

---

# Self-Organizing Map for Industry Condition Monitoring

---

**Johnathan DiMatteo**

School of Computer Science

University of Waterloo

Waterloo, ON, N2L 3G1

`jdimatteo@uwaterloo.ca`

proposal due: October 23; report due: December 3.

## Abstract

1 The following is a proposal for a graduate research project for the University of  
2 Waterloo Computer Science class “CS 680: Introduction to Machine Learning”.  
3 The goal is to use machine learning to estimate the operating condition of an as-  
4 set (called the health score). In particular, the algorithm will be implemented and  
5 evaluated on an industrial furnace fan critical to a steel plant. This will allow op-  
6 erators in the steel mill to monitor the asset schedule efficient maintenance tasks,  
7 reducing unnecessary maintenance and unexpected failures. A proposed solution,  
8 using self-organizing maps, is used to provide asset owners with an indication of  
9 the health of a particular asset (known as *condition monitoring*).

## 10 Problem Statement

11 *To provide asset owners with a numerical value representing the asset’s operating condition or*  
12 *risk of failure.*

## 13 Motivation

14 An often overlooked part of an asset’s expenses is maintenance. A popular maintenance strategy is  
15 known as Preventative Maintenance (PvM), where maintenance is regularly performed on an asset  
16 while it is still in good condition to prevent it from breaking down unexpectedly. Over \$ 200 billion  
17 is spent on such maintenance every year in the United States and one-third is wasted on improper  
18 or unnecessary maintenance [1]. Even worse, no maintenance at all can lead to unexpected failures  
19 which in turn can cause serious economic consequences or injury. There exists a significant need to  
20 modernize maintenance techniques around the world to ensure safety, reliability, and efficiency.

21 During the Second World War, a British scientist named Conrad Waddington made a fascinating dis-  
22 covery about the maintenance of aircraft while working for the Royal Air Force (RAF). Previously,  
23 aircraft bombers had a notorious problem of breaking down - in fact the ideal serviceability in a  
24 squadron of bombers was only around 70-75% [2]. What he discovered was that preventative main-  
25 tenance methods actually increased the rate of unexpected failure. The process of more maintenance  
26 leading to more failures became known as the Waddington Effect as a result. By increasing the inter-  
27 val between maintenance cycles and eliminating all maintenance deemed unnecessary, Waddington  
28 was able to increase the effective flight hours of the RAF bomber fleet by 60% [2].

29 After this discovery, asset owners around the world tried to find the optimal time to repair an asset.  
30 This led to the invention of Predictive Maintenance (PdM), a philosophy that uses the actual oper-  
31 ating condition of assets to optimize operations [1]. For PdM to be effective, the asset’s operating

condition must be estimated. Estimating the asset's operating condition is the focus of the paper, which is referred to as the *health score*.

## 1 Background

### 1.1 Why Machine Learning?

One of the most common reasons PdM methods fail is a lack of continuous improvement and a lack of repeatability [3]. Additionally, equipment monitoring is a time consuming process, requires experts to identify failure patterns, and is expensive. Machine learning provides an automated approach that requires minimal asset knowledge, is inexpensive, can be trained on many assets, and can be re-trained as operating conditions change.

Machine learning techniques for predictive maintenance were considered not practical, too complex, or too time consuming. In particular, plant managers did not want to change their existing infrastructure (the software that handles data acquisition and analyzes it) to adopt the technology. But now Asset Performance Management (APM) software providers are growing and condition based maintenance is at the forefront. As they team up with cloud based data solutions, it becomes easy for asset operators to implement machine learning in their existing data infrastructures via a simple call to the cloud. Manufacturers around the world use APM technology from Bentley Systems, a global leader in APM capabilities according to a recent Gartner report [4]. The proposed solution will be deployed and maintained using APM software from Bentley Systems.

### 1.2 ArcelorMittal Dofasco

**ASK DAVID IF THIS IS OKAY** ArcelorMittal Dofasco is a steel company located in Hamilton, Ontario. They use Bentley's APM software and are eager for a machine learning solution to detect the operating conditions of various assets. In particular, they have offered a real data set of several industrial level furnace fans located in the Hamilton plant. Specifically, the data is composed of several smaller data sets, each representing various hours of operation. See the Appendix for full list of variables included.

In a steel making plant, called a steel mill, operations run almost 24/7 except when the mill is shut down once a month for repairs and maintenance. A failure of an asset leading to a shutdown at any other time results in severe costs. If the operating condition of the asset is known, then operators can determine whether or not it should not be repaired during that scheduled downtime, thus avoiding costly unexpected failures and the Waddington Effect.

A reheating furnace is used to raise the internal temperature of steel, so that it can be shaped into a final product. Setting the correct temperature is one of the most essential factors of product quality in the plant. If at any point the blast furnaces fail, the entire line must be shut down to allow operators to repair the furnaces. Operators suggest this could take X amount of time, resulting in approximately \$ of lost production

**GET A QUOTE?.**

## 2 Previous Work

Using machine learning for condition monitoring is not new. There are typically three approaches:

1. Supervised algorithms using sensor data and maintenance data.
2. Unsupervised algorithms using only sensor data.
3. Semi-supervised algorithms using only *healthy* sensor data.

Supervised methods require labeled data. Labeling is the process of associating an output with a set of variable values at a certain point in time. For example, if there is knowledge to when the failures occurred and there are sufficient failures, the data can be labeled as healthy or not. Unfortunately, there is just not enough failure data in practice to make this feasible. Instead, failure data can be estimated and artificially generated [5]. But even then the algorithms are limited to a binary outputs.

78 Using  $k$ -nearest neighbours (kNN) to estimate asset condition directly is difficult due to noise, and  
 79 requires domain specific knowledge to choose appropriate variables. One paper improves on kNN  
 80 methods for detecting the levels of severity for cracks in gears [6]. The disadvantage in these ap-  
 81 proaches is that it requires significant data in a variety of conditions, and it uses classification to  
 82 identify a severity level instead of a numerical value.

83 Researchers in [7] train an autoencoder on healthy imaging data to get a feature representation in a  
 84 smaller number of dimensions. Utilizing the Support Vector Machine (SVM), they learn a decision  
 85 boundary to identify anomalies. This approach is interesting for high dimensional data but focuses  
 86 on point anomalies, where as we are interested in projecting asset health over time.

87 A Self-Organizing Map (SOM) is a type of neural network that is often used as a dimensionality  
 88 reduction technique as it produces a low dimensional representation of the training samples [8]. The  
 89 two most relevant papers on using SOMs for condition maintenance were used on aircraft engines  
 90 and ball bearings [9] [10]. The SOM can be visualized directly as a tool to identify failures as in  
 Figure 1, where datapoints 517-607 are clearly anomalous. Alternatively, Huang et al. were able to

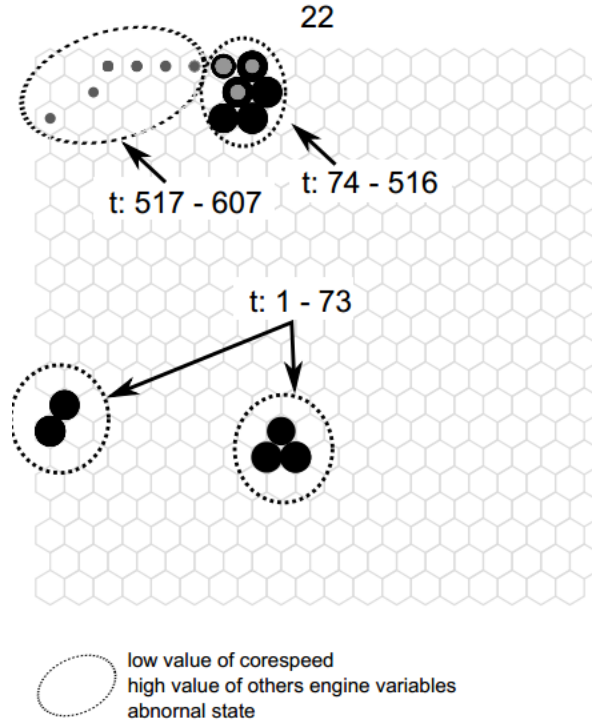


Figure 1: SOM map being used to identify anomalous behaviour [9].

91 use the minimum quantization error between a test point and the closest neuron of the SOM as the  
 92 health indicator for ball bearings [11].  
 93

$$Q = \min_k \|D - B_k\| \quad (1)$$

94 Where  $Q$  is the minimum quantization error,  $D$  is a test set observation, and  $B_k$  is the weight vector  
 95 of the  $k^{th}$  closest neuron of the SOM. But, using quantization error can be improved as it is sensitive  
 96 to noise. Researchers Tian et al. use a SOM and a  $k$ -nearest neighbours algorithm (see Figure 2)  
 97 in combination with the Euclidean distance between the test data and the healthy data to develop a  
 98 more robust health score.

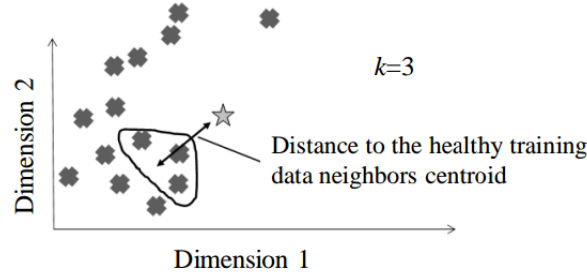


Figure 2: A health score can be used as the distance from a test point to a cluster (or a neuron in a SOM) [10].

### 99 3 Proposed Work

#### 100 3.1 Implementation

101 The project will solve PdM hurdles for operators in modernized workplaces by providing an ac-  
 102 curate estimate of asset health. The proposed method is to replicate the results achieved by other  
 103 papers using SOMs for condition monitoring. Specifically, to train a SOM and get a health score by  
 104 calculating how much a test data point deviates from normal operation. The higher the health score,  
 105 the lower the operating condition of the asset. Various SOM implementations will be examined,  
 106 such as the [Kohonen 1.1.2 Python package](#) (documentation for this package is sparse). If these do  
 107 not provide enough customization then a SOM will be self-coded using Tensorflow. Customization  
 108 may be required to allow the comparison of various distance metrics to get a health score. As dis-  
 109 cussed, this has been done before but is rarely seen in industry today. The use of these algorithms  
 110 on real and significant data is challenging and relevant. For example, the data will be provided in  
 111 hourly chunks throughout the year (instead of one continuous dataset). Futhermore, pre and post  
 112 processing techniques will be explored to improve performance, useability and interpretability. De-  
 113 cision boundaries will be learned as a final step (if we have enough failure data) to organize failures  
 114 into different classes and severity levels.

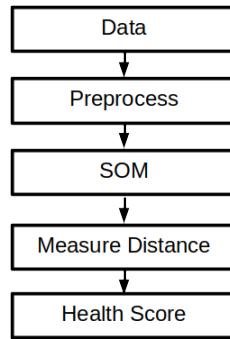


Figure 3: Implementation steps from Data.

#### 115 3.2 Evaluation

116 An expected significant challenge will be the model’s ability to generalize, given that the datasets  
 117 span only a days worth of operation in total. To avoid overfitting, the number of dimensions will be  
 118 kept low and regularization will be used where possible. To judge the models ability to generalize  
 119 and to evaluate performance, the algorithm will be tested on two datasets: unhealthy operation (1)  
 120 and healthy operation (2). Precision was chosen as a metric for (2) to minimize the number of  
 121 false positives. This is especially important because this is an experimental project with the steel  
 122 manufacturer, so it is important that the algorithm does not do more harm than good. For (1), the  
 123 algorithm should produce a value indicative of failure at least 75% of the time.

124 Performance evaluation of the algorithm from a PdM context will be performed against each of these  
125 three scenarios.

- 126 1. Run-to-Failure: how will the proposed solution compare if no maintenance is performed?
- 127 2. Preventative: how will the proposed solution compare if PvM is performed (the current  
128 process of the plant)?
- 129 3. Other Predictive Methods: how will the proposed solution compare if another PdM solution  
130 is performed?

131 Other metrics for evaluation will include the following: cost, latency (time between failure and  
132 an indication of failure from the algorithm) and training time. Finally, the algorithm will also be  
133 scrutinized by a subject matter expert (SME) in terms of usability and interpretability.

### 134 **3.3 Timeline**

135 The time from the proposal submission date and the final due date is six weeks.

- 136 1. Week 1: Choose a SOM algorithm or develop one.
- 137 2. Week 2: Choose a SOM algorithm or develop one.
- 138 3. Week 3: Test algorithm on data, visualize.
- 139 4. Week 4: Choose and evaluate various distance metrics.
- 140 5. Week 5: Tune model, test, repeat. Try preprocessing methods.
- 141 6. Week 6: Report writing and final touchups.

## References

- [1] R.K. Mobley. *An Introduction to Predictive Maintenance*. Plant Engineering. Elsevier Science, 2002. ISBN 9780080478692. URL <https://books.google.ca/books?id=SjqXzxpAzSQC>.
- [2] Philip M. Morse. OR in World War 2. Operational Research Against the U Boat. *Science*, 184(4144):1364–1365, 1974. ISSN 0036-8075. doi: 10.1126/science.184.4144.1364. URL <http://science.sciencemag.org/content/184/4144/1364>.
- [3] Douglas Hart. Predictive Maintenance - Avoiding the Ten Most Common Pitfalls. Technical report, Emerson, 02 2017.
- [4] Nicole Foust and Kristian Steenstrup. *Market Guide for Asset Performance Management Software*. Jun 2018.
- [5] Vasilis A. Sotiris, Peter W. Tse, and Michael G. Pecht. Anomaly detection through a bayesian support vector machine. *IEEE Transactions on Reliability*, 59:277–286, 2010.
- [6] Yaguo Lei and Ming J Zuo. Gear crack level identification based on weighted k nearest neighbor classification algorithm. *Mechanical Systems and Signal Processing*, 23(5):1535–1547, 2009.
- [7] Philipp Seebck, Sebastian Waldstein, Sophie Klimscha, Bianca Gerendas, Ren Donner, Thomas Schlegl, Ursula Schmidt-Erfurth, and Georg Langs. Identifying and categorizing anomalies in retinal imaging data. 12 2016.
- [8] Teuvo Kohonen. Self-organizing map. *Proceedings of the IEEE*, 78:1464–1480, 1990.
- [9] Etienne Côme, Marie Cottrell, Michel Verleysen, and Jérôme Lacaille. Aircraft engine health monitoring using self-organizing maps. In *Industrial Conference on Data Mining*, pages 405–417. Springer, 2010.
- [10] Jianrong Tian, Michael H. Azarian, and Michael Pecht. Anomaly detection using self-organizing maps-based k-nearest neighbor algorithm. 2014.
- [11] Runqing Huang, Lifeng Xi, Xinglin Li, C Richard Liu, Hai Qiu, and Jay Lee. Residual life predictions for ball bearings based on self-organizing map and back propagation neural network methods. 21:193–207, 01 2007.
- [12] O. Geramifard, J. . Xu, C. K. Pang, J. H. Zhou, and X. Li. Data-driven approaches in health condition monitoring a comparative study. In *IEEE ICCA 2010*, pages 1618–1622, June 2010. doi: 10.1109/ICCA.2010.5524339.

<b>Variable Name</b>	<b>Description</b>
F1S F1SFIBV Overall (g RMS pk)	Furnace 1 south fan inboard bearing vibration
F1S F1SFOBV Overall (g RMS pk)	Furnace 1 south fan outboard bearing vibration
F1S F1SSMIBV Overall (g RMS pk)	Furnace 1 south fan motor inboard bearing vibration
F1S F1SMOBV Overall (g RMS pk)	Furnace 1 south motor outboard bearing vibration
F1S F1NFIBV Overall (g RMS pk)	Furnace 1 north fan inboard bearing vibration
F1S F1NFOBV Overall (g RMS pk)	Furnace 1 north fan outboard bearing vibration
F1S F1NMIBV Overall (g RMS pk)	Furnace 1 north fan motor inboard bearing vibration
F1S F1NMOBV Overall (g RMS pk)	Furnace 1 north fan motor outboard bearing vibration
F1S North Fan Impellor side bearing temp	Furnace 1 north fan outboard bearing temp
F1S North Fan Motor side bearing temp	Furnace 1 north fan inboard bearing temp
F1S South Fan Impellor side bearing temp	Furnace 1 south fan outboard bearing temp
F1S South Fan Motor side bearing temp	Furnace 1 south fan inboard bearing temp

Table 1: Variable names and descriptions provided by subject matter expert at ArclorMittal Dofasco plant.