DATA SCIENCE IN EXPERIMENT

William Phelps

Christopher Newport University/Jefferson Lab





Data Science

- Heavily involves
 - Data
 - Statistics
 - Machine Learning
 - In particular this talk will be heavily biased towards ML
- This talk will focus on projects carried out by users and staff



Roadmap

- Deep Learning Partial Wave Analysis
- Machine Learning Detector
 Optimization
- Data Science Group Projects
 - Hydra
 - AIEC





t time



Partial Wave Analysis

- A python-based software framework designed to perform Partial Wave and Amplitude Analysis with the goal of extracting resonance information from multi-particle final states.
- In development since 2014 and has been significantly improved with each revision. Version 3.4 out now.
- Efficient amplitude analysis framework including multithreading and CUDA support
- Optimizers include: Minuit, Nestle (or add your own!)
- NIM Paper in development!

Website: https://pypwa.jlab.org GitHub: https://github.com/JeffersonLab/PyPWA



Group Members Carlos Salgado (NSU/Jlab)

Mark Jones (NSU) William Phelps (CNU/Jlab) Michael Harris (NSU) Andru Quiroga (CNU) Bruna Goncalves (NSU) Nathan Kolling (CNU)

Former Group Members

Josh Pond Stephanie Bramlett Brandon DeMello

PWA using Neural Networks

- Generate datasets using decay amplitudes (linear combination of spherical harmonics) with the following quantum numbers
 - L = 1,2,3
 - *m* = 0,1
 - $\epsilon_{R} = -1, +1$





Tools of the Trade

- Python 3.9 Anaconda
 - Keras/TensorFlow NN Libraries
 - Pandas/Numpy Data Handling
 - Matplotlib Visualization
- Many GPU nodes that Scientific Computing division has available
 - · Either through Jupyterhub or interactively using slurm to request a node
 - Several machines with 4 NvidiaTitan RTX GPUs and some with 14 NvidiaT4 GPUs

```
test = pd.read_csv("TRAIN/TRAIN.csv")
labels = pd.read_csv("TRAIN/TRAIN_labels.csv")
activation = 'relu'
model = Sequential()
model.add(Dense(units=1000, activation=activation, input_shape=(3600, )))
model.add(Dense(units=1000, activation=activation))
model.add(Dense(units=1000, activation=activation))
model.add(Dense(units=1000, activation=activation))
model.add(Dense(units=2))
model.compile(optimizer=adam(lr=.001), loss='mean_squared_error', metrics=['accuracy'])
```

model.fit(test, labels[labels.columns[1:]], epochs=300, batch_size=256, validation_split=0.2)







python

Early Results

- We compare the intensity function and compare it to the model prediction
- Model Architecture:
 - I28xI28 2D histogram as input
 - 9x128 Dense Layers Relu activation
 - 9 production amplitudes as output
- In order to deal with the vast amounts of data we used generators to generate data for each epoch on the fly



TensorFlow K Keras



Useful Tools: Generators, Complex Valued Deep Learning

Autoencoder for PWA

Unsupervised learning!



Autoencoders for PyPWA

Phi [rad]

Phi [rad]

- Encoder portion is a standard MLP, but without labels!
- Decoder is a PyPWA model that takes in production amplitudes and produces a histogram
- Autoencoders dramatically improved the accuracy!
- Even works well for noisy data



Mass Dependent Autoencoder work for PWA



The Mass-Dependent Generator



Results

- With a CONV3D input to our autoencoder we see a good agreement with the generated data and inference from our neural networks
- Shown on the right are three different tests with randomly generated data/resonances



PWA Summary

- We have been able perform PWA "fits" with neural networks
- Autoencoders dramatically improved the performance
- Future work includes uncertainty quantification, and we are currently investigating Bayesian Neural Networks, dropout during inference, and Variational Autoencoders

Many thanks to the EPSCI and Data Science group at JLab! David Lawrence, Thomas Britton, Malachi Schram, Kishansingh Rajput

Al-optimized design of the ECCE tracker

C. Fanelli, Z. Papandreou, K. Suresh

NEW

C. Fanelli, Z. Papandreou, K. Suresh et al (ECCE), Al-assisted Optimization of the ECCE Tracking System at the Electron Ion Collider https://arxiv.org/pdf/2205.09185.pdf (submitted to NIM-A)



Detector technology being finalized: more realistic material budget, simulation details (e.g. background), and constraints to be included

Plan to encode all this in the existing framework

Opportunity to explore other optimization strategies and scalability; more advanced distributed pipelines/workflow

Multiple objectives: different physics analysis

In few months we will make a decision on the SW framework (solution adopted by ECCE so far VS ATHENA one). With Regina I am planning to explore both coupling to Fun4All and to DD4Hep **NEW**

Al-optimized design of the ECCE tracker

C. Fanelli, Z. Papandreou, K. Suresh



Navigate the Pareto front solutions (different detector design points) interactively

https://ai4eicdetopt.pythonanywhere.com

(credits K. Suresh)

https://eic.ai

Ongoing Projects – Data Science Group

- Hall A SoLID: Kishansingh Rajput
 - Particle ID
- Hall D CCP: Nikhil Kalra
 - Particle ID
- Hall D AIEC calibration: Diana McSpadden
 - Online calibration using Gaussian Processes
- Hydra: Thomas Britton
 - Online Monitoring

row

40

20

0

Thomas Britton

Hydra – Detector Monitoring

- Plenty of detector issues that are not alarmable in the traditional sense but still detectable
- Every run produces an initial >22 plots. More thorough monitoring is performed offline and produces >109 plots. With a run lasting ~2-3 hours every day there are between ~175 and 875 plots to look at every day. Desire Al to augment monitoring.







LED Light pulser left on

HydraRun also saw the FDC problem, which I probably would have missed inspecting it by eye.

Diana McSpadden

AIEC: Stabilizing Gain in the Central Drift Chamber

- Accelerate the calibration from months to minutes.
 - L. Gain Correction Factor: CDC Voltage Gain calibration
 - 2. Time to Distance: track fitting calibration
- Calibration is required to provide reliable PID for physics analysis
- Considerations:
 - External environmental conditions (temperature, pressure)
 - Changing beam conditions (current)
- Features used:
 - Atmospheric pressure
 - Temperature
 - Current drawn from CDC high voltage board.





Summary

- There are many projects going on at the lab but I was only able to mention a few
 - Al Townhalls have been held in 2020 and 2021 and are a great way to see what the community is doing
- If you are interested in participating keep an eye out for Jlab/EIC hackathons and workshops (<u>https://eic.ai</u>)



AIEC: Stabilizing Gain in the Central Drift Chamber

Why?

- Accelerate the calibration from month(s) to minute(s).
 - L. Gain Correction Factor: CDC Voltage Gain calibration
 - 2. Time to Distance: track fitting calibration
- Calibration is required to provide reliable PID for physics analysis
- Considerations:
 - External environmental conditions (temperature, pressure)
 - Changing beam conditions (current)



Goals: Near real-time control and calibration

- Traditional Calibration constants are generated per run (~2 hours) due to changing environmental conditions
 - Requires significant computing time, attention from experts

Goal: Using environmental and experimental data:





Environmental and Experimental Data

- Began project with 122 available features
- Feature evaluation reduced to 3.
- Selected features are readily available from EPICS system in near real time.
- Training and Test Data:
 - 536 runs from 2020 run period
 - 65 runs from 2021 run period
 - "Pressure balanced" train and test sets:
 - **Equal** proportions of low, medium, and high atmospheric pressures
 - Training: 480 runs
 - Test: 121 runs



Example Shapely value analysis for feature important evaluation.



Data Science Technique

Gaussian Process

• **3 features** (after feature importance evaluation of over

100 features):

- Atmospheric pressure
- Temperature
- Current drawn from CDC high voltage board.
- **GP** calculates the PDF over **all admissible functions** that fit the data
- GP provides standard deviation: used for uncertainty quantification



Cosmic Test: Can we stabilize gain in near real-time?

- Split the CDC into 2 halves
 - One side with **fixed HV** (traditional setup)
 - Let a Gaussian Process control the other side

Cosmic Experiment:

- •Update the HV every 5 minutes using GP
- Completely autonomous
- All information logged in database
- Remote monitoring of HV settings



Schematic of downstream view of CDC, with straws HV control status indicated.

Results: Yes, we stabilized gain

 Compared MPV of the peak height, ADC units (no magnetic field = no dE/dx) on cosmic muon values during each cosmic run for both sides of the CDC.



 Peak heights from Gaussian Process side of the CDC show dramatic reduction
 in pressure dependence

Ongoing Work: Using UQ

- Gaussian Process provides
 uncertainty quantification.
- Do not want to adjust high voltage to an "uncertain" value.
- We only apply a new calibration if the uncertainty is within the 3% of ideal gain correction factor otherwise, we extract to the closest prediction within tolerance

This method will be implemented in the upcoming GlueX CPP run.



Hydra

•Plenty of detector issues that are not alarmable in the traditional sense but still detectable



Every run produces an initial >22 plots. More thorough monitoring is performed offline and produces >109 plots. With a run lasting ~2-3 hours every day there are between ~175 and 875 plots to look at every day. Desire AI to augment monitoring.

Wednesday, June 15th, 2022



Hydra (an anecdote)



•The labeler was instructed by the detector expert to label any plot containing **fewer than 100k events as "NoData"**. This is one example of several in which the labeler labeled as "Good" and the **A.I. predicted "NoData"**...the true label given the number of events

Real world result: HydraRun also saw the FDC problem, which L probably would have missed inspecting it by eye. There were no alarms



BECAL PID Study A. Quiroga, W. Phelps, C. Fanelli, and J. Huang

- Higher pion rejection compared to conventional methods when considering high electron efficiency (~95%)
- Work is in progress (started on thanksgiving)
- Interface with ECCE software: reco-track, track projection, 7x7 calorimeter towers near track (track-based clustering by AI)
 [Link to details]
- Many models tried: MLPs, CNNs, Multi-Input models, Autoencoders, GANs.
- Ongoing hyperparameter tuning on 14 GPU nodes









BECAL PID Study A. Quiroga, W. Phelps, C. Fanelli, and J. Huang

- Higher pion rejection compared to conventional methods when considering high electron efficiency (~95%)
- Work is in progress (started on thanksgiving)
- Interface with ECCE software: reco-track, track projection, 7x7 calorimeter towers near track (trackbased clustering by AI) [Link to details]
- Many models tried: MLPs, CNNs, Multi-Input models, Autoencoders, GANs.
- Ongoing hyperparameter tuning on 14 GPU nodes







HRISTOPHER NEWPORT