Machine Learning

Yitong

4.29 update

- Tested the pulse classification plots & algs on PREM webpage
 - Most plots are showing properly
 - The plots with [SS] are not showing up -> special characters [] → change [] to ____
 - Re-group the pulse plots with their S1, S2, SPE, MPE, SE sub-plots \rightarrow change titles
 - Fraction algs are returning value 0 on PREM webpage
- The Hitchhiker's Guide to Machine Learning
 - NERSC videos on basic ML, Deep Learning, Neural networks
 - LZ ML Jupiter notebook tutorials
 - Search for the possibility of ML + data quality checks



Progression of laziness





Training



Too generalized

Too accurate to a specific training data set

Optimal Training

Loss Function Y[predict] vs. Y[real]



ML Alg

Deep Learning



Deep Learning







More hidden layers \rightarrow more chaining of non-linearity \rightarrow increase the complexity of the model



Learning actions from the environment

Initial Questions (data quality)

- What are we trying to model?
 - MDC3 simulation data (normal distribution data with label), 5 main populations? Each plot?
 - Possibly real data later (normal + anomalous data without labels)

What are we trying to classify?

- Normality vs. anomalies:
 - 1. Shift of the current population
 - 2. Disappearance of the current population
 - 3. Appearance of the new population

How are we going to classify anomalies?

• 1. Compare the test data against the trained normal set



Later Questions (data quality)

- How to evaluate that the model is doing a good job?
 - Choose the training epoch \rightarrow optimal fit
 - Computation time: quick check
 - A certain tolerance on some anomalies (i.e. noise...)
 - The determination of both abnormality of current model (shift & disappearance) & novel model (appearance)
- How do we know if this model will be productive on the data that we've never seen before?
 - Tests on real data and see if the output matches our expectations

Anomaly Detection Deep Learning

- To model normal behavior first, and then exploit this knowledge to identify anomalies.
 - An anomaly score is assigned to each data point → measurement of the deviation from the normal behavior
 - Determine the threshold: scores above a given threshold are tagged as anomalies; below this threshold are tagged as normal
- Three branches (supervised, unsupervised, semi-supervised)
- Five models (AutoEncoder; Variational AutoEncoder; GAN; Sequenceto-Sequence Model; One Class SVM)
- Data Quality & Physics Analysis (signal/background)

Branches

- Supervised: labels are available for both normal & anomalous sets
 - Both normal and anomalous datasets are labeled
 - Don't have labeled anomalous data
- Unsupervised: no labels for the training set
 - Both normal and anomalous datasets are not labeled
 - Learn by finding structure within the input features
- Semi-supervised: labels are only available for the normal set, but not the anomalous set
 - With the assumption: most data points within an unlabeled dataset are normal
 - Large amounts of unlabeled data + small amounts of labeled data
 - Detection of both known and previously unseen anomalies

Supervised Learning Training nml nml nml nml abn abn abn abn nml nml nml nml abn abn abn abn nml nml nml abn abn nm1 abn abn Unsupervised Learning Training ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? Semi-supervised Training ? ? ? ? nml abn nml nm1 ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ?

Autoencoders

- The model is trained to minimize the reconstruction error
- Dimensionality reduction technique



Autoencoder

Autoencoder



To model normal behavior

- Semi-supervised approach
- Train the autoencoder on normal data samples (sim data)
- Model learns a mapping function to reconstruct normal data samples successfully

To identify anomalies

- Reconstruction error score
 → Anomaly score
- Flag error if above a certain threshold score