

Machine Learning

Defect Prediction on Production Line

Team:

AIT BACHIR Romuald

ALLEMAND Fabien

FLORET Arthur

JARDOT Charles



Abstract

The primary goal in detecting anomalies using machine learning in the context of a industry is to develop a model able to detect whether an item is defective or not using data related to said item. This model would be a major improvement for the industry since it takes time to manually test every single part.

As expected from production lines' data, defective and non-defective pieces are not equally represented. Since the imbalance between classes harm the model's learning process, balancing them is one way to deal with this issue.

The Naïve Bayes classifier was chosen to answer classify items on the production line because it was evaluated as the most performant model at this task among various other machine learning models. The results demonstrate that such a simple algorithm can already predict defects, which is encouraging for the future when more data is collected and heavier models are implemented.

Contents

1	Introduction	1
2	Data Analysis	1
2.1	Data Description	1
2.2	In-Depth Analysis	1
3	Data Preparation	4
3.1	Missing Values	4
3.2	Balance Classes	5
4	Model Selection	6
4.1	Evaluation Method	7
4.2	Naive Bayes Classifier	8
4.3	k-Nearest Neighbors	8
4.4	Random Forest Classifier	9
4.5	Multilayer Perceptron	9
4.6	Novelty Detection	11
5	Model Fine Tuning	12
6	Conclusion	14

List of Figures

1	Distributions of the input features	2
2	Features pairwise correlation matrix	2
3	PCA results (projection of the training population on the three most relevant dimensions) . .	3
4	PCA results (projection of the training population on the two most relevant dimensions) . .	3
5	Distribution of the <i>angle_2</i> feature with items grouped by class	4
6	Box-plot of the <i>snap_ring_final_stroke</i> feature	4
7	Distribution of the <i>capuchon_insertion</i> feature with items grouped by class	5
8	Confusion matrix of a k-Nearest Neighbors classifier with the base dataset	6
9	Confusion matrix of a k-Nearest Neighbors classifier with the dataset balanced by removing individuals	6
10	Confusion matrix of a k-Nearest Neighbors classifier with the dataset balanced by duplicating individuals	7
11	Naive Bayes Classifier Evaluation	8
12	k-Nearest Neighbors Evaluation	9
13	Random Forest Classifier Evaluation (with <i>class_weight</i> set to "balanced")	10
14	Random Forest Classifier Evaluation	10
15	Best parameters for the MLP	11
16	MLP Classifier Evaluation	12
17	One Class SVM Classifier Evaluation	13
18	Naive Bayes Classifier Evaluation with tuned parameters	13

1 Introduction

The goal of this project is to develop an AI based solution for a real world industry problem faced by Valeo [1].

Valeo [2] is an industry leading French automotive supplier. In order to stay competitive, the company wants to develop a system that is able to identify defects on products before testing.

Four files containing relevant data are available. The data are mainly values measured during production on different mounting stations as well as additional measures performed on test benches.

The target is to find the best prediction: $\text{Output} = f(\text{inputs})$

2 Data Analysis

The goal of this section is to take a first look at the data in order to understand what it represents, then proceed to an in-depth analysis of said data.

2.1 Data Description

Data is contained in four different comma-separated values (csv) files: inputs and output for both training and testing. The very first line of each file corresponds to the headers. The following lines contain informations about a given product where each product is identified by a unique code:

ID = PROC_TRACEINFO = It's a unique code given to the product.

Example: I-B-XA1207672-190701-00494.

- XA1207672 is the reference.
- 190701 is the date: here 01st of July of year 2019.
- 00494 is the unique code given to the product, whatever it happens, the product will have this id number frozen forever.

This number is increased by 1 each time we process a new product, every 12 seconds. So for example: I-B-XA1207672-190701-00495 is the next product.

Input features are measures (numerical values) collected on different assembly stations with the sensors or devices connected to Programmable Logic Controllers which are storing all of them to keep the full quality traceability. Examples : OP070_V_1_angle_value, OP120_Rodage_I_value...

This is the result value of OP130 (test bench). Value 0 is assigned to OK samples (passed) and value 1 is assigned to KO samples (failed). This is the combined result of multiple electrical, acoustic and vibro-acoustic tests.

2.2 In-Depth Analysis

Let us put the testing data aside for now and take a closer look at the training data.

The training data consists of 13 measures (inputs features) made during production line on just shy of 35000 products. These measures correspond to various forces, angles, screwing torques... By taking a look at the distribution for each feature (Figure 1) we can clearly see that every one of them has a different shape.

A quick glance at the pairwise correlation matrix reveals that there are no obvious correlations between any features (Figure 2).

Looking at the results of a PCA (Figure 3 4) or at multiple histograms (with items grouped by class for each feature, as displayed in Figure 5), there is no feature that clearly helps separating the two classes.

An in-depth look at the box-plots (Figure 6) reveals that there might be some outliers in the dataset.

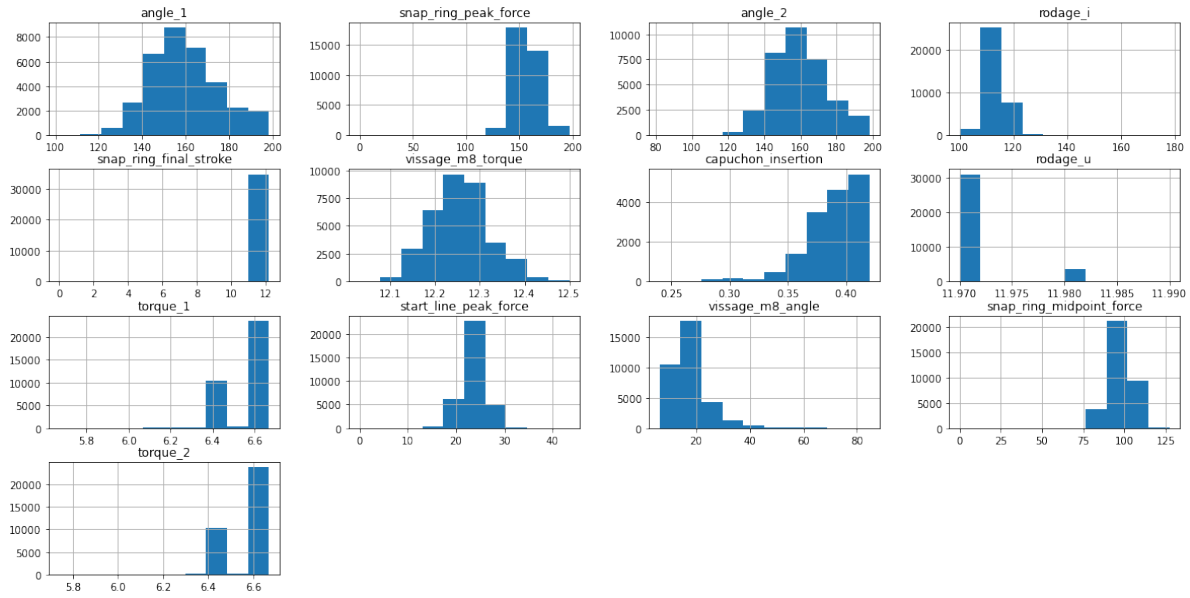


Figure 1: Distributions of the input features

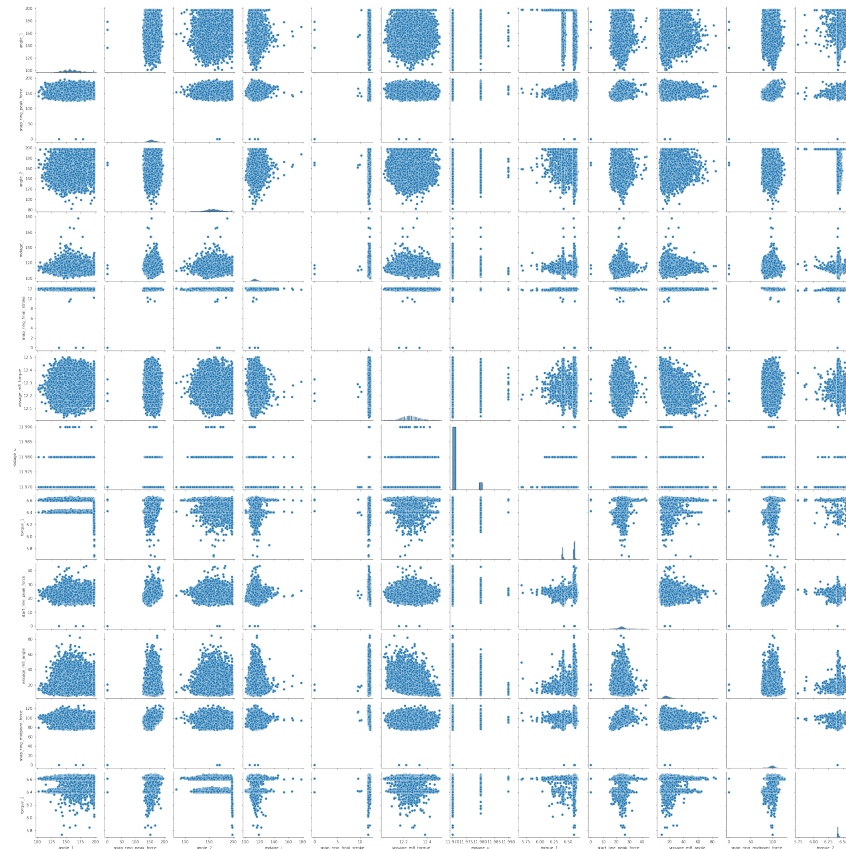


Figure 2: Features pairwise correlation matrix

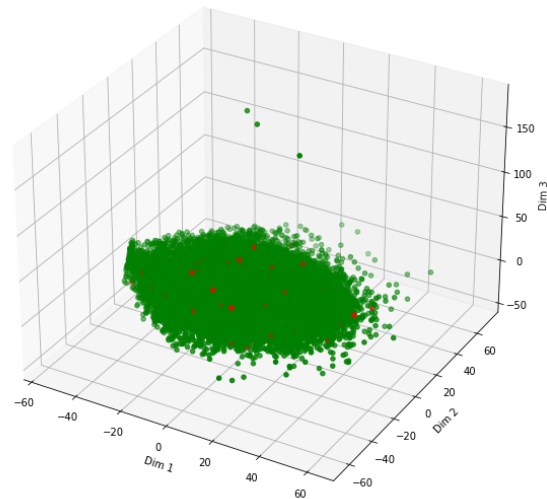


Figure 3: PCA results (projection of the training population on the three most relevant dimensions)

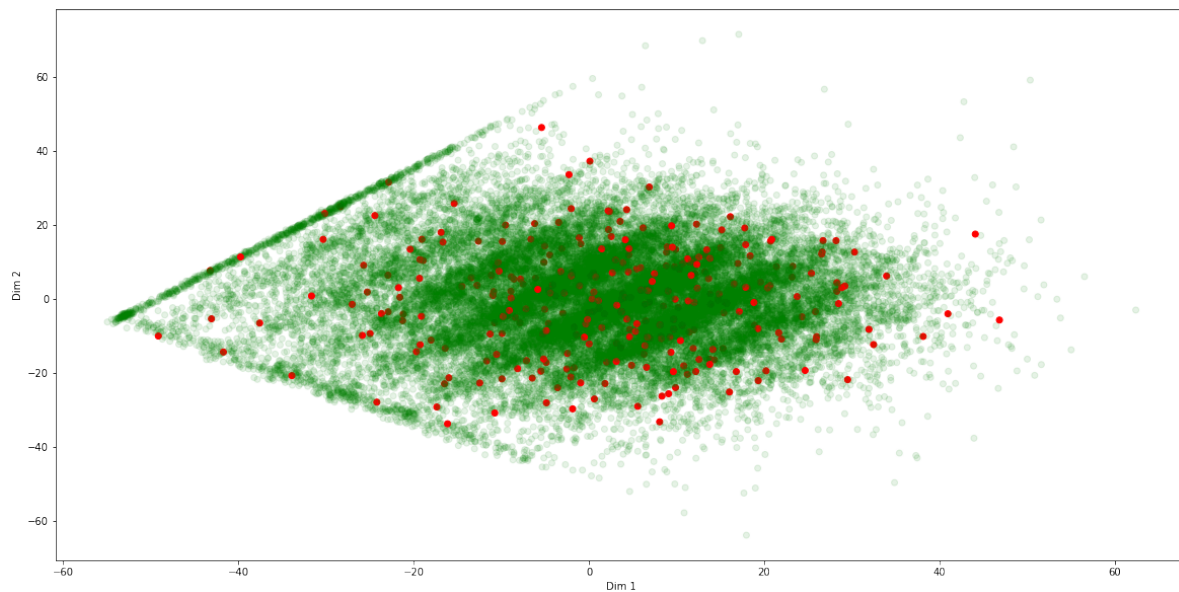


Figure 4: PCA results (projection of the training population on the two most relevant dimensions)

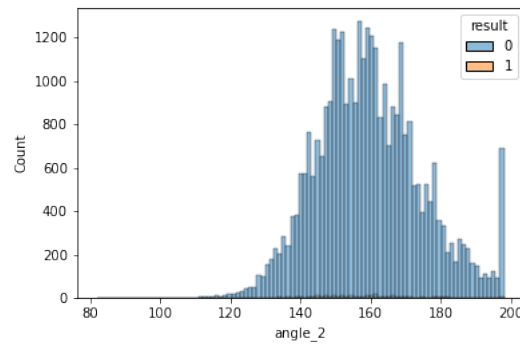


Figure 5: Distribution of the *angle_2* feature with items grouped by class

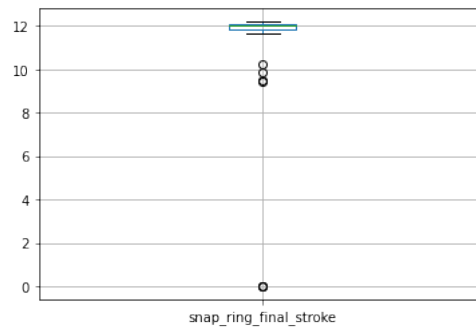


Figure 6: Box-plot of the *snap_ring_final_stroke* feature

One feature requires more attention. *capuchon_insertion* is the only feature containing missing values. In fact, more than 50% of the sample population are missing this value.

The training data also contains the result value of the test bench for each product tested on the production line. As expected, there are only a few of defective items. Although this is good for the company, it means that **less than 1%** of the training population are defective items. The effectives of the two classes are very unbalanced which could really impact our models.

3 Data Preparation

The data analysis performed in the previous section revealed some flaws in the dataset. The two most noticeable are the presence of missing values and the population imbalance between the two classes. These two flaws are known for harming the learning of most ML models. Different attempts at solving these issues are presented in this section.

3.1 Missing Values

The previous section mentions missing values in the dataset. This notion is self-explanatory: some individuals have no value entered for some feature.

We noticed that these missing values only concerned the *capuchon_insertion* feature but their amount was not negligible. Among 34515 items, 18627 of them (i.e. 53% of the population) do not have any value describing this dimension.

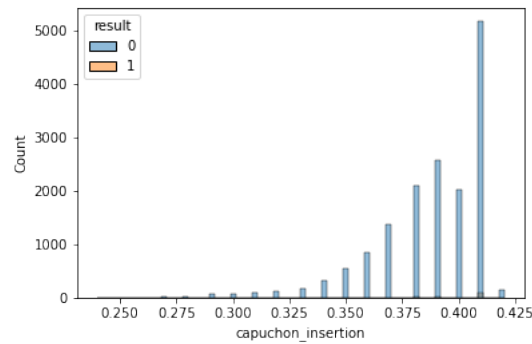


Figure 7: Distribution of the *capuchon_insertion* feature with items grouped by class

The first idea to prevent this issue would be to fill the missing values with an estimation like the mean or the median of the values that are entered for this feature.

However this method has some drawbacks. First, it makes the assumption the data could have been measured but on the production line the *capuchon_insertion* was probably not measured because the item was already broken. Then it makes an assumption on the value of the parameter which can potentially compromise the data. What if the defective and valid items have a different mean value for this feature? Of course we could also fill the data with an estimation of the value for each class but knowing the classes are unbalanced it might be hazardous. Further analysis show 110 defective items among the 305 present in the training data are missing this value. While 36% of missing values does not look too bad, estimating the mean or the median with only 195 individuals starts to get tricky.

Another way to deal with this issue is to simply drop this feature from the dataset. Figure 7 clearly shows there is no distinction between defective and valid items regarding *capuchon_insertion*. Hence it is not a valuable feature when it comes to classify the parts on the production line.

3.2 Balance Classes

The following paragraphs deal with class balancing, explaining the importance of balanced classes and three different approaches to rectify the imbalance.

As said in the previous section, a mere 0.008% of the population is defective individuals. In other words, the class of individuals which needs to be detected is truly under-represented in the dataset. Without any pre-processing such a dataset is very unlikely to give good results with most of the ML models. As shown in Figure 8, there are not enough defective individuals for the classifier to properly learn how to distinguish a defective individual from a valid individual. As a result, the classifier simply over-fit the valid class.

Remarque Over-fitting the over-represented class is very hazardous as the accuracy of the classifier will remain really high even though the classifier is not classifying anything. Hopefully the confusion matrix shows the classification details.

The easiest method to deal with unbalanced classes is to remove valid individuals from the training dataset. This process is fairly straight forward (valid individuals can be randomly selected in the dataset) and does not require a lot of processing (selected individuals are simply deleted).

Figure 9 clearly shows that the classifier is not over-fitting the valid class anymore. In that sense, the results are better but they are still far from what is expected from a performant classifier.

This method is very arguable for two reasons. First, deleting individuals obviously implies a loss of data. Even though the individuals are randomly selected, some important features can disappear. In other words, the remaining data is not necessarily representative of the global population. On top of that, the resulting dataset

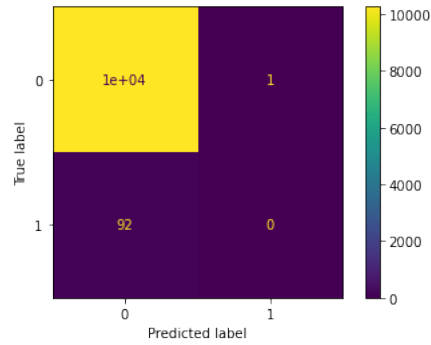


Figure 8: Confusion matrix of a k-Nearest Neighbors classifier with the base dataset

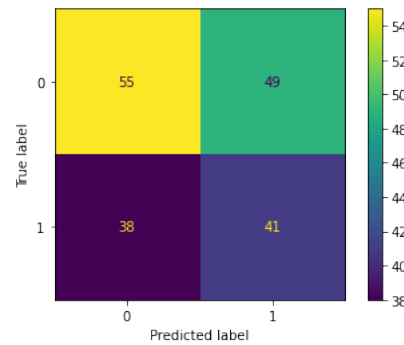


Figure 9: Confusion matrix of a k-Nearest Neighbors classifier with the dataset balanced by removing individuals

is smaller meaning the model will have less instances to learn on. In this problem, this will dramatically impact the results because there is already few data.

Of course it is also possible to work the other way around and artificially create individuals instead of deleting others. The easiest method is actually to duplicate defective items until there are as many as the valid ones. This compensate the imbalance between the two classes without knowing anything about the features values. Figure 10 shows the effect on the classifier is significant.

Another way that has been implemented during this project is to slightly modify individuals during the duplication. This way new individuals are actually new individuals (not only duplicates) adding information to the overall dataset. This method assumes individuals from the same class are close to one another on every dimension (i.e. regarding every feature) and changing a feature's value by a small amount will not induce a change of class. The way it has been implemented is surely not ideal but rather intuitive. A random amount of features are selected during each "duplication" and the value is modified by a very small percentage (typically 0.005).

The final method in order to balance classes consists in leaving the classes imbalanced in the dataset but during computation attributing a more significant weight to the individuals from the under-represented class. Technically speaking this is equivalent to duplicating individuals but it requires less memory, it is easier for computation and is already implemented for some algorithms from well known ML libraries.

4 Model Selection

Choosing the ML model that works the best for a given problem is a key factor for good results but it can be very time consuming. The following pages present the results obtained with several ML algorithms found in

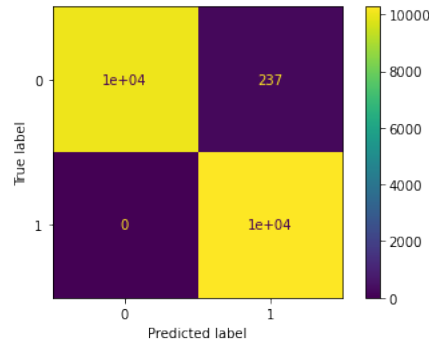


Figure 10: Confusion matrix of a k-Nearest Neighbors classifier with the dataset balanced by duplicating individuals

the Scikit-Learn library [0]. Each method will be briefly introduced and tested using the training dataset. The resulting classifier are evaluated using three different tools:

- Confusion matrix
- Precision, recall and F1 score
- ROC AUC Score

These three tools are presented and explained in Section 4.1.

Knowing the data is labeled all the methods but the last one are based on supervised learning.

Moreover it has already been demonstrated in Section 3.2 that balancing classes dramatically improves results so the ML models are tested with an augmented training dataset containing duplicated individuals.

4.1 Evaluation Method

There are many tools to analyse the performances of ML models. Three tools will be thoroughly used later. The goal of the following paragraphs is to give a quick understanding at what they represent and how to use them properly.

The confusion matrix is a tool specifically used for classification problems. This matrix simply represents the results of the classification made by a given classifier. For instance in Figure 10, the classifier classified 237 valid items as defective (top right).

By looking at a confusion matrix, it is fairly easy to see what the classifier did wrong by comparing it to the ideal matrix which would contains zeros for the false positive (top right) and false negative (bottom left) emplacements.

The values contained in the confusion matrix also allow to compute the precision, recall and f1 score of a classifier:

- Precision: how accurate is the classifier when classifying a defective item (closer to one is better)
- Recall: how good is the classifier at finding defective items in the population (closer to one is better)
- f1 score: the harmonic mean of precision and recall (closer to one is better)

Precision and recall are linked values meaning there is always a tradeoff: if precision rises, recall tends to decrease. So it is important to adjust precision and recall according to the application. In this problem, a company would most likely prefer to double check an item rather than selling a defective one so a strong recall would be ideal, at the expense of precision.

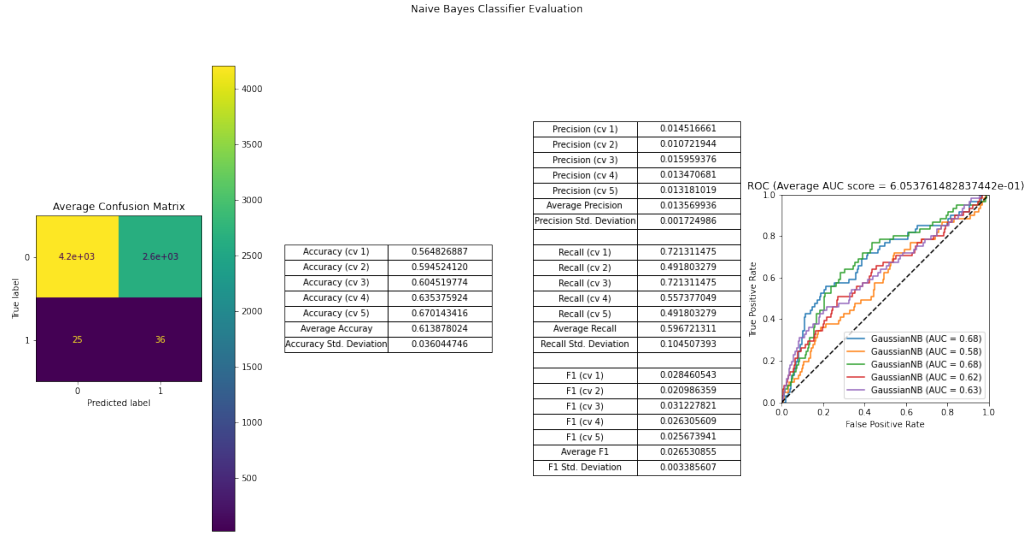


Figure 11: Naive Bayes Classifier Evaluation

Finally, Receiver Operating Characteristic (ROC) curve plots the true positive rate against the false positive rate (i.e. sensitivity against specificity, Figure 18). It allows to distinguish the classifier from a random classifier. It is generally associated to the ROC Area Under the Curve (AUC) score. A ROC AUC score over 0.5 means the classifier is better than a random classifier. An ideal classifier would show a ROC curve following the top left edge of the graph and a ROC AUC score of one.

4.2 Naive Bayes Classifier

Naive Bayes classifiers use Bayes' theorem making the "naive" assumption that every pair of features are independant given the value of the class variable [4] (i.e. $P(x_i|y, x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n) = P(x_i|y)$) leading to the following equality:

$$P(y|x_1, \dots, x_n) = \frac{P(y)P(x_1, \dots, x_n|y)}{P(x_1, \dots, x_n)} = \frac{P(y) \prod_{i=1}^n P(x_i|y)}{P(x_1, \dots, x_n)}$$

The following results correspond to a Gaussian Naive Bayes [0] classifier that is very common when dealing with numerical continuous values. However it also assumes that every features follows a Gaussian distribution.

The confusion matrix in Figure 18 shows poor results. This is expected because Figure 1 demonstrates some features do not follow a Gaussian distribution. In addition, the assumption that features are independant from one another might not be relevant with this dataset. For instance, the *snap_ring_peak_force* might increase with the value of *anlge_1*.

4.3 k-Nearest Neighbors

The k-nearest neighbors algorithm works by classifying an individual by looking at its neighbors [0]. The classifier simply look at the k nearest neighbors of the individual in the training population and attributes the class that correspond the best to the neighbors (the neighbors vote for their class, each vote is weighted either evenly accros the neighbors or according to the distance between the individual and the voting neighbor [0]).

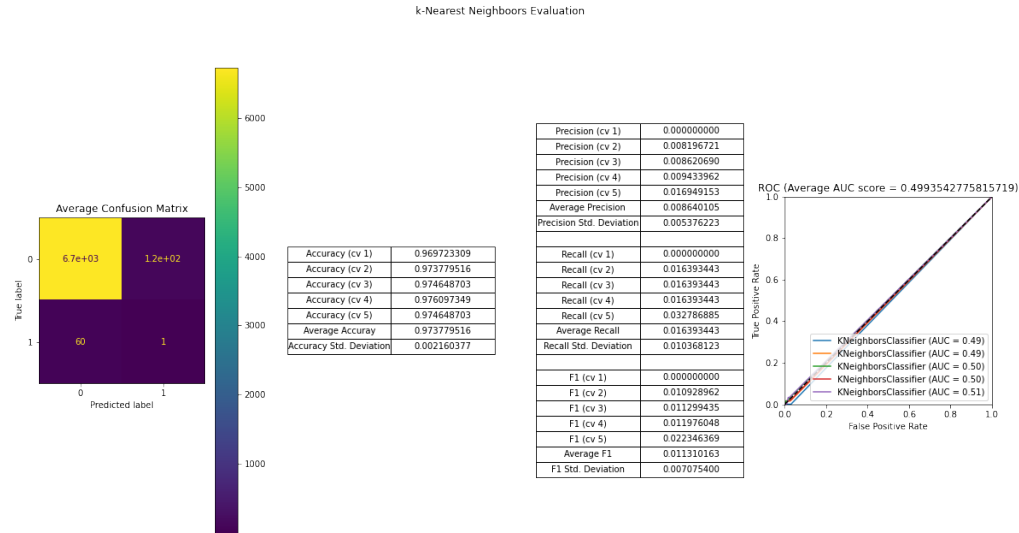


Figure 12: k-Nearest Neighbors Evaluation

According to Figure 12, this algorithm creates a model that is close to a random classifier. Using cross-validation, the average precision is lower than 1% and the average recall is lower than 2% (with very little deviation in each case). This means that, no matter the training set, the classifier is no able to distinguish a defective item from a correct one.

4.4 Random Forest Classifier

Random Forest Classifier works by creating multiple Decision Trees on various subsets of the data, the decision of the forest for classification is simply the class that has been predicted by most of the trees [0]. A Decision Tree is built by finding the rule which divides the best the population of each node, this can be the rule that isolate most individuals from the same class (this method is called information gain and is based on the concept of entropy)[0]. A node cannot be divided when there are not enough individuals (in the node or in the to-be-created child-node), then the leaf is labeled with the class with the most representants on the node.

Random Forest Classifiers have a *class_weight* attribute in Scikit-Learn in order to address imbalance between classes[6]. However it seemed to perform worse than training the forest with the dataset containing duplicates (Figure 13). The ROC AUC score on Figure 14 shows just-above-random results. But looking at the precision and recall shed the light on a useless model. Not being able to detect any defective item makes it a very bad candidate.

4.5 Multilayer Perceptron

Multilayer Perceptron is part of the non-linear supervised learning branch. MLPs are characterised by iterative learning on the perceptrons present in the network layers [0]. Several arguments can be adjusted in order to modify the shape and parameters of the MLP:

- `hidden_layer_sizes` : number of neurons in each hidden layer
- `activation` : activation function for the hidden layer

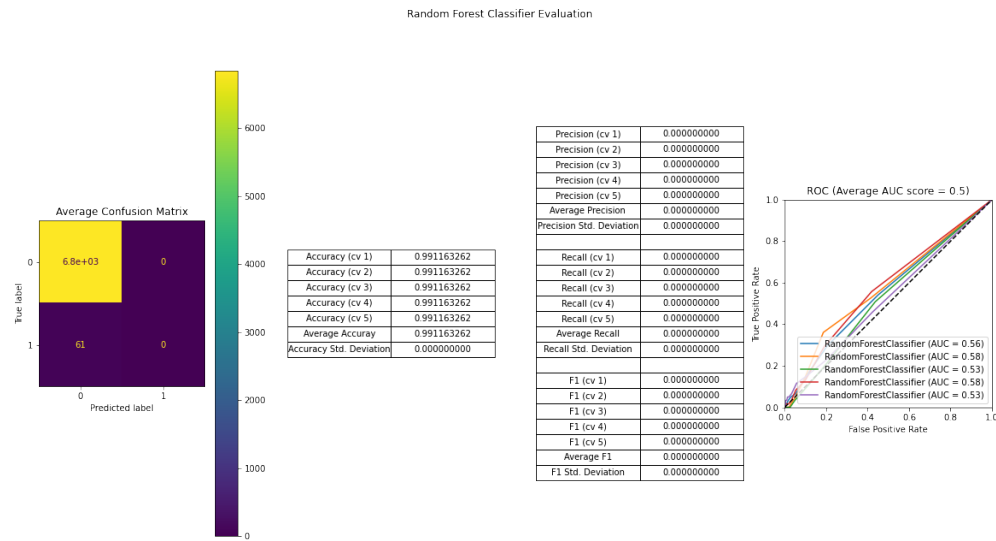


Figure 13: Random Forest Classifier Evaluation (with *class.weight* set to "balanced")

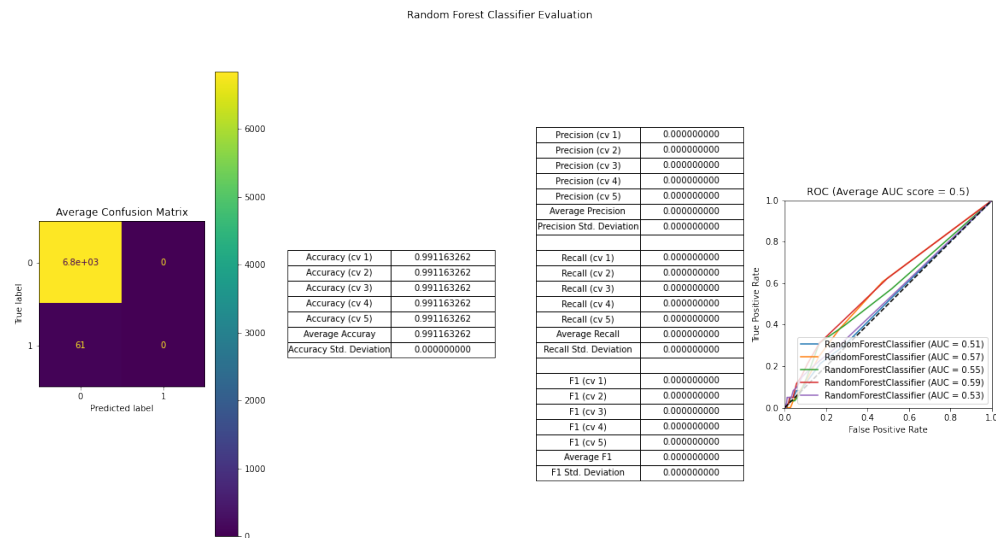


Figure 14: Random Forest Classifier Evaluation

Hyperparameters	Selected value
hidden_layer_sizes	(200,150,100,50)
activation	hyperbolic tangent
solver	stochastic gradient descent
alpha	10^{-3}
batch_size	128
learning_rate	adaptive
max_iter	2000
random_state	1
early_stopping	True
n_iter_no_change	6

Figure 15: Best parameters for the MLP

- solver : solver for weight optimization
- alpha : strength of the L2 regularization term
- batch_size : size of minibatches for stochastic optimizers
- learning_rate : learning rate schedule for weight updates
- max_iter : maximum number of iterations
- random_state : random number generation for weights and bias initialization
- early_stopping : early stopping indicator (terminate training when validation score is not improving)
- n_iter_no_change : maximum number of epochs to not meet tolerance for the optimization improvement

The MLP that gave the best results is composed of an input layer with 12 perceptrons and an output layer with one perceptron.

It is important to note that in the sklearn module, it is not possible to modify the cost function, which by default is the LogLoss function. The Logloss function minimizes the inverse of the log likelihood:

$$LL = -\frac{1}{n} \sum_{i=0}^n (y_i \log(p_i) + (1 - y_i) \log(1 - p_i))$$

The loss is calculated for each batch, n is the value of the batch size, y_i is the output of the neural network and p_i is the associated probability.

There is no method to anticipate the most relevant hyperparameters. The only way is to test several combinations by changing only one hyperparameter at a time and to observe the performance of the resulting network on our data. The best parameters obtained during testing are presented in Figure 15.

Figure 16 shows deceiving results. Despite the time invested in training the network, it is not able to perform better than a random classifier on average (ROC AUC score below 0.5 on Figure 16). This is another bad candidate for the problem.

4.6 Novelty Detection

Anomaly or Novelty detection is a subset of machine learning particularly fitting to unbalanced class problems [5]. Knowing only the data and the estimated proportion of defects, the model will train itself and learn what a valid individual is supposed to look like [3]. It will then be able to classify new individuals: everything that is different than a valid item will be classified defective. Consequently, this means that the model is able to detect defects that are not depicted in the training dataset. This ability might be very be a advantage for

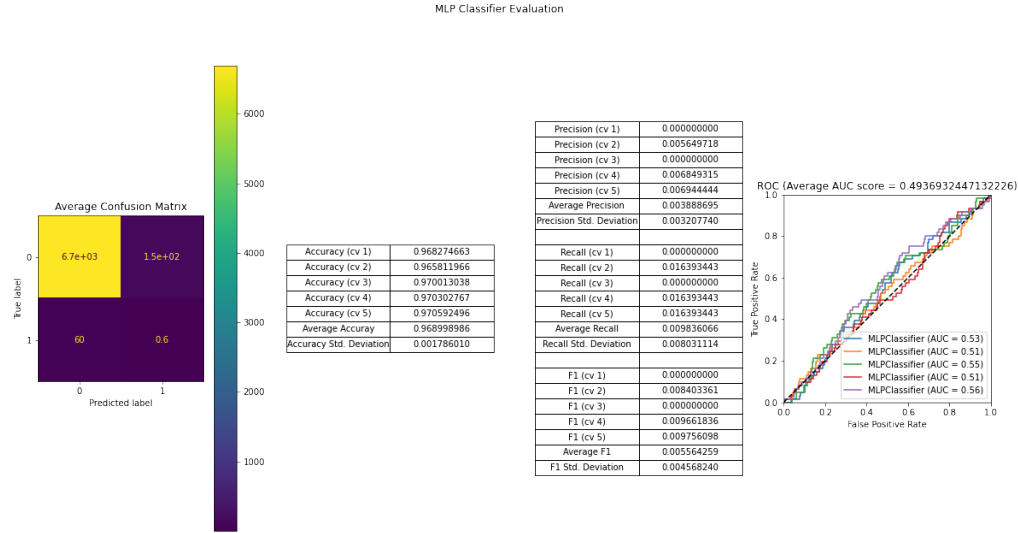


Figure 16: MLP Classifier Evaluation

this application as the model will stay relevant even when new defects occur. For instance, the tools on the assembly line will probably wear down leading to the apparition of new defects.

The model choosen for novelty detection is a One Class SVM classifier [0] with a radial basis function kernel. This classifier did notn perform well during testing with an average precision lower than 1% but it has an impressive average recall of 0.95 with only 2% standard deviation during cross validation.

5 Model Fine Tuning

Having trained and evaluated multiple models adjusted with different sets of parameters, it is now time select the best one and fine tune it in order to improve its performances.

According to previous testing the best model for this particular task is the Naive Bayes Classifier. With a ROC AUC score of 0.6 it managed to achieve an average precision of 1% and a fairly high average recall of 60% all be it with a little bit of deviation between runs (Figure ??).

Using a grid search method [7] it is possible to find the best combination of parameters for a classifier. The Gaussian classifier from de Scikit-Learn library unfortunately has only one. Still, performing a grid search on this parameter which correspond to the portion of the largest variance of all features that is added to variances for calculation stability [4]. Testing a hundred value evenly spaced on a log scale between 1 and 10^{-15} and using the ROC AUC score, it found that a value close to 10^{-7} is the best choice for the *var_smoothing* parameter.

Nevertheless, the results did not improve at all if not deteriorate. The avrage ROC AUC score still hovers around 60%, the average precision looks similar but the average recall dropped form 59% to 53% meaning the model is not as good at finding defective items on the production line.

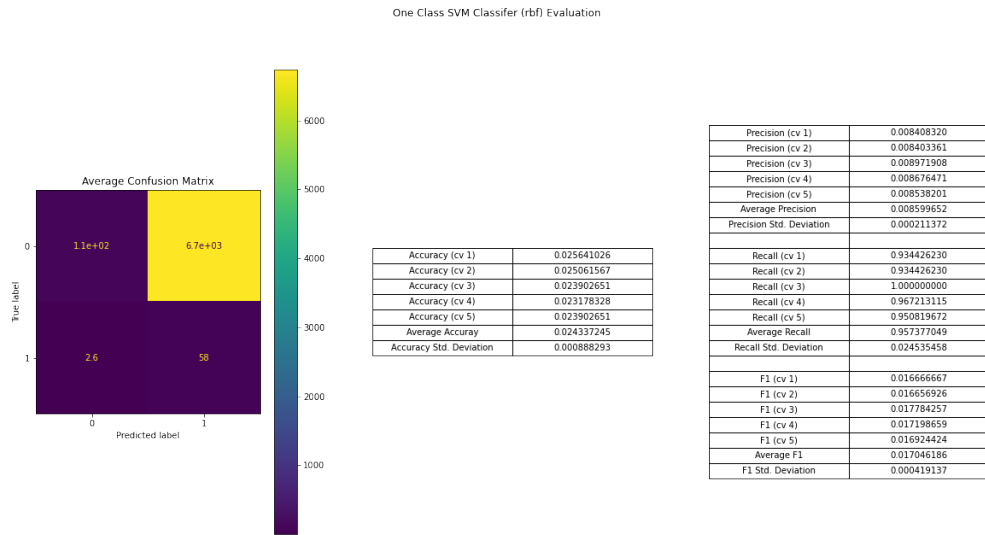


Figure 17: One Class SVM Classifier Evaluation

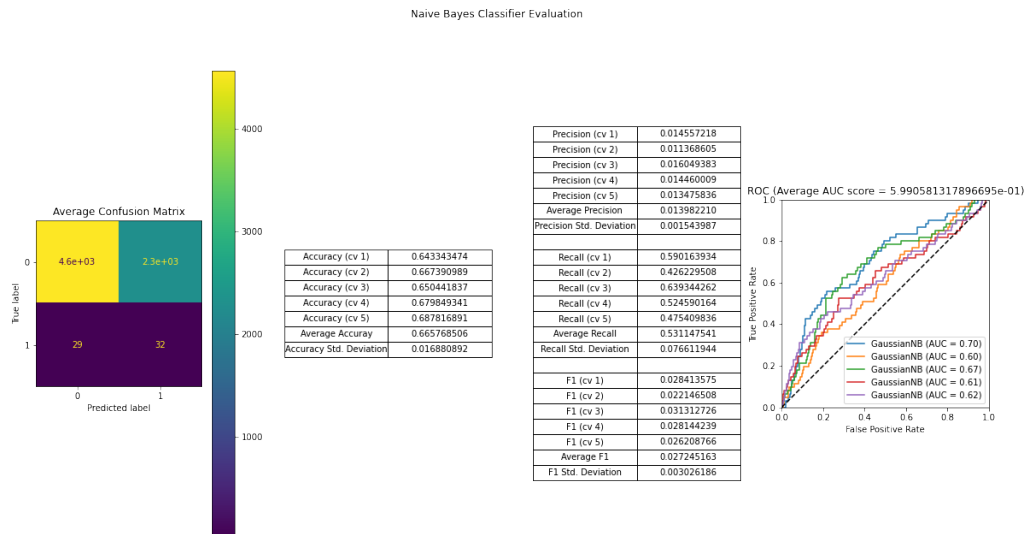


Figure 18: Naive Bayes Classifier Evaluation with tuned parameters

6 Conclusion

Detecting defective items on a production line based on measurements is not an easy task. It requires a good understanding of the situation, a good amount of critical thinking and a lot of testing.

After a proper analysis of the available data, some features were removed and classes needed to be balanced.

Several ML models were trained and the best candidate is the Naïves Bayes Classifier with a ROC AUC score of 60%. Eventhough it has a low precision, it shows a fairly high recall, a good quality for a model aiming at detecting defects on a production line.

To sum up, so as to detect anomalies on production line, the Naïves Bayes classifier was the final choice concerning the projet. Classifier optimization enhance the performance of the model even if it remains perfectible. The main problems have been in avoiding overfitting due to the unbalanced class, which reduce the classifier performance. To improve it, others ways such as Novelty Detections or SMOTE was tried. However, they don't provide enough good results so as to select them. Hyperparameters such as "kernel" or class weight could be the key to solve the problem with a better model.

References

- [1] Defect prediction on production line by valeo. <https://challengedata.ens.fr/challenges/36>. Accessed: 2023.
- [2] Valeo. <https://www.valeo.com/fr/>. Accessed: 2023.
- [3] Grab N Go Info. One class svm for anomaly detection — unsupervised machine learning. https://www.youtube.com/watch?v=0IkFnHpUUjE&ab_channel=GrabNGoInfo. Accessed: 2023.
- [4] Scikit-Learn. Naive bayes. https://scikit-learn.org/stable/modules/naive_bayes.html. Accessed: 2023.
- [5] Scikit-Learn. Novelty and outlier detection. https://scikit-learn.org/stable/modules/outlier_detection.html. Accessed: 2023.
- [6] Scikit-Learn. `sklearn.ensemble.randomforestclassifier`. <https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html>. Accessed: 2023.
- [7] Scikit-Learn. `sklearn.model_selection.gridsearchcv`.. Accessed: 2023.
Scikit-Learn. `sklearn.naive_bayes.gaussianmb`. https://scikit-learn.org/stable/modules/generated/sklearn.naive_bayes.GaussianNB.html. Accessed: 2023.
Scikit-Learn. `sklearn.neighbors.kneighborsclassifier`. <https://scikit-learn.org/stable/modules/generated/sklearn.neighbors.KNeighborsClassifier.html>. Accessed: 2023.
Scikit-Learn. `sklearn.neural_network.mlpclassifier`.. Accessed: 2023.
Scikit-Learn. `sklearn.svm.oneclasssvm`. <https://scikit-learn.org/stable/modules/generated/sklearn.svm.OneClassSVM.html>. Accessed: 2023.
Scikit-Learn. Supervised learning. https://scikit-learn.org/stable/supervised_learning.html. Accessed: 2023.
Wikipedia. Decision tree learning. https://en.wikipedia.org/wiki/Decision_tree_learning. Accessed: 2023.
Wikipedia. k-nearest neighbors algorithm. https://en.wikipedia.org/wiki/K-nearest_neighbors_algorithm. Accessed: 2023.
Wikipedia. Random forest. https://en.wikipedia.org/wiki/Random_forest. Accessed: 2023.