Results

Metrics:

- Defined by a Gaussian fit of the $E_{\text{predicted}} E_{\text{dep}}$ distributions
- <u>*Resolution*</u>: σ/E_{dep} , where σ standard deviation of the fit
- <u>Bias</u> μ/E_{dep} , where μ mean of the fit

Parameterization:

$$\frac{\sigma}{E_{\rm dep}} = \sqrt{\left(\frac{a}{\sqrt{E_{\rm dep}}}\right)^2 + b^2 + \left(\frac{c}{E_{\rm dep}}\right)^2}$$

Models' pred. time and memory usage:

	BDT	FCDNN
Pred. time, sec/100k	3.5	17
Size, MB	50	12



13/16

Calibration sources

• We consider three calibration sources (i.e. sources with well-known signal):

Source	Type	Radiation
$^{241}Am - ^{13}C$	γ	neutron + 6.13 MeV
⁶⁰ Co	γ	1.173 + 1.333 MeV
⁶⁸ Ge	e^+	annihilation 0.511 + 0.511 MeV

• Can be used for **validation**: the agreement between the expected source spectra and the spectra reconstructed from the real calibration data will indicate the robustness of the algorithms' prediction



• There is an additional bias, caused by the different event topology for gamma sources as opposed to positrons in the training dataset, which was corrected using values predicted by the models on pure gamma events

A. Gavrikov (HSE+JINR)	DLCP 2022	2022-07-06

14/16

- Energy reconstruction using the information collected by PMTs
- Aggregated features approach
- The following ML models are used: BDT, FCDNN
- As a result *achieved*:
 - High **quality** 3% @ 1 MeV, requared for physics goals of JUNO
 - **②** Great **computation speed**, thanks to a small set of aggregated features (in $10^4 10^5$ times faster than traditional methods)
- Considered *three calibration sources* for the future evaluation of the models on the real data

Publications:

- A. Gavrikov, et al. arXiv: 2206.09040 (2022)
- A. Gavrikov, et al. EPJ Web Conf. 251 (2021), 03014
- Z. Qian, et al. NIM-A 1010 (2021), 165527
- Onferences:
 - Poster at the «15th Pisa Meeting on Advanced Detectors», INFN, 22.05.22 28.05.22
 - Poster at the «56th meeting of the PAC for Particle Physics», JINR, 24.01.22
 - Talk at the vCHEP2021 conference, CERN, **17.05.21 22.05.21**