





Search for ultra-high energy photons with Machine Learning in the Auger Prime era

Ezequiel E. Rodriguez Supervisors: Dr. Brian Wundheiler, Prof. Dr. Ralph Engel



About me...

- Finished my undergraduate studies in 2019 in Bahia Blanca
- Worked as a Machine Learning model validator (1.6 years) and data scientist (6 months).
 For the most part engaged in fraud detection (small signals hidden in background).
- Started my PhD and joined the DDAp and HIRSAP in Abril





Photon search with the Underground Muon Detector (UMD)

Motivations:

Why photons?

- Tracers of the highest-energy processes in the Universe
- · Point back to their sources
- One of the "messengers" in the Multi-Messenger approach

The objective:

Estimate or constraint the diffuse photon flux at low energies $\rightarrow \log_{10} \left(\frac{E_{\rm MC}}{\rm eV}\right) \le 17.5$



Multi-variate analysis in photon searches



Combined with Fisher's Linear Discriminant
 Analysis

0.008 Data Photon Sim. (Test) Proton Sim. (Test) Photon Sim. (Training) Proton Sim. (Training) Photon Sim. (Test) Proton Sim. (Test) Photon Sim. (Training) Proton Sim. (Training) 0.007 (normalized) 0.002 0.004 Entries (normalized) 0.8 0.6 Entries 0.003 0.2 0.00 500 600 700 800 900 1000 1100 $\log_{10}^{0} (S_{b}^{1} [VEM])$ X_{max} [g cm⁻²] Data Photon Sim. (Test) Proton Sim. (Test) Photon Sim. (Training) Proton Sim. (Training) 0.35 0.3 Entries (normalized) 0.25 0.2 0.1 0.05 16 2 4 6 8 10 12 14 18 20 Nstations

Auger Coll., ApJ 933 125, 2022

- Event observables sensitive to photon signatures
- Combined with Boosted Decision Trees



Current strategy - Embracing uncertainty

- Represent events in tabular form
- Employ both *event-level* observables and *station-level* observables
- LDF-like observables are defined by binning the distance to the shower axis according to the core resolution

	rec_energy	rec_zenith	n_counters	n_stations	saturation_flag	mu_45	mu_75	mu_105	mu_135	mu_165	
0	1.11e+17	27.56	3.0	13.0	1.0	NaN	NaN	NaN	NaN	0.07	
1	1.21e+17	28.08	4.0	12.0	1.0	NaN	NaN	NaN	NaN	0.28	
2	1.05e+17	26.21	4.0	13.0	1.0	NaN	NaN	0.38	NaN	NaN	
3	1.11e+17	26.67	6.0	12.0	1.0	NaN	NaN	NaN	0.21	NaN	
4	9.82e+16	27.41	5.0	12.0	1.0	NaN	NaN	NaN	0.10	NaN	
5	1.10e+17	26.98	6.0	13.0	1.0	NaN	NaN	NaN	NaN	0.28	
6	1.15e+17	27.89	4.0	13.0	1.0	NaN	NaN	NaN	NaN	NaN	
7	1.14e+17	27.56	7.0	11.0	1.0	NaN	NaN	NaN	0.07	NaN	
8	1.15e+17	27.28	3.0	12.0	1.0	NaN	NaN	NaN	NaN	NaN	
9	1.13e+17	27.80	4.0	13.0	0.0	NaN	0.8	NaN	NaN	NaN	
10	8.32e+16	30.78	3.0	11.0	2.0	NaN	0.6	NaN	NaN	NaN	
11	8.69e+16	30.20	6.0	10.0	1.0	NaN	NaN	NaN	NaN	NaN	
12	9.12e+16	29.75	6.0	10.0	1.0	NaN	NaN	NaN	0.10	NaN	
13	8.10e+16	29.98	3.0	12.0	1.0	NaN	NaN	NaN	NaN	NaN	
14	9.00e+16	30.10	3.0	12.0	1.0	NaN	NaN	NaN	NaN	0.07	

The dataset

Showers are from the Prague Library (GAP 2018-043) $E_{\rm MC} \propto E^{-1} \, {\rm and} \, \theta_{\rm MC} \propto \sin \theta \, \cos \theta$

Hadronic models: EPOS-LHC + FLUKA

Conservative background: protons with muon-deficit

Quality Cuts:

- $16.5 \le \log_{10} \left(\frac{E_{\rm MC}}{\rm eV}\right) \le 17.5$
- $\theta_{\rm MC} \le 45^{\circ}$

Events: approx. 45k (balanced classes)

Missing values: 70%

Initial stratified split

<u>Training set</u> (2/3 original dataset)

- Model training
- Hyperparameter tuning
- Model Selection

<u>Testing set (1/3 original dataset)</u> • Unbiased Performance Estimation

Feature extraction

Normalized muon densities

$$\rho_{\mu \text{ norm}}^r = \frac{\rho_{\mu}^r}{\rho_{\mu \text{ ref}}^r}$$

Parameters are fitted from a separate set of <u>CORSIKA simulations</u> for fixed combinations of $E_{\rm MC}$ and $\theta_{\rm MC}$

$$\rho_{\mu_{\text{ref}}}^{r} = R_{\mu}^{r}(\theta) \left[\frac{E}{10^{16.7} \text{eV}}\right]^{k_{\mu}^{r}}$$
$$R_{\mu}^{r}(\theta) = ax^{2} + bx + c$$
$$x = \sin^{2}\theta - \sin^{2} 30^{\circ}$$

Distributions from Offline reconstructions



Ezequiel E. Rodriguez – ezequiel.rodriguez@iteda.cnea.com.ar

The best contender – eXtreme Gradient Boosting (XGBoost)

- Introduced in 2016 (22k citations until last week)
- Highly optimized tree-based ensemble
- Handles missing values with Sparsity-aware Split Finding



Input: *d*. feature dimension Also applies to the approximate setting, only collect statistics of non-missing entries into buckets $qain \leftarrow 0$ $G \leftarrow \sum_{i \in I}, g_i, H \leftarrow \sum_{i \in I} h_i$ for k = 1 to m do // enumerate missing value goto right $G_L \leftarrow 0, \ H_L \leftarrow 0$ for j in sorted(I_k , ascent order by \mathbf{x}_{ik}) do $G_L \leftarrow G_L + q_i, \ H_L \leftarrow H_L + h_i$ $G_R \leftarrow G - G_L, \ H_R \leftarrow H - H_L$ $score \leftarrow \max(score, \frac{G_L^2}{H_L + \lambda} + \frac{G_R^2}{H_D + \lambda} - \frac{G^2}{H + \lambda})$ end // enumerate missing value goto left $G_R \leftarrow 0, H_R \leftarrow 0$ for j in sorted(I_k , descent order by \mathbf{x}_{ik}) do $G_R \leftarrow G_R + g_i, \ H_R \leftarrow H_R + h_i$ $G_L \leftarrow G - G_R, \ H_L \leftarrow H - H_R$ $score \leftarrow \max(score, \frac{G_L^2}{H_r + \lambda} + \frac{G_R^2}{H_n + \lambda} - \frac{G^2}{H + \lambda})$ end end

Algorithm 3: Sparsity-aware Split Finding

Input: *I*, instance set of current node **Input**: $I_k = \{i \in I | x_{ik} \neq \text{missing}\}$

Output: Split and default directions with max gain

Chen, T., & Guestrin, C. (2016, August). *Xgboost: A scalable tree boosting system.*

Model Training

Stratified 5-fold Cross-Validation (CV) Low-level optimization task



<u>Mean</u> Area Under the Curve (AUC) for Background Rejection (BR) Vs Signal Efficiency (SE)



Savina, P., Bleve, C., & Perrone, L. , PoS (ICRC2019), 414.

Training – Hyperparameter tuning – Model selection



Out-of-sample performance estimation with bootstrap



As we contemplate events with lower primary energy performance drops

Ezequiel E. Rodriguez - ezequiel.rodriguez@iteda.cnea.com.ar

Out-of-sample performance estimation with bootstrap



As we contemplate events with higher zenith angle performance drops

Summary and Outlook

- Developed observables as input for to MVA classifiers
- Initial model development with "traditional" Machine Learning models
- Estimation of model performance

3	2022 2023				
Name	Begin date	End date	ov 'Dec 'Jan 'Feb 'Mar 'Apr 'May 'Jun		
MODEL EXPLORATION	11/9/22	12/30/22			
SENSITIVITY / ROBUSTNESS	2/1/23	4/28/23			
BACKGROUND STUDIES	3/31/23	6/1/23			
PHOTON SEARCH	6/1/23				

Ezequiel E. Rodriguez - ezequiel.rodriguez@iteda.cnea.com.ar

Backup slides

Ezequiel E. Rodriguez – ezequiel.rodriguez@iteda.cnea.com.ar

Bayesian Optimization - Tree Parzen Estimator (1)

- Sequential Model-Based Optimization(SMBO) -> Tree Parzen Estimator (TPE)



Figure 1: The pseudo-code of generic Sequential Model-Based Optimization.

$$p(x|y) = \begin{cases} \ell(x) & \text{if } y < y^* \\ g(x) & \text{if } y \ge y^*, \end{cases} \quad \gamma = p(y < y^*)$$

Bayesian Optimization - Tree Parzen Estimator (2)



Ezequiel E. Rodriguez – ezequiel.rodriguez@iteda.cnea.com.ar

Results from Bayesian Optimization



Ezequiel E. Rodriguez - ezequiel.rodriguez@iteda.cnea.com.ar