# DS 6040: Bayesian Machine Learning: Project Proposal

H. Diana McSpadden (hdm5s)

## The Team

I will be working solo on this project. This is a risk and I am hoping for access to Dr. Basener, School of Data Science tutors, and even coworkers if I need active brainstorming and feedback throughout the project.

## Problem Statement and Background

Current work classifying waveforms from the particle accelerator beams utilizes KNN and Siamese models to identify "normal" and "anomalous" waveforms.  These waveforms are of interest because non-normal waveforms are symptomatic of malfunctioning equipment. My project dataset contains normal and anomalous waveforms produced from a time span before malfunctioning hardware was demonstrably malfunctioning, thus, classification is predictive of equipment degradation.

KNN models have been used with success in classifying these waveforms[1]. Existing Siamese models are better at finding anomalies than KNN models[2]; however, uncertainty quantification (UQ) is necessary because new/unknown presentations of anomalies in waveforms are expected, and UQ measures may identify waveform classifications with little certainty/similarity to previously trained waveforms. UQ has been built into the Siamese models; however, comparison of model uncertainty using Bayesian ML provides a more robust understanding of strengths and weakness of the data, models, and classification.
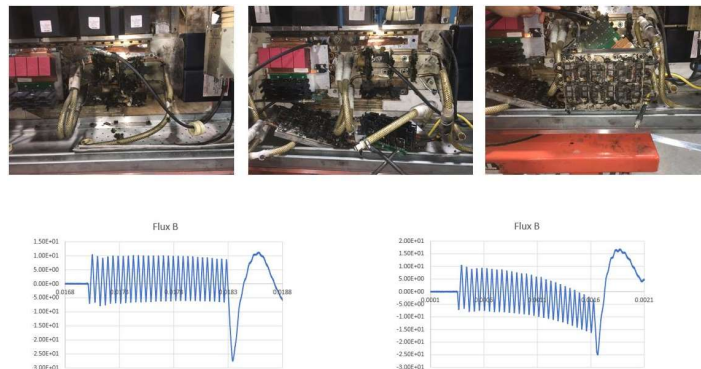


*Figure 1: Above: Oak Ridge National Lab photos of a damaged detector (it "exploded")  caused by errant beam. Below: system generated waveforms.*

## The Data

The data source is curated from ~25,000 "normal" waveforms, and ~3,000 "errant" waveforms provided in 32-bit array .npy files.

---

[1] Miha Rescic, Rebecca Seviour, Willem Blokland, Predicting particle accelerator failures using binary classifiers, Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment, Volume 955, 2020, 163240, ISSN 0168-9002, https://doi.org/10.1016/j.nima.2019.163240.
[2] Paper in process of publication.

Files are divided into existing train, validation, and test datasets. Train/test/validation splits must be preserved for comparison with other AI/ML methods.

The data are time series wave forms for a specific pulse duration. The waveforms are labeled "normal" or "anomalous" by filename encoding.

- Normal waveforms: "…ftype00000… .npy" files
- Anomalous waveforms: "…ftype00110… .npy" files

The .npy file are not public, so I will not be shared these files when I submit my project.

Exploratory data analysis will be performed on training, validation and test datasets.

## Goals of Analysis

1. Train a Bayesian classifier to identify normal and anomalous waveforms with training and validation data.
2. Test the classifier on test waveforms and produce performance metrics for the classification.
3. Use Bayesian methods to measure uncertainty for waveform classification.
   a. Ideally, generate UQ for all model parameters.
   b. Ideally, generate a similarity or uncertainty score for each test classification.
   c. Ideally, compare model UQ and prediction UQ to additional candidate models (e.g. the existing Siamese model)
4. Drawing on model, and classification UQ's, describe the strengths and weakness of the Bayesian model for comparison to UQ from other model types (e.g. the Siamese Model)

## What I Will Produce

1. A 5-minute project presentation
2. A four-page project report containing the problem, approach, results, and conclusions.
   a. Report will include:
      i. Description of the problem
      ii. Data evaluation
      iii. Any feature extraction and engineering
      iv. Modeling approaches
      v. Statistical tests
      vi. Model performance analysis including confusion matrices, ROC curves, and standard deviation error bar plots.
3. Code and model files as ipynb, .py, and .pkl