# Predictive Analytics Solution

April 28, 2018

Shangwu Yao
shangwuy@andrew.cmu.edu

## 1   Introduction

Predictive analytics could be used to proactively anticipate gas-related safety hazards, which could reduce fatalities, injuries and other accidents. The most important considerations for predictive analytics in safety-related situation are accuracy and reliability of the analysis algorithms. In order to deal with that, I propose to use a method that adopts abnormality detection using hand-picked rules and machine learning-based hazard prediction algorithms, traditional engineering methods are used to detect abnormalities in sensor data and machine learning algorithms are used to learn from past experiences, the combination of them would be a reliable solution.

## 2   Solution

### 2.1   Machine learning-based failure detection

Machine learning-based failure detection methods have the advantage of learning from past experiences without human supervision, with enough training data, machine learning-based algorithms might achieve better results than human experts.
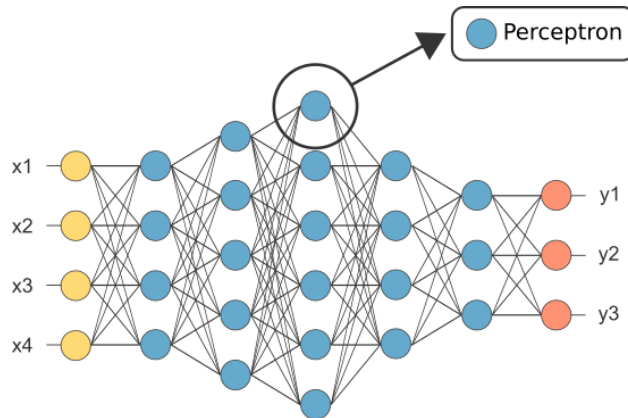


Figure 1: Multi-layer perceptron, image borrowed from https://www.neuraldesigner.com/blog/perceptron-the-main-component-of-neural-networks

For prediction of gas-related safety hazards, information about an upcoming hazard could be just in a very short time range that only spans one or two time steps in the collected data, or could in a relatively long period of time, so it is important to choose a model that could use temporal information. One could use incorporate such information by using a multi-layer perceptron model that takes a fixed window of information in several consecutive time steps as input, or a Long short-term memory model[1] that takes information of a single time step as input.
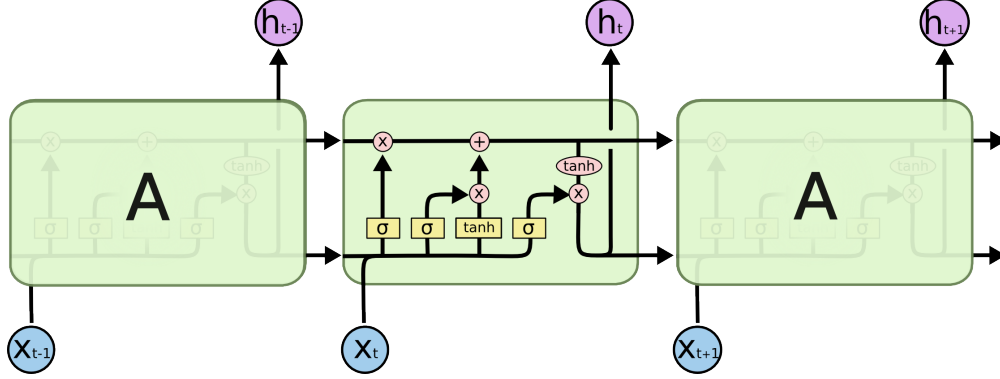
Figure 2: Long short-term memory model, image borrowed from http://colah.github.io/posts/2015-08-Understanding-LSTMs/

Machine learning algorithms are suitable to use in this failure detection situation, as the alarms that got triggered in the end could be used to label data. So, with input data being the average and peak gas level and temperature at different time steps, and output being whether the alarm will be triggered or not, a simple classifier could be used to solve this problem.

The problem is that training data for this problem will be highly imbalanced, because most of the training data will end up with the alarm not being triggered, then a classifier might just learn to predict the alarm not being triggered to everything, without even look at the training data at all. The solution to that is to use hard-negative mining, there are many possible ways to do this: one could use Support Vector Machine[2] and adopt the update rule mentioned in [3], or keep a fixed positive to negative ratio of 1 : 3 as Girshick[4] or selectively update examples with higher loss as Shrivastava et, al.[5].

A real-time ensemble of multiple model would also help to improve prediction accuracy[6] and thus improve the reliability of predictive analytics.

## 2.2 Abnormality detection

Abnormality detection is quite straightforward compared to machine learning-based hazard prediction, but requires some more hand-picked features by human. For example, in the alarm data, the average and peak amount of gas level as well as temperature at different time are provided, and a sharp rise or drop of those values might indicate a sudden change in the environment, which could lead to accidents. Also, a steady increase or decrease of gas level might be another sign of upcoming hazards. Creating these rules for detecting abnormalities might require human expert knowledge and a lot of hard-coding, but this is a reliable supplement to machine learning-based algorithms, as machine learning algorithms are not designed to deal with data that is very different from the training data.

# 3 Conclusion

Accuracy and reliability are the two most important considerations in predictive analytics for gas-related safety hazards. I presented a reliable method that adopts abnormality detection using hand-picked rules by human experts and machine learning-based hazard prediction algorithm that could learn from past experiences. Machine learning models such as Support Vector Machine, Multi-layer perceptron and Long short-term memory could be used to incorporate temporal information, and hard-negative mining could be used to deal with the imbalance of positive and negative data in this specific problem.

# References

[1] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.

[2] J.A.K. Suykens and J. Vandewalle. Least squares support vector machine classifiers. *Neural Processing Letters*, 9(3):293–300, Jun 1999.

[3] Pedro F Felzenszwalb, Ross B Girshick, David McAllester, and Deva Ramanan. Object detection with discriminatively trained part-based models. *IEEE transactions on pattern analysis and machine intelligence*, 32(9):1627–1645, 2010.

[4] Ross Girshick. Fast r-cnn. *arXiv preprint arXiv:1504.08083*, 2015.

[5] Abhinav Shrivastava, Abhinav Gupta, and Ross Girshick. Training region-based object detectors with online hard example mining. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 761–769, 2016.

[6] Thomas G. Dietterich. Ensemble methods in machine learning. In *Multiple Classifier Systems*, pages 1–15, Berlin, Heidelberg, 2000. Springer Berlin Heidelberg.