



Credit: iStock: yoh4nn

# A machine learning approach to quality control

Most of the products we take for granted, from laptops to jet engines, contain huge numbers of components assembled in multiple stages by different manufacturers. Quality control is vital throughout the assembly process with rigorous testing required at every step to ensure that the next customer in the chain – up to and including the end-user – does not receive faulty or sub-standard goods.

With large numbers of components and assembly stages, testing every part 'by hand' is time-consuming and expensive. We wanted to know if we could use machine learning to reliably predict failure rates and, by doing so, reduce the costs of testing, improve quality control and reduce delivery times.

To do this we needed a real industrial problem to work on. Step forward a global electronics company specialising in the manufacture of frequency inverter drives, used to control and regulate the speed of electric motors. Each product combines software with two or three electronic printed circuit board assemblies

(PCBAs). Every PCBA has four quality control stages, each of which includes up to 5,000 control tests or steps.

If a faulty product is sent to a customer, a number of costs are incurred. A replacement product has to be shipped, customer service time and resources are needed to support the complaint and



**Read more** about the machine learning approach used by the authors, in their paper: 'Cost-sensitive learning classification strategy for predicting product failures' in Expert Systems with Applications Vol. 161 2020 doi.org/10.1016/j.eswa.2020.113653

there is the potential for reputational damage.

Our task was to find a flexible classification method that could take into account the cost of supplying faulty goods and trade it off against the quality (or number of faulty goods) sent to the customer. Our challenge is that the data is voluminous, complex and often incomplete. Tests are sometimes – and entirely legitimately – skipped where they do not affect the outcome but this does create holes in the dataset.

Due to the highly optimised industrial process, the number of faulty products is very low. Thus, the data is 'imbalanced' where 99% of the products are classified as good and only 1% are faulty. For classification problems, this is a data characteristic that algorithms tend to struggle with.

While there are possibly hundreds of algorithms capable of dealing with datasets such as these, we needed one that could also introduce a cost dimension to the problem so that manufacturers could see if this approach would save them money.

Our solution uses a cost-sensitive classification strategy which we modified in order to address this specific industrial problem. It assumes that there is an interdependence between each stage of the manufacturing process and, therefore, we can predict the final stage by analysing data collected during previous production stages.

Our two-step process starts with the application of 'feature engineering' methods to prepare the data, converting it from multidimensional to two dimensions and using some mathematical tricks to allow for missing data. We were then able to classify the data – predicting whether a

product will be good or faulty - using our algorithm.

When classifying new products we ended up with four possible outcomes:

- ▶ A product that is faulty but classified as good
- ▶ A product that is faulty and is classified as faulty
- ▶ A product that is good but classified as faulty
- ▶ A product that is good and classified as good.

Our interest lies in the undetected faulty products that get shipped to the customer. If a faulty product is correctly classified as faulty it will be spotted and repaired immediately, at minimal cost. Similarly, if a good product is misclassified as faulty, it will be checked and found to be good straightaway – again with minimal cost.

However, a previous analysis had found that the cost of a false negative (in other words, sending a faulty product to a customer) is roughly 20 times higher than the cost of a false positive (a good product misclassified as faulty). We used this number ( $C=20$ ) as our cost parameter but we also compared it with a higher number ( $C=100$ ) to see how that would affect the outcome.

The results were interesting. When  $C = 20$ , 98% 'sensitivity' (where 100% is no faulty goods are sent to the customer) can be achieved at a cost reduction of 25% against checking of individual products. Even when  $C$  is 100, a similar set of results was achieved. With sensitivity still at 98%, the cost savings are 23% instead of 25%.

We then checked our results by using the same algorithm against 25 real-world datasets with different levels of

imbalance. The results confirmed that the approach is robust and is also flexible enough to be used when the cost is not specified.

Our research suggests that taking a machine learning approach to quality control can be beneficial for manufacturers when the numbers of faulty products are low. In those circumstances, the cost-savings of using an algorithm can outweigh the risks and costs associated with supplying sub-standard products.

---

## Authors

Flavia Dalia Frumosu, Technical University of Denmark  
Abdul Rauf Khan, Technical University of Denmark  
Henrik Schioler, Aalborg University  
Murat Kulahci, Technical University of Denmark  
Mohamed Zaki, University of Cambridge  
Peter Westermann-Rasmussen, Danfoss Drives A/S