Boosted Decision Tree Model for LEGEND

Henry Nachman PIRE-GEMADARC Collaboration Meeting Academica Sinica – Taipei, Taiwan 02 June 2023



Outline

- 1. Brief Introduction
 - I. What is $0\nu\beta\beta$?
 - II. What is LEGEND?
- 2. Parameter Extraction
 - I. Waveform Selection
 - II. Pulse Shape Parameter
- 3. BDT Model Training and Validation
- 4. Results, Discussion, and Conclusion





Neutrinoless Double Beta Decay - $0\nu\beta\beta$

$(A, Z) \to (A, Z+2) + e^- + e^-$

- Exchange of virtual light neutrinos results in no net neutrinos from this process
- Thus, net matter production
- If this process exists, it would be extremely rare orders of magnitude longer than the age of the universe.



3





LEGEND

- LEGEND is a next generation collaborative experiment searching for $0\nu\beta\beta$ using high purity germanium detectors.
- Currently in the 200kg phase, LEGEND plans to use ~ 1000 kg of ^{76}Ge to search for this elusive process.
- Innovation for this high mass and low-background experiment will allow LEGEND to achieve a discovery potential at a half-life of $10^{28}\rm{years}$





Background in LEGEND-1000

- To achieve its discovery potential L-1000 strives for extremely low background.
- Only ~3 $0\nu\beta\beta$ signals would constitute a 3σ discovery
- Many background reduction methods:
 - Low background materials
 - Active veto liquid argon scintillator
 - Analysis based background discrimination
 - Boosted Decision Tree (BDT) Model



Expected sources, and levels of backgrounds in L-1000





Inverted Coaxial Point Contact (ICPC) Detectors



Larger size of ICPCs makes many pulse shape characteristic correlations (discussed later) more prominent





Single-Site Events



Multi-Site Events



Pipeline Goals -

- 1. Develop a Machine Learning (ML) software tool that can reject multisite events using only the raw waveform data.
- 2. Create a fast (real time) analysis tool for efficient detector characterization analysis.
- 3. Explore and develop explainability tools to learn from and inform physics from the ML model.

Pipeline Steps -

- 1. Calibration and event selection
- 2. Pulse Shape Parameter extraction
- 3. Boosted Decision Tree (BDT) training and validation



Spectrum Calibration



Proxy Events

- To train a model to classify events as single-site (SS) or multi-site (MS) we need large samples of each
- Of course $0\nu\beta\beta$ signals are not exactly at our disposal.
- Use data from the Th-228 calibration as proxy to represent the types of signals we would expect to see.
 - Double-Escape Peak (DEP) for Single-Site/ $0\nu\beta\beta$ -like
 - Single-Escape Peak (SEP) for Multi-Site/background-like







Single-Site Proxy Event

2614 keV

THE UNIVERSITY of NORTH CAROLIN/



Multi-Site Proxy Event

2614 keV



Double Escape Total Collected Energy: 2614-511-511 = 1592 keV Single Escape Total Collected Energy: 2614-511 = 2103 keV



Event Extraction



Pulse Shape Parameter

LEC

THE UNIVERSITY of NORTH CAROLINA of CHAPEL HILL

	144	
	145	
n(energies)}")	146	Returns:
	147	- y_out: mean
	148	
	149	[tau1, tau2, f] = popt
	150	y_out = np.zeros(wfArray.shap
.187, # 228Th -> 208Tl (85%)	151	max_amp = 10
1173.24, # 68Co	152	
1332.5, # 60Co,	153	<pre>for wf in range(wfArray.shape[0]):</pre>
1592.5, # 228Th DEP	154	<pre>wf_in = wfArray[wf,:]</pre>
2103.5, # 228Th SEP	155	<pre>wf_in = (wf_in-wf_in.min())/(wf_in.max()-wf_</pre>
2614.53]) # 228Th -> 208Tl (99.8%)	156	wt_out = np.zeros(len(wt_in))
DEP" in targetPeak:	157	# Defines the constant terms
<pre>print(f"Calibrating on 228Th DEP")</pre>	158	const1 = 1 / tau1
peakIndex = 2	159	const2 = 1 / tau2
"SEP" in targetPeak:	160	# Defines the exponential terms
<pre>print(f"Calibrating on 228Th SEP")</pre>	161	expl = np.exp(-1 / taul)
peakIndex = 3	162	exp2 = np.exp(-1 / tau2)
"FEP" in targetPeak:	163	
<pre>print(f"Calibrating on 228Th FEP")</pre>	164	trac = t
peakIndex = 4	165	# Sets initial of output to same as input
	166	wf_out[0] = wf_in[0]
	167	e1 = e2 = wf_in[0]
s = np.array([583.187, # 228Th -> 208Tl (85%)	168	es = 0
2614.53]) # 228Th -> 208T1 (99.8%)	169	
	178	tor 1 in range(1, len(wt_in), 1): # iterates
.ns, var = pgh.get_hist(energies, bins=1000)	1/1	ei += wf_in[i] - ez + ez*consti
eaks, cal_peaks, cal_pars = pgc.hpge_find_E_peaks(hist, bins, var, peaks)	172	es += wt_in[i] - ez - es*constz
	1/3	ez = wr_in[i]
:h_peaks(data_pks, cal_pks, plotBool=plotBool):	174	wf_out[1] = e1 - frac*e3
	175	and a second of the reduced and
th uncalibrated peaks with literature energy values.	172	maxing = np.argmax(wr_in) # index or max
	170	pre_wr = wr_in[:maxinu] * waveform before ma
litertools import complications	170	sook - as max(uf in)
scipy.stats import linnegress	100	peak = np.max(wt_in)
	191	valo7 = 0.97*pask
(s = len(cal_pks) if len(cal_pks) < len(data_pks) else len(data_pks)	182	difface = no absolute(wf in - val97)
cate - combinations(names(lan(cal ake)) - a ake)	183	closestFromPeak = diffArr[neakInd:].argmin()
sets = combinations(range(len(cai_pks)), n_pks)	184	ind97 - neakInd+closestEnomPeak
_sets = compinations(range(ien(data_pks)), n_pks)	185	$v_{out}[wf] = (n_{out}(wf_{out}(ind97):1))$
ann hast m hast h - nn inf Nona Nona	185	Jonefull - (ubrace/ur_onef(ruga/)/1//
i cal sat in noumenate(cal sats):	187	return no mean(v out)
i, cai_set in enumerate(cai_sets):	188	()_ord
cal = cal nke[list(cal set)] # lit anargies for this set	189	***************************************
car - car_pha[rist(car_set/] # rit energies for this set	198	
s data set in data sets:	191	
	192	def dp0Vis(popt, wfArrav):
data perilist(data set)] # uncal energies for this set	193	[tau1, tau2, f] = popt
are professioned action and a chergina intersister	194	wfinAdi = np.zeros(wfArray.shape)
linregress(data, v=cal)	195	wfCorr = np.zeros(wfArray.shape)
(m * data + b))**2)	196	max amp = 10
		trash = []
	198	for wf in tqdm(range(wfArran,
ser. m. b	199	desc="
	200	CD3
Henry Nachman	201	wf in = wf 14
	202	if wf

What are Pulse Shape Parameters?

- Lots of information about a detection event can be gleamed from the shape of the event waveform.
- Many of these parameters were developed for analysis of previous Ge detector experiments, namely the MAJORANA DEMONSTRATOR.
- PSD Extracted in this pipeline:
 - Current Amplitude $\stackrel{\circ}{\multimap}$ Energy (A/E)
 - Delayed Charge Recovery (DCR)
 - Late Charge (LQ80)
 - Drift Time 10%, 50%, 99.9% (TDRIFT)



Pole-Zero Correction

- Raw waveforms exhibit an RC decay in the tail from the electronic circuit used to read out charge from the detector
- Can be fit with a two-term • exponential decay function and corrected



Pole-Zero Correction Visualization



THE UNIVERSITY of NORTH CAROLIN CHAPEL HIL

Drift Time





Current Amplitude



• Divide by DAQ energy to get : A/E

THE UNIVERSITY of NORTH CAROLINA





Tail Slope – Delayed Charge Recovery



Late Charge





High LQ Event





Parameter Correlations

- Traditional Analysis:
 - Each pulse shape parameter is considered independently with little to no consideration for multi-parameter correlation.
 - Requires fine-tuning of each individual parameter (time consuming).
- Machine Learning Analysis:
 - Simultaneously analyzes all parameters taking into account their correlations.
 - Does not require fine-tuning of individual parameters (fast analysis, ~30 sec training time)
- Motivates the use of a Boosted Decision Tree machine learning model which can train on many different parameters simultaneously.





Boosted Decision Tree



What is a Boosted Decision Tree Model

Henry Nachman

- Supervised machine learning model
 - Machine Learning Computer algorithms that learn from large datasets to fit models to observed patterns.
 - Supervised = labeled. Essentially guiding the machine to specified classifications
- Decision Tree
 - Network of Boolean (T/F) decisions through which events are classified
- Boosting
 - Growing multiple trees in a series, each fitting to the residuals of the previous.

Fig: Screenshot of part of one trained tree.



Distribution Matching

- The DEP and SEP datasets have different distributions of parameters which need to be remedied to avoid unintentional bias:
 - 1. DATASET SIZE if the SEP (Multi-Site) has many more events than DEP, the model will have a statistical bias towards classifying things as multi-site.
 - Need to remove ENERGY DEPENDENCE from parameters. Our datasets come from two peaks of different energy – not a property of single-site and multi-site events. Otherwise, model will simply train on energy and will be unrealistically good at classifying events.
 - 3. Encourages the model to investigate the correlations between parameters, rather than just differences in distributions of single parameters.





How to Match Distributions

- 1. Bin the data from both sets (singlesite & multi-site)
- 2. Keep only the data that is consistent between both sets
- Drastically cuts the totally number of events available for training.
- Don't want to match A/E parameter as dividing by E should remove Energy dependence and other distribution differences are likely signatures of multi-sitedness







Model Explainability - SHAP

- Explain how the machine learned from the data rather than blindly trusting the results
- SHAP SHapley Additive exPlanations game theory principles to assign each parameter a value based on its influence over the model's classification, for each event.



External Dependencies

My pipeline makes use of external software frameworks including:

- PYGAMA for peak fitting
 - Developed by LEGEND collaboration
- LightGBM for model training
 - Guolin Ke et al.
- SHAP for Shapley explainability study
 - Scott Lundberg, Su-In Lee



BDT Distribution

Mix Data	Pre- Matching	Post- Matching
SEP (Multi-Site)	35420	3899
DEP (Single-Site)	14730	3899
Trees	47	

• Model validated on 30% of training data.

THE UNIVERSITY of NORTH CAROLIN/





ROC – mix data

 Measure of binary classifier – Perfect classifier has an AUC = 1

 Note: the A/E measure shown here is without drift time correction or other methods used in the traditional LEGEND analysis





SHAP Feature Importance - Mix

- Made by plotting Shapely values for every event.
- Color of point represents the ultimate classification.
- "Important" features have a dumbbell shape
- Notice the anti-correlation between
 - TDRIFT50:A_DAQE
 - DCR:A_DAQE





Discussion and Conclusion

- BDT does outperform (raw) A/E discriminator but likely does not rival fully-tuned A/E in traditional analysis
- Fast analysis tool Whole pipeline can be run in a matter of minutes from raw data to trained BDT model.
 - Particularly important as LEGEND ramps up in detector characterization with \sim 250 detectors to be analyzed in a short time frame.

Next Steps

- Multi-detector support.
- Further explainability studies.
- Incorporation of Solid State Detectors pulse shape simulation data.







Questions?

Correspondence: Henry.Nachman@unc.edu







Why is the Universe Matter Dominated?



Standard Model of Particle Physics

- Within the Standard Model total baryon and lepton number is conserved.
- For every matter particle, we expect to see an antimatter particle
- There must be a process that breaks the balance matter and antimatter to account for observed asymmetry



Standard Model of Elementary Particles

Image Source: Wikipedia



THE UNIVERSITY

Enter Neutrino – ν

• Neutrinos are:

CHAPEL HIL

- Neutral subatomic particles
- Fermions (half-integer spin)
- Come in 3 flavors



- Importantly
 - Observations of neutrino oscillations show that neutrinos have MASS.
 - Discovery for which Takaaki Kajita & Art McDonald received the 2015 Nobel Prize in Physics.



Majorana Particle

 These characteristics of neutrinos suggest it could be a Majorana Fermion – or its own antiparticle.

$\nu = \overline{\nu}$

- If the neutrino is Majorana, this could allow for baryogenesis (from neutrino leptogenesis)
- One possible explanation for matter antimatter asymmetry.



Ghosthunting for Neutrinos

- Lack of charge, and extremely low interaction cross-section make direct neutrino measurements very difficult.
- Instead we can use the theorized process known as *Neutrinoless Double Beta Decay* to probe Majorana nature of neutrinos.
- Observing $0\nu\beta\beta$ would prove that neutrinos are Majorana.





Standard Model – Double Beta Decay

$$(A, Z) \to (A, Z+2) + e^- + e^- + \bar{\nu}_e + \bar{\nu}_e$$



Observed in 11 different naturally occurring isotopes.

Among them: (Spoiler Alert) ^{76}Ge







Large Enriched Germanium Experiment for Neutrinoless ββ Decay



THE UNIVERSITY

CHAPEL HILI





Innovation toward LEGEND-1000



MAJORANA DEMONSTRATOR (MJD) : PPCs, low noise electronics



GERDA : LAr veto, water shield





LEGEND-200 : Now taking data



LEGEND-1000 : Conceptual design development continuing

Characterization

 Prior to installation, detectors undergo tests to confirm appropriate energy resolution, pulse shape reconstruction and other qualities







Characterization Data

- During characterization: detectors are biased (given a voltage gradient) and take data with a Th-228 source.
- Data for this model is taken over two different runs:
 - Top Source
 - Side Source
- All data used in this model from a single ORTEC ICPC detector (V01387A)
- Characterization performed by Morgan Clark at Oak Ridge National Lab on June 17, 2021.





Scaling up for L-1000

- As the LEGEND experiment prepares for L-1000 many detectors will need to be efficiently characterized
- At peak ~ 2 detectors/week
- Motivates a need for a fast analysis tool
 - Such as the BDT model I developed.



Fit to calibration peaks



Single-Site Proxy Event



Multi-Site Proxy Event

2614 keV





Single Escape Total Collected Energy: 2614-511 = 2103 keV

What Parameters are Extracted?



Current Amplitude





Energy

- Three different methods for estimating the energy of an event were developed
 - Energy from data acquisition (DAQ) system simple trapezoidal filter
 - Offline trapezoidal filter
 - Maximum of the waveform
- Ultimately DAQ Energy had a comparable distribution to the other methods and was chosen for its simplicity in integration into the pipeline





What makes it Boosted?

- A single decision tree model is a "poor learner" on its own.
- Boosting growing multiple trees in a series, each fitting to the residuals of the previous.
- Validation Metric Binary Log Loss function
 - Minimizing residuals of mischaracterized events







Model Settings

Hyperparameter Name	Value	
Number of Iterations	Early Stopping: 10	
	$(\max 2000)$	
Validation Metric	Binary Log-loss	
Learning Rate	0.0744	
Number of Leaves	73	
Bagging Frequency	62	
Minimum data in Leaf	26	
Max Bin	542	
Drop Rate	0.330	
Min Gain to Split	0.536	
Boosting Algorithm	Gradient-based One-Side Sampling (GOSS)	

Build trees until the metric does not improve for 10 iterations



ŁE



Training

- 2 models : 1 with only top data, 1 with a mix of top and side data (with data augmentation)
- Data is split into training : validation subsets with 70:30 ratio.

Top Only	Pre-Matching	Post-Matching
SEP Size	17710	2062
DEP Size	7365	2062
Trees	76	

Mix Data	Pre-Matching	Post-Matching
SEP Size	35420	3899
DEP Size	14730	3899
Trees	47	





Data for Training

- Multi-Site Proxy : Single Escape Peak (SEP)
- Single-Site Proxy : Double Escape Peak (DEP)
- 2 Datasets : Top Source (Flood) Data & Side Source Data

Dataset	# Events in SEP	# Events in DEP
Top Source	25300	10521
Side Source	4712	2097





Top and Side Source Parameter Distribution



A/E Distributions

Data Augmentation

- Much more data exists for top source vastly differently sized datasets can result in unintentional bias in the model.
- SMOTE-NC: Artificially synthesize new data points by interpolating between existing points in 6-dimensional parameter space.





Background Events & Sideband Subtraction



60

BDT Distribution – Top data





Receiver Operator Curve – Top data

Visualization of binary classification

- True Positivity Rate (rate of events correctly classified as SS) vs False Positivity Rate (falsely classified as SS)
- Perfect classifier would have a step function with an AUC = 1
- Random classifier would be a diagonal line with AUC = 0.5

Note: the A/E measure shown here is without drift time correction or other methods used in the traditional analysis





SHAP Feature Importance - Top

- Made by plotting Shapely values for every event.
- Color of point represents the ultimate classification.
- "Important" features have a dumbbell shape





Next Steps

- 1. Multi-detector support In Progress
 - a. Increases total data yield
 - b. Helps eliminate possible bias from "quirky" detectors
- 2. New (and improved) parameters
- 3. Explainability expansion
 - a. Multivariate correlation studies
- 4. Other ML models particularly those that look at the whole waveform, rather than just the provided parameters.



Avenues for Publication

- Start with an internal LEGEND technical document.
 - Let the larger collaboration know about the work I am doing and provide an opportunity for guidance, and feedback
- 1. Depending on this package's use in future characterization work possible incorporation into a paper on LEGEND's detector characterization efforts.
- 2. Contingent on reliable results on further explainability development particularly multivariate analysis – possible standalone ML paper



Binary Log Loss

$$H_p(q) = -\frac{1}{N} \sum_{i=1}^{N} y_i \cdot log(p(y_i)) + (1 - y_i) \cdot log(1 - p(y_i))$$

Binary Cross-Entropy / Log Loss









٦ĒĽ





