
INSIGHTS FROM APPLIED MACHINE LEARNING FOR SAFEGUARDING A PUREX REPROCESSING FACILITY

PROCEEDINGS OF THE INMM & ESARDA JOINT VIRTUAL ANNUAL MEETING
AUGUST 23-26 & AUGUST 30-SEPTEMBER 1, 2021

Nathan Shoman*
Sandia National Laboratories[†]
nshoman@sandia.gov

Benjamin B. Cipiti
Sandia National Laboratories[†]
bbcipit@sandia.gov

Thomas Grimes
Pacific Northwest National Laboratory
thomas.grimes@pnnl.gov

Benjamin Wilson
Pacific Northwest National Laboratory
benjamin.wilson@pnnl.gov

Randall Gladen
Pacific Northwest National Laboratory
randall.gladen@pnnl.gov

ABSTRACT

The International Atomic Energy Agency (IAEA) is seeking technologies that would improve the ability of safeguards to meet demands resulting from increased growth in nuclear industry. There are many different areas that could be improved, however, nuclear material accountability (NMA) is of particular interest. NMA for large throughput facilities often require expensive destructive assay (DA) measurements to meet safeguards goals. It would then be desirable to improve NMA through cheaper unattended measurements such as process monitoring (PM) and non-destructive assay (NDA). These sensors have high uncertainties that currently prohibit their direct use in a material balance. However, machine learning could provide a path towards effective utilization of these higher uncertainty measurements. There is a sizable body of literature on data driven approaches for anomaly detection, however, there are many unique considerations for application to nuclear facilities. This work provides concrete discussion and experimentation on several of the potential barriers to deployment of machine learning for improved NMA. Many of these barriers must be resolved before any effective real-world deployment can be realized.

1 Motivation

The International Atomic Energy Agency's (IAEA) Department of Safeguards published a R&D roadmap for requested technologies that will help their safeguards implementations remain effective and efficient [1]. One of the top R&D priorities, T.1.R1, calls for the development and introduction of an integrated system of instrumentation data processing and review, with high level of automation and with unified user interface. This is further reinforced by other R&D goals listed in the roadmap (e.g. section V.4) that call for enhancing safeguards effectiveness monitoring and evaluation.

One area that could see a significant benefit from improved safeguards is nuclear material accounting (NMA) at large throughput facilities. These facilities often require destructive assay (DA) measurements to achieve the low uncertainties required to meet safeguards objectives. High precision measurements are often costly and place a large burden on the IAEA. Safeguards at large throughput facilities could be improved through utilization of cheaper unattended measurement systems such as non destructive assay (NDA) and process monitoring (PM). These sensors are not directly used in the material balance calculation due to their higher measurement uncertainties. However, it is possible that data driven approaches, such as machine learning, could more effectively leverage these measurements to improve safeguards.

*Corresponding author

[†]Sandia National Laboratories is a multimission laboratory managed and operated by National Technology and Engineering Solutions of Sandia, LLC, a wholly owned subsidiary of Honeywell International Inc., for the U.S. Department of Energy's National Nuclear Security Administration under contract DE-NA0003525. SAND2021-9220C

Machine learning (ML) has seen a rapid increase in popularity due to the incredible results achieved in some domains. Many of the most visible successes have consisted of tasks often seen in computer vision and natural language processing. Although less visible, there has also been a wide range of success in deploying ML solution for anomaly detection. This work will highlight efforts made by the authors to identify and solve roadblocks for applied ML to improve NMA.

2 Background

The general goal of international safeguards is the timely detection of diversion of significant quantities (SQ) of nuclear material for weapons purposes and deterrence of such diversion by the risk of detection [2]. The trivial approach would be to simply count the items of interest. Indeed, for facilities where nuclear material is contained in discrete items, random sampling of the items is used to ensure material is present and untampered. However, detection of diversion becomes more challenging for large throughput facilities where material is in a bulk form (e.g. solutions at reprocessing facilities or powders at enrichment facilities). The next several sections will introduce properties of bulk material measurements, the tools traditionally used to detect loss in the measurements, and how machine learning differs.

2.1 Traditional Nuclear Material Accounting

A cornerstone of the traditional NMA approach is the material balance (MB) [3], which is sometimes called material unaccounted for (MUF) or inventory difference (ID). The MB is a statistical quantity that corresponds to a physical area within a nuclear facility which is informed by a subject matter expert to help reach the overall safeguards goal. The MB is calculated at regular intervals of time called the material balance period (MBP). The MBP is also determined based on safeguards requirement and expert opinion. The MB for any given time t (usually a multiple of the MBP) with several measurement locations i (where n is the total number of locations for a particular measurement type) is a simple calculation [3] described as follows: $MB_t = (\sum_{i=1}^{n_{in}} I_{i,t-1} + \sum_{i=1}^{n_{in}} Tin_{i,t} - \sum_{i=1}^{n_{out}} Tout_{i,t}) - \sum_{i=1}^{n_{in}} I_{i,t}$.

Terms include the total input transfer ($\sum_{i=1}^{n_{in}} Tin_{i,t}$), the total output transfer ($\sum_{i=1}^{n_{out}} Tout_{i,t}$) and total inventory ($\sum_{i=1}^{n_{in}} I_{i,t}$) at time t . Ideally $MB_t = 0$ for the case where no material has been removed. However, measurements used to calculate a MB always have some associated error which result in a non-zero MB even for a no-loss case. Typically safeguards measurements at location i at time t have a multiplicative error model described by Equation 1. There are several additional statistical tests and transforms that help identify losses even with non-zero MBs, which will not be discussed here.

2.2 Measurement Error

Conventional NMA utilize measurements of bulk materials combined with statistical tests in an attempt to detect potential material losses. The probability of detection will never be 100% due to imperfect measurements. For safeguards, the measurement error model is often multiplicative, as described in Equation 1.

$$\begin{aligned}
 M_{i,t} &= G_{i,t}(1 + S_i + R_{i,t}) \\
 &\text{where} \\
 S_i &\sim \mathcal{N}(0, \delta_S^2) \\
 R_{i,t} &\sim \mathcal{N}(0, \delta_R^2)
 \end{aligned} \tag{1}$$

$M_{i,t}$ is the measured quantity at location i and time t and $G_{i,t}$ is the true quantity of interest. This error model usually applies to all of the individual safeguards measurements.

The short-term systematic (i.e. epistemic) error, S_i , arises from measurement conditions or settings which are not changed for some period of time and vary in an unpredictable way and is difficult to reduce (e.g. calibration curves). In contrast, the random (i.e. aleatory) error $R_{i,t}$ varies in an unpredictable way under repeatable conditions and can be reduced through repeated measurements.

2.3 Machine Learning

Machine learning refers to algorithms that learn to perform a task without being explicitly programmed to do so. ML encompasses many traditional techniques like support vector machines (SVM) [4], random forest [5], and k-Means [6]. This work focuses on a specific subset of machine learning called deep learning. Perhaps one of the most well known deep learning methods is the feed forward neural network [7]. Each network is comprised of individual units called

artificial neurons. These neurons receive, process, and output signals to other connected units. These networks can perform complex tasks when careful selection of an objective function and optimization method is made.

Optimization of a neural network occurs in a process called training. Training examples are passed through a neural network to produce a prediction (i.e. $\hat{y} = f(x, \theta)$). The error between the prediction and actual value is then used to tune the neural network weights and biases (θ) such that the prediction improves with respect to the given training example. Larger neural networks with more parameters naturally require more training examples which can be problematic for some applications. Later this work will consider the impact of training dataset size on the performance of one particular deep learning method applied to safeguards data.

3 Problem Statement

The material balance (i.e. MUF or ID) will have some variance that results from uncertainty in the individual measurement components. Simple experiments can show that increases in the MB variance will cause a reduction in detection probability for a loss regardless of the specific anomaly detection method. Therefore, it would be desirable to develop an algorithm that could use higher uncertainty measurements more effectively than traditional NMA in order to detect a material loss. Such an algorithm would reduce costs and improve overall facility safeguards. This work hypothesizes that a machine learning algorithm could identify subtle changes in observed signals to detect anomalous behavior. In contrast to the traditional approach, the hypothetical ML algorithm would not directly quantify actinide inventories and would only classify observations as normal or off-normal.

The goal of this work is to describe how various parameters impact the performance of the proposed ML approach. Most evaluations in this work considered four different diversion scenarios at a fixed location. The loss scenarios are identified using the numbers 1 to 4 where 1 is the easiest to detect and 4 is the most difficult. Results in this work will often be reported by analyzing one of these loss scenarios. More information on how each algorithm trained on the scenarios is provided in section 3.2.

3.1 Process Modeling

Real-world nuclear facility data can be difficult to obtain so this work relies on process modeling to develop a training dataset. The focus of this work is to improve NMA at bulk facilities. A PUREX reprocessing facility was chosen due to the complex facility operations as an exemplar to test the various machine learning approach. Specifically the Separation and Safeguards Performance Model (SSPM) [8] is used to generate training, testing, and validation data. Although a PUREX model is used here, there exists a family of SSPM-based models which include UREX+, electrochemical reprocessing, fuel fabrication, and enrichment facilities. All SSPM models track elemental and isotopic material flows through various unit operations.

Performance of the ML algorithms developed here were evaluated on data generated from the SSPM. In practice, data would be collected from a variety of NDA and PM sensors for use in a ML algorithm. One potential application would be direct use of observed gamma spectra with a ML algorithm (as opposed to inversion of the spectra to derive mass). However, the computational overhead of converting simulated mass to simulated gamma spectra for a large set of data can be expensive. This work directly uses mass from SSPM with an error model similar to what would be seen in a gamma spectra (i.e. Equation 1) to more effectively explore the applied ML space. This is a reasonable proxy for what would be observed in practice as gamma spectra will also have random and systematic error terms. The evaluated SSPM datasets also include multiple isotopes which would be representative of multiple gamma peaks from different isotopes present in an observed gamma spectra.

3.2 Baseline Machine Learning Approach

There are two general categories of machine learning algorithms; supervised and unsupervised. Supervised methods require explicit examples of classes or behavior. An example of a supervised algorithm would be a convolutional neural network [9] that attempts to classify images as either cat or dog. In contrast, unsupervised algorithms do not require examples of classes and tend to classify data based on some readily observable metric such as data similarity. An example of an unsupervised algorithm would be using k-means to classify an image as a cat or non-cat.

Unsupervised algorithms do not rely on explicitly labeled examples of abnormal behavior and instead must use some proxy metric. The proposed approach assumes that a neural network can be used to learn a function $y = f(x, \theta)$ when given some data from SSPM, x , and a set of learned parameters, θ , to accurately predict some other quantity resulting from normal facility operation, y . During a material loss the learned function should no longer provide a good estimate of the measured quantity y resulting in a high prediction error $\hat{y} - y$. Consequently, the prediction error of the neural

network should be low during normal conditions and high under off-normal conditions. A graphical representation of the approach is shown in Figure 1.

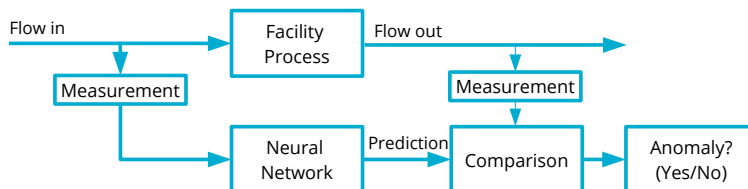


Figure 1: Proposed setup for applied ML for NMA

A threshold for the prediction error must be developed to discriminate between a "high" and "low" prediction error which is not initially defined. An overall metric of facility status can be difficult to define as the neural network generates several predictions corresponding to multiple isotopes and locations. It is important to determine what combination of prediction errors should be cause for alarm. A second machine learning algorithm, isolation forest [10], is employed to determine if prediction errors are sufficiently high to warrant flagging as a potential material loss event. This is performed by optimizing false alarm probability. Conceptually this algorithm, which is also unsupervised, assumes that anomalies should be few and different with respect to the bulk of the residuals. Therefore, anomalies should be easier to isolate spatially from the rest of the dataset compared to normal points. The algorithm works by partitioning the dataset into segments to generate a "forest" of decision trees.

Since the unsupervised approach has no knowledge of material losses, it is assumed that a certain percentage of the normal data is anomalous. This facilitates the development of a decision boundary and results in classification of some normal data with the largest measurement errors as abnormal. These misclassifications should be caused by random errors and variation and consequently randomly distributed through the sequence of observed data. However, under loss conditions, it is expected that there will be many anomalous classifications in a short window of time which would be an indicator of material loss. This work considers the unsupervised approach whereas the supervised approach is discussed in a separate work.

4 Data Requirements

Machine learning algorithms are often data hungry and require datasets that are representative of their deployment environments. Given the difficulties in obtaining real-world data it is important to understand how machine learning data requirements impact performance. The following sections aim to explore the relative importance of these requirements.

4.1 Facility Discretization

It is important to consider how data will be represented to a machine learning algorithm as representation can have a significant impact on performance. Early experiments attempted to predict the facility behavior using a single neural network with measurements from multiple locations and observations of multiple isotopes. However, it was discovered that there are unique behaviors in some unit processes that require special consideration. For example, some batchwise processes, such as the product evaporator, had trouble utilizing the same preprocessing techniques that were used with continuous operations (eg pulse separation column). This resulted in the creation of several regional neural prediction models rather than a single prediction model. This approach allows for more effective handling of unique properties of the various unit operations. Additionally, the regional models ensure that there is accurate predictions for the entirety of the MBA with the added benefit of some localization capabilities within the MBA.

4.2 Dataset Sizes

Collecting large amount of training data from actual facilities could will be expensive and difficult. Consequently, it is important to quantify the amount of training data needed for this ML-based approach. A parametric study was conducted to consider the impact of the training dataset size on performance. The size of the training dataset will impact both parts of the unsupervised approach; the neural network prediction and the Isolation Forest classification. The impact on overall system performance is considered (i.e. performance when both algorithms train on smaller datasets).

As a baseline, this work assumes 100 operational years of training data. Specifically, years of training data refers to operational years simulated by the underlying process model. Note that this is not merely 100 iterations of the same operational year (i.e. same pattern with different errors applied), but unique simulations. There are a wide range of potential facility patterns due to the many combinations of input fuel that could be selected. While likely not practical in a real-world deployment without a sufficiently robust process model, this quantity was chosen to have a sufficiently large baseline to infer optimal performance.

4.2.1 Combined Performance

This parametric study on dataset sizes considers the joint impact of dataset size on overall system performance. That is, the training dataset size for both the neural network predictor and isolation forest residual classifier were reduced together. This approach is more robust than considering the performance of each algorithm independently. During training, the classifier will be impacted by both a reduced training dataset size and lower quality data resulting from poorer neural network predictions. The joint performance of reduced datasets for training are summarized in Figure 2.

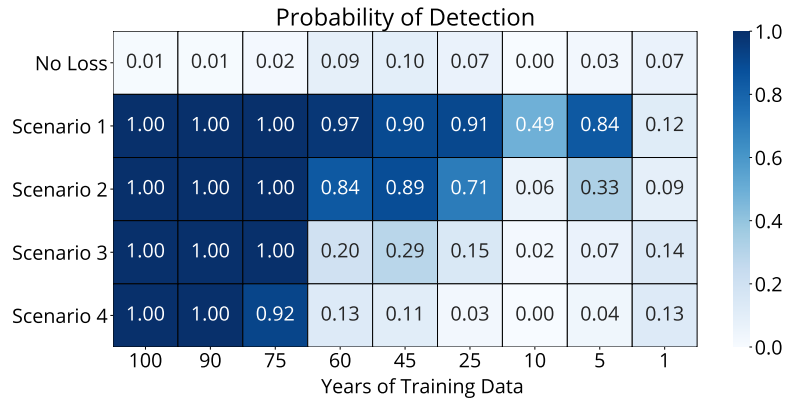


Figure 2: Probability of detection for several material loss scenarios with varied training dataset sizes. A probability of 1 indicates a 100% probability of detection.

It should be noted that in practice twice the quantity of data reported in Figure 2 will be required to achieve similar results. That is, one set of data would be needed to train the neural network and another would be required to generate training data for the isolation forest classifier. Simply training the neural network on a set of data then using residuals generated from the same set of training data would be biased. In effect, to obtain the performance of the 75 years of training data column, 150 years would be required.

Abrupt scenarios generally fare better than protracted scenarios when the size of training datasets are reduced. This is expected behavior as more subtle material losses become indistinguishable from normal behavior as the level of noise generated from imprecise machine learning predictions increase. Sharp drop offs are often observed rather than gradual decreases in probability of detection. For example, when the amount of training data is reduced from 75 years to 60 years there is a sharp drop in probability of detection for Scenario 3. This is because many of the runs no longer meet the static threshold that signals anomalous behavior.

4.3 Unexpected Transients

Facilities might sometimes have unexpected transients that deviate from expected behavior but are not malicious in nature. During these transients it may be very difficult to detect a second anomaly (e.g. facility transient in conjunction with a simultaneous material loss), however, an effective ML-based approach should be able to maintain the same level of performance after the facility has recovered from the transient. Two upset cases with independent causes and at different locations were simulated to evaluate the proposed ML approach. The performance is reported as the accuracy of the neural network in terms of prediction error which has a strong relation to anomaly detection capabilities. For example, high prediction errors after a transient has occurred will imply lower probabilities of detection for a material loss.

Generally, facility transients can be grouped into one of two categories. The first category includes transients that do not significantly alter the behavior of the facility itself. Examples in this category might include small leaks or changes to

product quality. The second category includes transients that do alter the behavior of the facility itself in some way. This may be a transient that causes a new equilibrium level in a surge tank or a unit operation that stops working and causes changes to operational timing. Only the second category will be discussed here as they have a stronger impact on overall system performance.

Both upsets considered fit into the second category where the facility behavior is changed. Specifically, the first upset case shown in Figure 3 represents a transient where the baseline of an area is changed. This could be a tank that has a new equilibrium level, for example. The neural network exhibits desirable behavior in that there is a strong error during the transient itself. However, after the transient has ended, the neural network provides erratic predictions and does not return to the baseline performance level. This is because the transient caused a new baseline for which the neural network was not trained on. Thus, after this transient, the training distribution no longer matches the evaluation distribution.

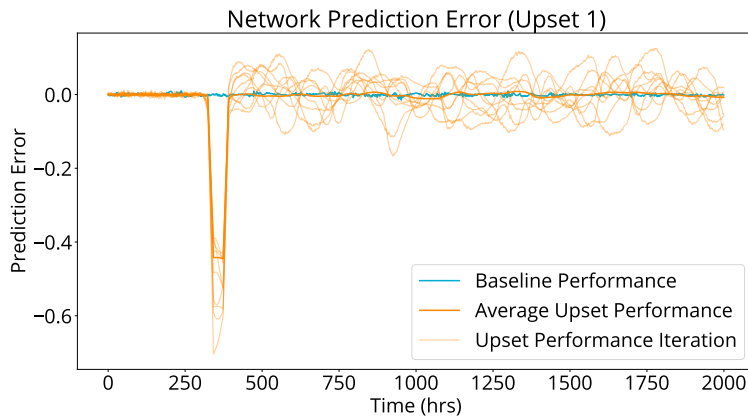


Figure 3: Prediction error of a neural network for a single isotope during a facility transient that changes the baseline of a facility section

The next upset case represents a facility transient that alters the behavior of a unit operation. This could be damage to an operation that causes sub-optimal performance or poor performance caused by damaged sensors. Figure 4 shows the prediction error when such a transient has occurred. Similar to the other cases, there is a sharp increase in prediction error at the start of the transient. However, after the transient has ended, there remains a large prediction error. This is because the facility is no longer operating according to the behavior the neural network learned in training.

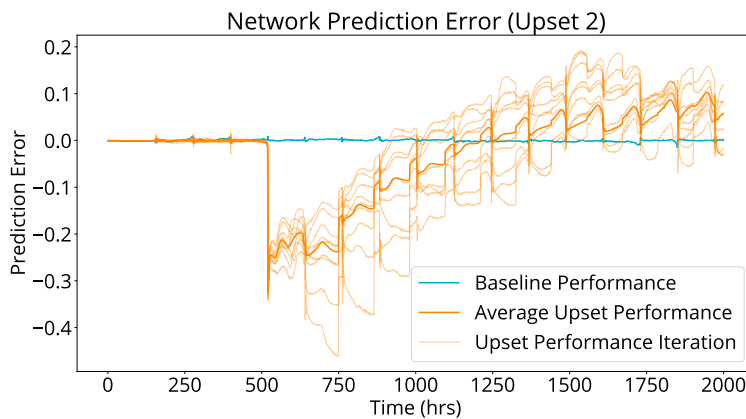


Figure 4: Prediction error of a neural network for a single isotope during a facility transient that changes facility behavior

The ML approach shows problematic post-transient behavior in the form of poor enduring performance. If uncorrected, this would degrade anomaly detection capabilities for a significant length of time after the transient has passed.

5 Measurement Error Impact

Detection of anomalous behavior (i.e. material loss) in measurements that have been contaminated with errors requires careful analysis. Many trivial approaches for change and/or anomaly detection fail to perform adequately when evaluating data contaminated with errors. It is helpful to reframe anomaly detection as a mean shift detection problem. Under normal operation measurements will have some mean and variance that is determined by the measurement technology. However, during a material loss, the mean of the feature will shift to a new mean. The magnitude of the shift will be dependent on the material loss itself.

Detection of this mean shift is dependent on the variance of the measurement itself. Consider the shift in the mean of a distribution from μ to μ^* . It stands to reason the probability of detection of that mean shift is a function of the variance σ such that $\lim_{\sigma \rightarrow \infty} PD(\mathcal{N}(\mu_t \rightarrow \mu_t^*, \sigma^2)) = \text{FAP}$ where FAP is the false alarm probability. It is important to note that no algorithm or approach can escape this phenomena. Certain approaches can fuse data to reduce σ , but they cannot overcome the limitation that process variance reduces the probability of detection for an anomaly.

Measurement error has a unique impact on unsupervised machine learning approaches. Recall that many ML approaches require large training datasets. It is reasonable to assume that a training dataset for any ML algorithm for a nuclear facility would be comprised of data from multiple measurement campaigns each with a different sensor calibration. The resulting aggregate training dataset would then capture changes in a signal of interest that arise from both facility variation and measurement error which leads to a larger variance than a dataset that captures facility variation alone. This causes a reduction in anomaly detection performance which often leads to worse overall performance than a traditional safeguards approach, even when both approaches utilize data with the same magnitude of measurement error. This is illustrated in Figure 5. Practical impacts are large and if not corrected for will completely degrade the ability of the ML approach to detect anomalies. More discussion in this phenomena is presented in a separate work.

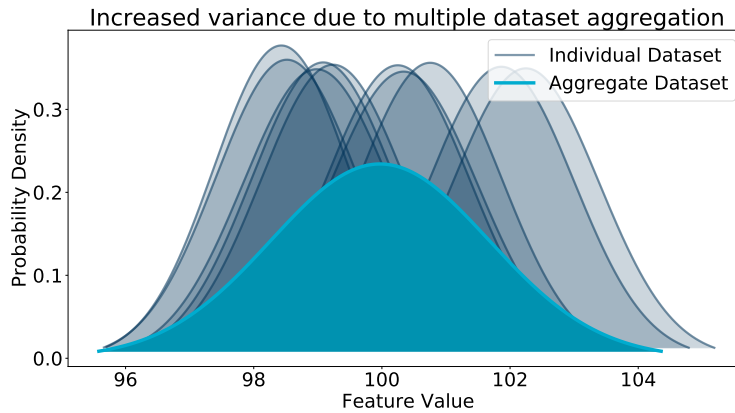


Figure 5: Probability density functions for multiple normal datasets

6 Threshold Decision

Recall that the proposed unsupervised machine learning approach discussed in Section 3.2 has to convert arbitrary prediction errors into a classification of normal/off-normal. Isolation Forest is used to perform this task. The expected behavior is that occasionally a normal point (labeled 1) will be classified as abnormal (labeled -1). However, during a constant material loss, it is expected that there would be a dense clustering of off-normal observations. Figure 6 shows the expected response of the second stage to a variety of scenarios. The abrupt losses tend to have a consistent cluster of off-normal classifications during the loss. As the loss becomes more protracted, the off-normal readings become more sparse. This is expected because the protracted loss decreases in magnitudes as length increases which represents a smaller shift from the baseline normal behavior.

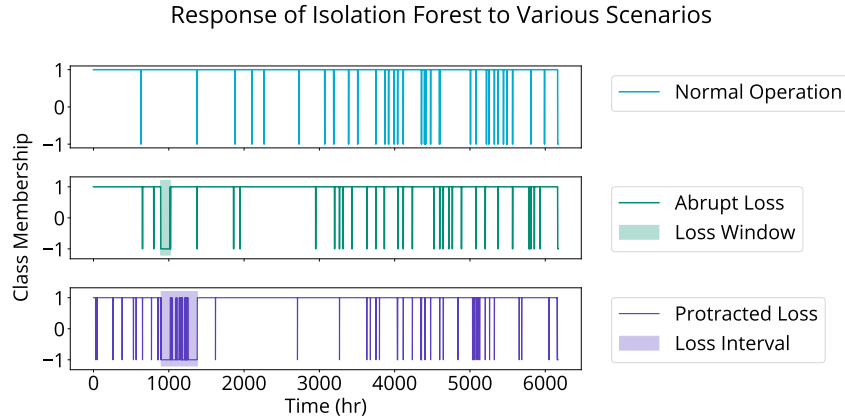


Figure 6: Isolation Forest response to various conditions

An intuitive metric would be to consider the density of off-normal classifications in a given window. However, determining this threshold can be challenging as there are two parameters to tune (length of window and density of off-normal classifications) and only one constraint (false alarm probability). Also maximizing for probability of detection would result in information leakage, as in practice, there would be no examples of material loss to optimize against. Since there are more parameters than constraints this leads to multiple possible combinations of alarm windows and density for a specific FAP threshold. Fortunately this ends up having little practical impact on the performance of the algorithm in most cases. A grid search was performed to develop a list of combinations of window density and window length to reach a FAP of 5%. These parameters were then used as evaluation parameters on the isolation forest algorithm. The largest impact on performance, summarized in Figure 7, occurred for shorter abrupt losses where the difference between the best parameters and worst could cause a 23% difference in detection probability.

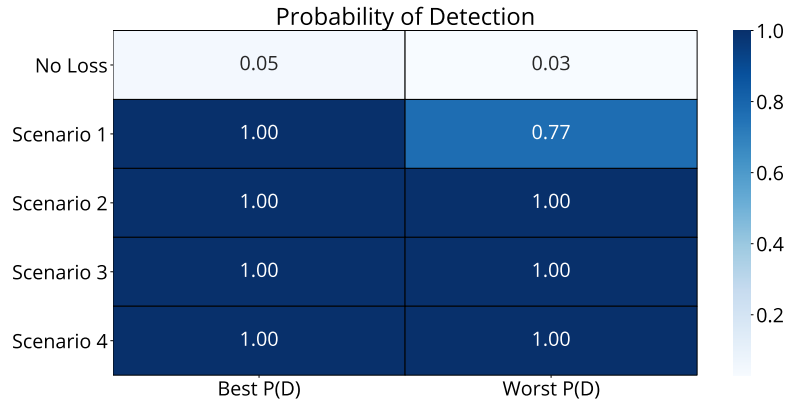


Figure 7: Probability of detection for several loss scenarios with varied thresholds

7 Machine Learning Performance

Given the constraints and limitations listed above it is possible to determine optimistic (i.e. ideal conditions) and pessimistic (i.e. worst conditions) performance levels of the proposed machine learning approach. Figure 8 summarizes the traditional and machine learning approaches under two situations. First, the 'uncalibrated' condition represents data as collected today. That is, sensors each have their own independent calibration curves (i.e. unique systematic errors). Next, the 'calibration' condition is considered which represents a situation where sensors are calibrated against each other. The calibration condition assumes that errors are present, but that the mean of sensors measuring the same feature is the same (i.e. systematic bias is the same for all sensors for a given area, but non-zero).

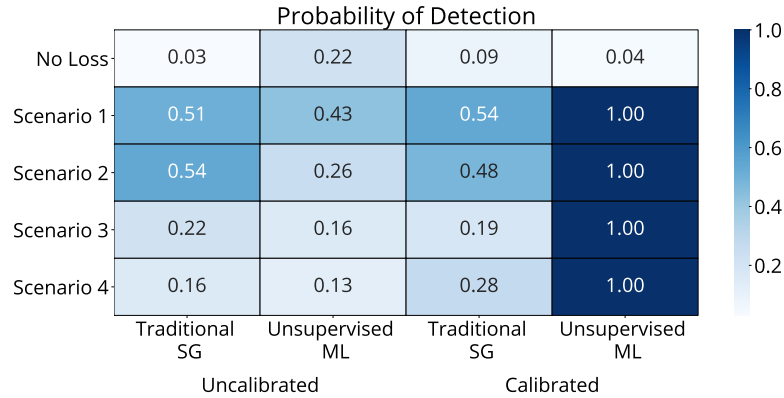


Figure 8: Detection probabilities for various scenarios

The change in the errors makes a substantial impact on the machine learning approach, but a negligible change on the traditional safeguards approach. This is related to some specific details of how statistical methods of traditional safeguards treat errors, which is addressed in a companion work. In this case the machine learning algorithm demonstrates performance well beyond what can be achieved with traditional safeguards. However, it is important to note that this change does not mitigate other problematic issues that can occur such as changes in baseline facility performance due to facility transients. Generally, while traditional safeguards exhibit lower performance in some situations than the machine learning approach, it tends to be more robust to measurement error and facility transients.

8 Conclusions

This work considered several practical matters that will impact the performance of an unsupervised machine learning algorithm. Current findings can be summarized as follows:

- Data representation (and facility partitioning) is important for effective performance
- Significant amounts of data are needed for adequate performance with this proposed machine learning approach
- Facility transients can transition observed signals to a new baseline which can cause poor performance from previously trained algorithms
- Measurement has a significant impact on the performance of a ML algorithm
- Correct selection of anomaly thresholds can have a significant impact on performance
 - Would require high quality synthetic or historical data

Machine learning has the promise to out-perform traditional safeguards, but is subject to several requirements. Measurement error tends to disproportionately impact the unsupervised machine learning approach and training data availability remains problematic. Given these limitations it is unlikely that ML approaches will wholly replace traditional safeguards in the near future. However, it could supplement existing NMA techniques. Further work is required to help mitigate or remove the barriers identified here.

9 Acknowledgments

This work was funded through the National Nuclear Security Administration’s Office of International Nuclear Safeguards.

References

- [1] United Nations, “IAEA Research and Development Plan: Enhancing Capabilities for Nuclear Verification,” Jan 2018.
- [2] United Nations, “IAEA Statute,” Oct 1956.

- [3] A. Goldman, R. Picard, and J. Shipley, “Statistical methods for nuclear materials safeguards: An overview,” *Technometrics*, vol. 24, no. 4, pp. 267–275, 1982.
- [4] C. Cortes and V. Vapnik, “Support-vector networks,” *Machine learning*, vol. 20, no. 3, pp. 273–297, 1995.
- [5] L. Breiman, “Random forests,” *Machine learning*, vol. 45, no. 1, pp. 5–32, 2001.
- [6] S. P. Lloyd, “Least squares quantization in pcm,” *IEEE Transactions on Information Theory*, vol. 28, pp. 129–137, 1982.
- [7] J. Schmidhuber, “Deep learning in neural networks: An overview,” *CoRR*, vol. abs/1404.7828, 2014.
- [8] B. B. Cipiti and N. Shoman, “Bulk handling facility modeling and simulation for safeguards analysis,” *Science and Technology of Nuclear Installations*, vol. 2018, 10 2018.
- [9] Y. LeCun, Y. Bengio, and G. Hinton, “Deep learning,” *Nature*, vol. 521, pp. 436–444, May 2015.
- [10] F. T. Liu, K. M. Ting, and Z.-H. Zhou, “Isolation forest,” *IEEE International Conference on Data Mining*, 2008.