



Regulatory Learning: how to supervise machine learning models? An application to credit scoring

Dominique Guegan, Bertrand Hassani

► To cite this version:

Dominique Guegan, Bertrand Hassani. Regulatory Learning: how to supervise machine learning models? An application to credit scoring. 2017. halshs-01592168v1

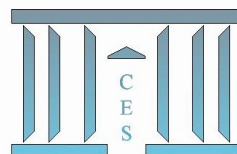
HAL Id: halshs-01592168

<https://shs.hal.science/halshs-01592168v1>

Submitted on 22 Sep 2017 (v1), last revised 23 Oct 2017 (v2)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



**Regulatory Learning: how to supervise machine learning
models? An application to credit scoring**

Dominique GUEGAN, Bertrand HASSANI

2017.34



Regulatory Learning: how to supervise machine learning models?

An application to credit scoring.

July 3, 2017

Dominique Guégan¹, Bertrand Hassani²

Abstract

The arrival of big data strategies is threatening the latest trends in financial regulation related to the simplification of models and the enhancement of the comparability of approaches chosen by financial institutions. Indeed, the intrinsic dynamic philosophy of Big Data strategies is almost incompatible with the current legal and regulatory framework as illustrated in this paper. Besides, as presented in our application to credit scoring, the model selection may also evolve dynamically forcing both practitioners and regulators to develop libraries of models, strategies allowing to switch from one to the other as well as supervising approaches allowing financial institutions to innovate in a risk mitigated environment. The purpose of this paper is therefore to analyse the issues related to the Big Data environment and in particular to machine learning models highlighting the issues present in the current framework confronting the data flows, the model selection process and the necessity to generate appropriate outcomes.³.

¹Université Paris 1 Panthéon-Sorbonne and labEx ReFi, CES, 106 bd de l'Hôpital, 75013 Paris, France, dguegan@univ-paris1.fr

²Group Capgemini, Université Paris 1 Panthéon-Sorbonne, University College London and LabEx ReFi, CES, 106 bd de l'Hôpital, 75013 Paris, France, bertrand.hassani@capgemini.com

³This work was achieved through the Laboratory of Excellence on Financial Regulation (Labex ReFi) supported by PRES heSam under the reference ANR10LABX0095. It benefited from a French government support managed by the National Research Agency (ANR) within the project Investissements d'Avenir Paris Nouveaux Mondes (investments for the future Paris New Worlds) under the reference ANR11IDEX000602.

1 Introduction

The latest trends in financial regulation are simplified models as the current ones were considered too complicated. One can remember the analysis of the Bank of England which shows that models tested on the same market segment, i.e. similar data, were producing results scattered by up to several standard deviations.

Unfortunately for regulators, the arrival of Big Data is going to have an even larger impact. Not because the models are necessarily more complicated but mainly because the data flows will differ from an entity to the next, from one moment to the next, or if a transfer learning strategy is implemented. Even if the data flow does not differ, the model used to capture their evolution will have to evolve dynamically and consistently with the data. Consequently, we already saw that the demands for new models and the decisions made relying on these models are growing exponentially. We expect this movement to amplify in the near future. Thus, the risks associated to the models will only go increasing.

In the meantime, the regulatory and legal considerations that support avoiding companies "too big to let fail" are also changing financial institutions relationships towards data. Indeed, (i) the arrival of the PSD2 in the EU (though the UK have been pioneer in this aspect) which demands financial institutions to transmit third party providers clients' data as soon as these customers requires them, (ii) the enforcement of the right to be forgotten and (iii) the use of non structured data will change the data environment quantitatively and qualitatively for which banks will have to evolve.

So the nebula referred to as big data covers a broad term for data sets so large or complex that traditional data processing applications are inadequate. Challenges include analysis, capture, cleansing, search, sharing, storage, transfer, visualization, and information privacy. The term often refers simply to the use of predictive analytics or other certain advanced methods to extract value from data. Seldomly does it refer to a particular size of data set, though this restrictive perception is probably the best. Accuracy in big data may lead to more confident decision making. Better decisions can mean greater operational efficiency, cost reduction and reduced risk.

Data analysis is the key to the future of financial institutions. Our environment will move from traditional to rational though the path might be emotional. Data analysis allows for looking at a particular situation from different angles. Besides, the possibilities are limitless as long as the underlying data is of good quality. Indeed, data analysis may lead to the detection of correlations, trends, etc. and can be used in multiple areas and industries. Dealing with large data sets is not necessarily easy. Most of the time it is quite complicated as many questions arise related to data completeness, the size of the data sets and reliability of the technological infrastructure.

The work requires parallel computing as well as multiple servers. To summarize, what is considered Big Data embraces the capabilities, the users objectives, the tools deployed and the methodologies implemented. It evolves quickly because what is considered Big Data one year becomes business as usual the next. Depending on the organization, the infrastructure to put in place will not be the same as the needs are not identical from one entity to another. Finally, in Big Data management there is no "one-size-fits-all", and any piece of information is a risk data.

This paper points out the necessary legal environment for Big data using models and data flows. Thus, the paper is organized as follows: in the next section we specify the legal and regulatory environment necessary to develop, which will open some research routes. In Sections three and four, we present the machine learning techniques we use. In Section five, we present a credit data application and discuss our results. A last section concludes and discuss the possibility to regulate such a topic and even more importantly the possibility to supervise it. We also address the cultural change required from both sides, i.e. the necessity of financial institutions and the regulators to adopt a better quantitative culture.

2 Legal Environment

As implied in the previous section, Big Data is a buzzword used frequently in private and public sectors, the press, and online media. Large amounts of money are being invested to make companies Big Data users. Governmental institutions are eager to experiment Big Data applications in the fields of crime prevention, intelligence, fraud, among others. Though the exact nature

and delineation of Big Data is still unclear, it seems likely that Big Data will have an enormous impact on our daily lives.

Though these impacts expected to be quite positive, there are naturally inherent risks to Big Data applications because it might result in discrimination, privacy violations, and model risk. Therefore an adequate framework is necessary to ensure that the beneficial uses of Big Data are promoted and facilitated, while the negative effects are mitigated. Van der Sloot and Van Schendel (2016) provides building blocks for developing such a framework by giving an overview of the experience in the use and regulation of Big Data in multiple countries. Though their analysis is mainly focused on the use of Big Data by governments, most findings are applicable to banks and financial institutions.

Their analysis shows that the dangers of Big Data may materialize in two areas: (i) a possible violation of the right to privacy or to data protection and (ii) the danger of discrimination and stigmatization: "Regarding the first point, it appears from underlying research that most countries are well aware of the risks to the privacy of citizens. With regard to the risk of discrimination and stigmatization, this appears to be true to a lesser extent" for governments. We believe that in financial institutions this might become a real issue (as shown with an example developed in Section 5). According to their analysis, it appears that, in most countries, the current regulations in the area of privacy and data protection are applied to Big Data processes: Germany with the distinctive personality right, the United States without an umbrella law for the regulation of privacy, but with sectoral legislation, and most other countries with relatively similar rules concerning privacy and data protection. Current legislation is generally applicable to Big Data, but it is necessary to formulate new rules regarding models and analytics. Now, we will focus on the technologies and show through an example the necessity of considering the regulatory learning framework.

3 Machine Learning Prerequisites: The Data

A point needs to be made absolutely clear before any further presentation. None of the methodologies presented in the following sections can be used if they are not fed by appropriate inputs.

Therefore, we will start this chapter defining data. Then we will discuss pre-processing these inputs.

Data is a set of qualitative and quantitative pieces of information. Data is engendered and obtained by both observation and measurement. It is collected, reported, analyzed, and visualized. Data as a general concept refers to the fact that some existing information is represented in some form suitable for processing. Raw data, or unprocessed data, is a collection of numbers and characters; data processing commonly occurs by stages, and the "processed data" from one stage may be considered the "raw data" of the next one. Field data is raw data that is collected in an uncontrolled environment. Experimental data is data that is generated within the context of a scientific investigation by observation and recording, this is data generated during analysis. It is particularly important to understand, particularly for scenario analysis, that the data used to support the process are not usually numeric values. Indeed, these are usually pieces of information gathered to support a story line, such as an article, media, incidents experienced by financial institutions and expert perceptions.

Data is any facts, numbers, or text that can be processed. Nowadays, organisations are gathering and accumulating vast and growing amounts of data in various formats and databases. We can split the data into three categories:

- operational or transactional data: sales, cost, inventory, payroll, and accounting
- non operational data: industry sales, forecast data, and macro economic data
- meta data - data about the data itself, such as logical database design or data dictionary definitions

Recent regulatory documents, for instance the Risk Data Aggregation (Palace (1996)), aim at ensuring the quality of the data used for regulatory purposes. However, one can argue that any piece of data could be used for regulatory purposes. In the long term, this piece of regulation should lead in the long term to better risk management. Indeed, Palace (1996) requires that the information banks used in decision-making process capture all risks accurately as well as timely. This piece of regulation sets out principles of effective and efficient risk management by pushing banks to adopt the right systems and develop the right capabilities instead of ticking regulatory

boxes to be compliant at a certain date.

Before being in a position to use this information, it has to be pre-processed, i.e. made usable by models. For that we can rely on data mining strategies (Hastie et al. (2009)). The purpose of data mining is to extract information from data sets and transform it into understandable structure given the ultimate use of the data. The automatic or semi-automatic analysis of large quantities of data allows to detect interesting patterns such as clusters (Everitt et al. (2011)), anomalies, and dependencies. Outcomes of the analysis can then be seen as the essence or the quintessence of the original input data, and may be used for further analysis in machine learning, predictive analytics and more traditional modelling. Until now, data mining was mainly used by companies with a strong consumer focus, in other words: retail, financial, communication, and marketing organizations (Palace (1996)). We recall some of the major elements of data mining. They are:

- Data capture: data is collected from various sources and gathered in a data base.
- Data pre-processing, i.e. before proper mining:
 - Data selection: Given the ultimate objective, particular data needs to be selected as these will be used for further analysis.
 - Data cleansing and anomaly detection: Collected data may contain errors, may be incomplete, inconsistent, outdated, erroneous, etc. These issues need to be identified, investigated and dealt with.
 - Data transformation: Following the previous stage, the data is not ready for mining, it requires transformations such as kernel smoothing, aggregation, normalisation, interpolation etc.
- Data processing is only possible once the data has been cleansed and fit for purpose. This step combines:
 - Outlier detection: finding an observation point that is distant from other observations.
 - Relationship analysis: gathering data with similar characteristics, classification and analysing interactions.

- Pattern Recognition: regression analysis, time series analysis, distributions etc.
- Summarization and knowledge presentation: This step deals with visualization, as one should beware that key aspects are not lost in the middle of the process.
- Decision making process integration: this step uses the knowledge obtained from the previous manipulations.
- The analysis. This is the ultimate objective of data mining.

Once the preliminary work has been achieved, it is time to analyse and process data. This will be addressed in the next section.

4 Machine Learning and Artificial Intelligence

Once the data has been correctly formatted for specific objectives, it can be used for prediction, forecasting, evaluation, in other words, for modelling. Now, we will describe some of the tools that can be used. We will use them in the example provided in the next Section by analysing the capability of a company to reimburse a loan.

Machine learning deals with the study of pattern recognition and computational learning theory in artificial intelligence. Machine learning aims at building algorithms that can learn from data and make predictions from them. This means it operates dynamically, adapting itself to changes in the data, not only relying on statistics, but also on mathematical optimization. Automation is a keyword for all this technology. The objective being to make machines mimic the cognitive processes exhibited in the human brain.

Machine learning tasks are usually classified into four categories (Russell and Norvig (2009)) depending on inputs and the objective. These are:

- Supervised learning (Mohri et al. (2012)): the goal is to infer a general rule from example data mapped to the desired output. The example data is usually called training data. This consists in couples of input object and desired output, or supervisory signal. Once the algorithm analyses the training data and infers the function, it can be used to map new examples, and generalize its use to previously unknown situations. Optimization algorithms

should perfectly react to new instances providing unbiased and accurate outcomes, e.g. a methodology with outcomes that are accurate once they can be compared with the real occurrences.

- Unsupervised learning: in this case no training data is given to the learning algorithm, consequently it will have to find patterns using the inputs. Unsupervised learning can actually be used to find hidden structures and patterns embedded within the data. Therefore, unsupervised learning aims at finding hidden patterns from unlabelled data (Hastie et al. (2009)). In the case of unsupervised learning, it is not possible to evaluate the quality of the solution as we have no benchmark to go off of.
- Semi-supervised learning: when the initial training information (i.e. data and/ or targets) is incomplete, an intermediate strategy is used.
- Reinforcement learning (Sutton and Barto (1998)): it is a program that interacts and evolves within a dynamic environment in which it is supposed to achieve a specific task. However, once again there is no training data and no benchmark. This approach aims at learning how to map situations to actions, so as to optimise a numerical function, i.e. the output. The algorithm has to discover which actions lead to the best reward signal by trying them. These strategies capture situations where actions may affect all subsequent steps with or without any delay. This might be of interest.

The machine learning goal is to make accurate prediction relying on the generalization of patterns originally detected and refined by experience. In the following, we briefly present the models (which are part of the machine learning approach) implemented in the credit risk application. That will be further referenced in the next section.

- Logistic Regression: this is a regression model where the dependent variable is categorical. In this paper we consider the case of a binary dependent variable. Cases where the dependent variable has more than two outcome categories may be analysed in multinomial logistic regression, or, if the multiple categories are ordered, in ordinal logistic regression.
- Least absolute shrinkage and selection operator or lasso: this is a regression analysis method performing both variable selection and regularisation to improve the prediction accuracy and the interpretability of the statistical model created. Originally formulated for least

squares models, the Lasso led to a better approximation of the behaviour of the estimator (including its relationship to ridge regression) and a much better selection of subsets. As well, the Lasso approach provides valuable information on the connections between coefficient estimates and so-called soft thresholding. As for standard linear regressions, if covariates are colinear, the coefficient estimates do not need to be unique (Tibshirani (1996)).

- Decision tree learning and Random Forest (deVile (2006)): this is used to predict the values of a target variable based on several inputs, which are graphically represented by nodes. Each edge of a node leads to children representing each of the possible values of that variable. Each leaf represents a realisation of the target variable given the input variables represented by the path from the root to the leaf. A decision tree may be implemented for classification purposes or for regression purposes, respectively to identify to which class the input belongs or to evaluate a real outcome (prices, etc.). Some example of decision tree strategies are Bagging decision trees (Breiman (1996)), Random Forest classifier, Boosted Trees (Hastie et al. (2009)) and Rotation forest.
- Artificial neural networks: these are learning algorithms that are inspired by the structure and functional aspects of biological neural networks. Modern neural networks are non-linear statistical data modelling tools. They are usually used to model complex relationships between inputs and outputs, to find patterns in data, and to capture the statistical structure in an unknown joint probability distribution between observed variables. Artificial neural networks are generally presented as systems of interconnected "neurons" which exchange messages between each other. The connections have numeric weights that can be tuned based on experience making them adaptive to inputs and capable of learning. Neural networks might be used for function approximation, regression analysis, time series, classification, including pattern recognition, filtering, clustering, among others.
- Support vector machines (Cortes and Vapnik (1995), Ben-Hur et al. (2001)): these are supervised learning models in which algorithms analyse data and recognize patterns. They are usually used for classification and regression analysis. Given a set of training data - each of them associated to one of two categories - the algorithm assigns new examples to one of each two categories based on fully labelled input data. The parameterization is quite

complicated to interpret. This strategy can also be extended to more than two classes, though the algorithm is more complex. Interesting extensions are (i) the support vector clustering which can be used as an unsupervised version, (ii) the transductive support vector machines which is a semi-supervised version, or (iii) the structured support vector machine.

These approaches have all been implemented in the following credit application and we summarise the outcomes obtained for each approach.

5 Credit Scoring Application

The objective of this Section is to illustrate the differences of the selected models (introduced previously) in the context of Big Data comparing the resulting Gini index of each model. The exercise has been performed at time t , and the results are nothing more than snapshots. They are not a predictive picture for the future. The discussions will refer to the choice of the models, the choice of the indicator to determine what could be the best model, and the question of the dynamics: how to introduce them in order to provide a figure of the future risks.

343 factors representing the credit repayment capability (for instance the turnover, the gross margin, the result, the number of employees, the industry turnover, etc.) of a set of 12 544 companies (SME) over the year 2016 have been considered for evaluation on their probability to default. We did this using the 6 models presented in the previous section: a logistic regression, a Lasso approach, a random forest (simple or considering gradient boosting), a neural Network approach and Support Vector Machine strategy (SVM). In order to rank the models with respect to companies credit worthiness, the Gini index (Gini (1936)) and the ROC curve (Hanley and McNeil (1982)) are being used. A last classification is provided through a Box Plot representation.

One of the stakes when working with Big Data is dimension reduction. Thus, to obtain a score associated to each company the most pertinent factors (or variables) characterising default risk have been selected. So, even if the exercise proposed here does not address all the questions and does not systemically perform the optimal search of these factors. In our approach the outcome of the most advanced machine learning models has been benchmarked with the most traditional

approach used within financial institutions, i.e. logistic regression. In parallel of the presentation of the outcomes, we describe in the following how the models of interest have been implemented.

1. The first model implemented is a logistic regression. To apply the logistic regression, the dependent variable is the defaults/non-default for the companies and we consider a set of 343 variables (the factors which can explain the behaviour of these companies). We retain 23 factors. These 23 variables have been selected following the elimination of the correlated variables and the variables not properly collected at a time. As discussed previously, the logistic regression will capture interaction between a dependent variable and various independent variables.

For each company we compute the Gini index and we associate to it the ROC curve that we represent on Figure 1. The regression model is adjusted considering 80% of the data (60% for the fitting and 20% for the cross validation) and then 20% of this data is used for test purposes. A ROC curve is then used to evaluate the quality of the model. Recall that the Gini index is equal to $2 * AUC - 1$ where AUC correspond to the "Area under the curve". Then just to check the quality of the model initially created we successively and randomly removed some variables. It is interesting to note that the AUC value kept increasing, so in our case less information led to a better adjustment. The curves represent for each cut-off points how many "good" are approved over "total good" and how many "bad" are approved over the "total bad". Therefore the perfect model is the one that has a perfect cut-off (100% good, and 0% bad approved) so the ROC is 1.

In Figure 1 we observe that the better ROC Curve is provided for $AUC = 0.7626$, because it corresponds to the best trade-off between the "good/total good" and "bad/total bad". It is interesting to note that the performance obtained in that last case has been obtained only with 17 variables. Thus, it is not the number of variables which is predominant but the relevance of the variables for the objective we have.

2. In order to be close to the reality of the possible defaults of the companies, we want to analyse in more detail the impact of the factors and then we choose to consider a Lasso ap-

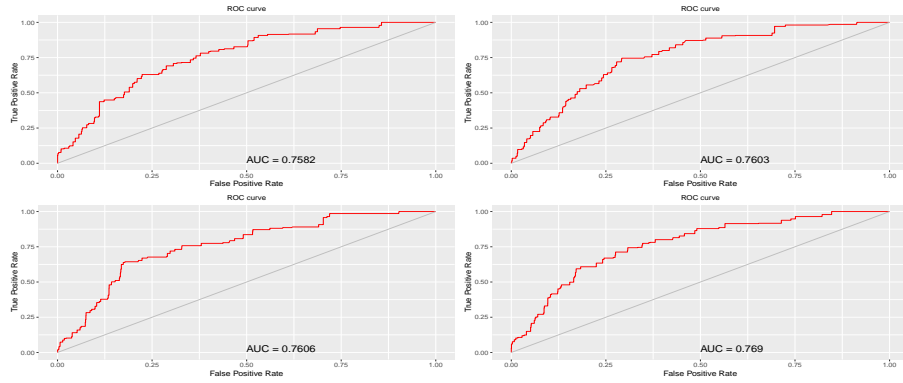


Figure 1: Results obtained implementing the logistic regression. We observe on this figures that for every time a variable is removed the Gini index increases. Note that the bottom right hand graph represents a logistic regression performed on the 343 variables simultaneously, and therefore without any variable selection. The curves represent for cut-off how many "good" are approved over "total good" and how many "bad" are approved over the "total bad". Therefore the perfect model is the one that has a perfect cut-off (100% good, and 0% bad approved) so the ROC is 1.

proach which consists in a penalisation function to avoid overfitting. This one penalisation function requires the fitting of an additional parameter usually referred to as λ ⁴.

The adjustment for that parameter is presented in figure 2. Each coloured line represents the value taken by a different coefficient in your model. λ is the weight given to the regularization term (the L1 norm), so as λ is getting close to zero, the loss function of the model approaches the OLS loss function. Consequently, when λ is very small, the LASSO solution should be very close to the OLS solution, and all of the coefficients are in the model. As λ grows, the regularization term has a greater effect and we see fewer variables in the model.

On Figure 3, we provide the ROC curve associated to the Gini index when we apply the Lasso regression to the best previous case, corresponding to $AUC = 0.7626$. Comparing

⁴Zou and Hastie (2005) elastic net strategy has been implemented considering a second quadratic penalisation function.

the two curves we observe that the percentages of the "good/total good" companies versus "bad/total bad" companies increases. Indeed, the curve increases quickly toward 1 being close to zero⁵.

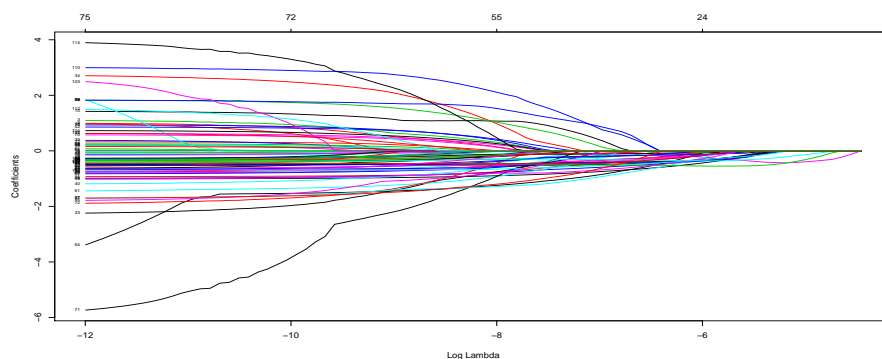


Figure 2: This figure represents the impact of the regularization on the coefficient and therefore the most appropriate λ parameter.

3. The third model created is a random forest. We used all the data available as random forests are supposed to capture linear and non linear dependencies. In our case, to evaluate the probability of default of customers, the model functions as a successive segmentation program. Each variable is split in two and then reinjected in a subsequent layer if their is some valuable information remaining. A random forest can be represented as the combination of a binary decision tree and a bootstrap strategy. Figures 4 and 5 provide us with an illustration of the random forest obtained considering the data used. In Figure 5, we have used a different process to optimize the random forest with less iterations, minimizing the error at each step, thus it appears more informative than what is represented in Figure

⁵in R, the function runs `glmnet` $n_{folds} + 1$ times; the first to get the lambda sequence, and then the remainder to compute the fit with each of the folds omitted. The error is accumulated, and the average error and standard deviation over the folds is computed. Note that `cv.glmnet` does NOT search for values for alpha. A specific value should be supplied, else `alpha=1` is assumed by default. If users would like to cross-validate alpha as well, they should call `cv.glmnet` with a pre-computed vector `foldid`, and then use this same fold vector in separate calls to `cv.glmnet` with different values of alpha. Note also that the results of `cv.glmnet` are random, since the folds are selected at random. Users can reduce this randomness by running `cv.glmnet` many times, and averaging the error curves.

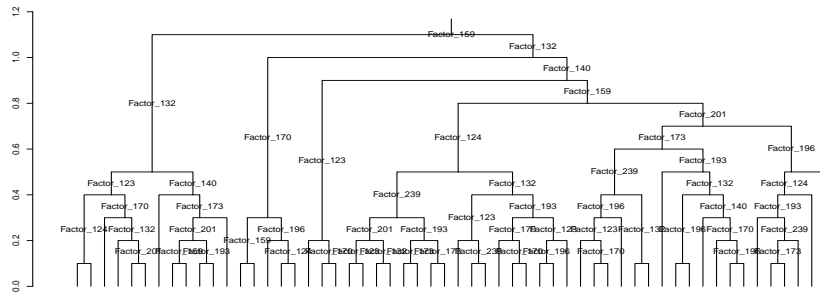


Figure 5: Tree obtained implementing a Random forest considering a second set of variables

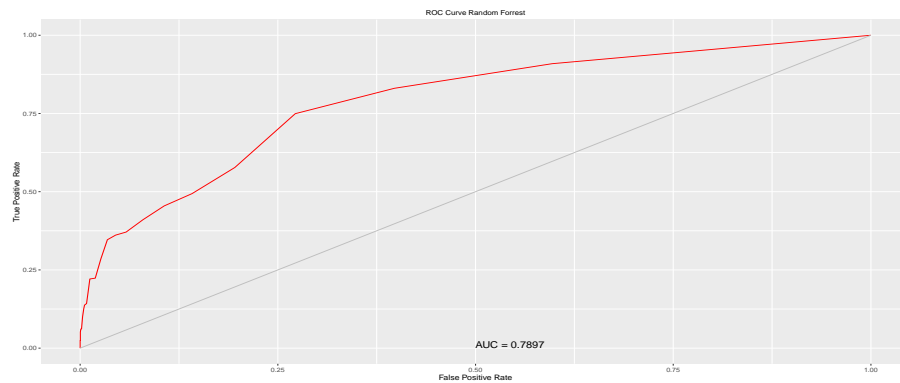


Figure 6: Results obtained implementing the Random Forest approach. The curves represent for cut-off how many "good" are approved over "total good" and how many "bad" are approved over the "total bad". Therefore the perfect model is the one that has a perfect cut-off (100% good, and 0% bad approved) so the ROC is 1.

the weight updating coefficient. For this approach, in Figure 7 we provide the ROC curve which still represents the impact of a choice of a learning function over another on the level of false positive versus true positive. This impact is non negligible and could lead financial institutions to a higher level of default. Indeed, we observe that we have increased the level of the ROC curve.

5. A Support Vector Machine approach has been implemented. Recall that the SVM is used here for regression purposes as it allows scoring customers evaluating their probability of

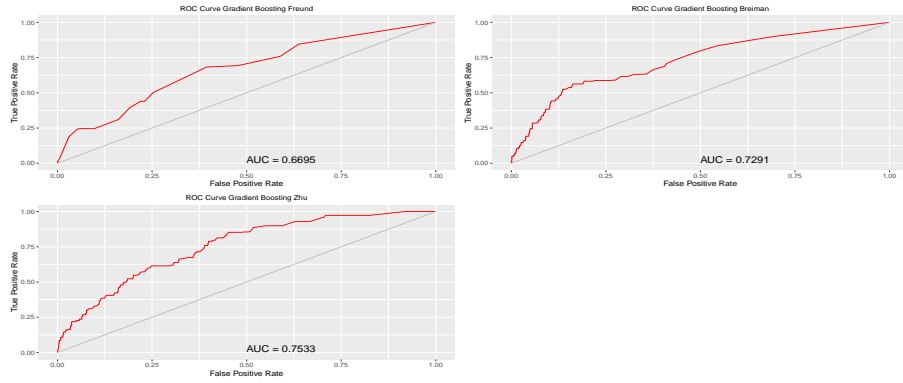


Figure 7: Results obtained implementing the Boosting approach how many "good" are approved over "total good" and how many "bad" are approved over the "total bad". Therefore the perfect model is the one that has a perfect cut-off (100% good, and 0% bad approved) so the ROC is 1.

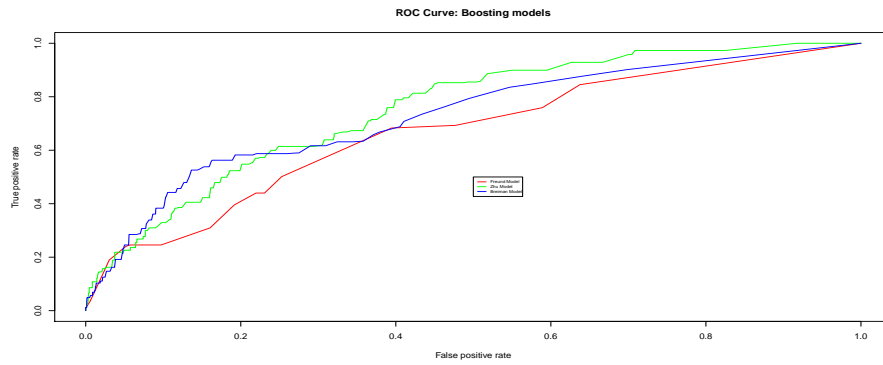


Figure 8: This figure allow comparing boosting approaches. The methodology is similar to the random forest. The difference is that it works to reduce the error at each step of the process. We compare three error functions: Breiman's function, Freund's function and *Zhu*' SAMME algorithm is implemented with $\alpha = \ln\left(\frac{1-err}{err}\right) + \ln(nclasses - 1)$

default considering different factors. It provides a classification of the companies which is a little different approach than the previous approaches. Here the factors selected are once again different from those used in the previous models. Indeed, the probability model for classification fits a logistic distribution using maximum likelihood to the decision values of all binary classifiers, and computes the a-posteriori class probabilities for the multi-class problem using quadratic optimization. The probabilistic regression model assumes (zero-

mean) Laplace-distributed errors for the predictions, and estimates the scale parameter using maximum likelihood. Here, Figure 10 represents the two area representing potential default and non-default.

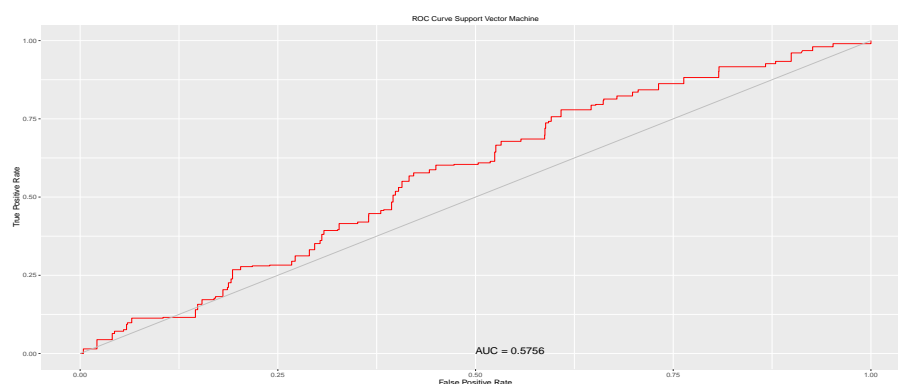


Figure 9: Results obtained implementing the Support Vector Machine approach how many "good" are approved over "total good" and how many "bad" are approved over the "total bad". Therefore the perfect model is the one that has a perfect cut-off (100% good, and 0% bad approved) so the ROC is 1.

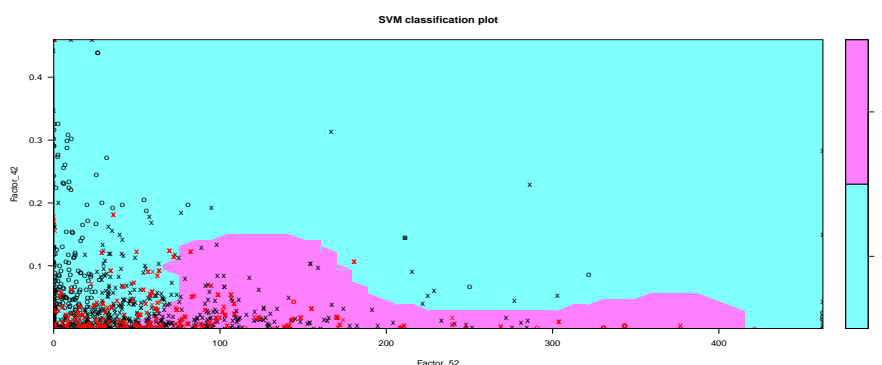


Figure 10: This figure illustrates the SVM methodology, which in our case provides us with the worst adjustment results.

6. Finally we have implemented a neural network. This one takes the same data set as for the SVM, and consider the same activation function as the logistic regression though in that case the way the factors are combined within the hidden layer is different. Indeed, the

weights associated to the factors considered in the learning is different, therefore it naturally results in a different set of outcomes. The ROC associated to the Neural Network approach is given in Figure 11.

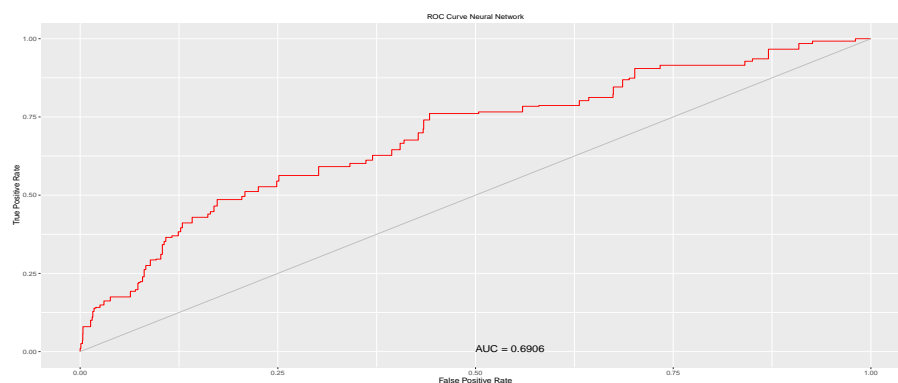


Figure 11: Results obtained implementing a Neural Network approach showing how many "good" are approved over "total good" and how many "bad" are approved over the "total bad". Therefore the perfect model is the one that has a perfect cut-off (100% good, and 0% bad approved) so the ROC is 1.

7. Finally, we provide in Figure 12 a way to compare the 6 different modellings computing the variance of the Gini index and its level of randomness, except for the logistic regression. On this figure we observe that on average, the Random Forest offers a better explanatory power relying on a larger number of variable, however, the improvement offered by the random forest compared to the logistic regression is not very large. Indeed, the logistic regression relies on a subset of the pool of variables used for the Random Forest. However, if we draw a parallel between this figure and the first one presented in that section, we may conclude that these results are only valid considering the current data set. Indeed, the first figure tells us that considering the set of data approved months ago does not provide the best adjustment as randomly removing some variables leads to better Gini indexes, while using the full set of 343 variables leads to an improvement as soon as a random forest is selected. This statement leads us to the conclusion that if the data sets and the types of variables change (for example using unstructured data), the order of the models may change again.

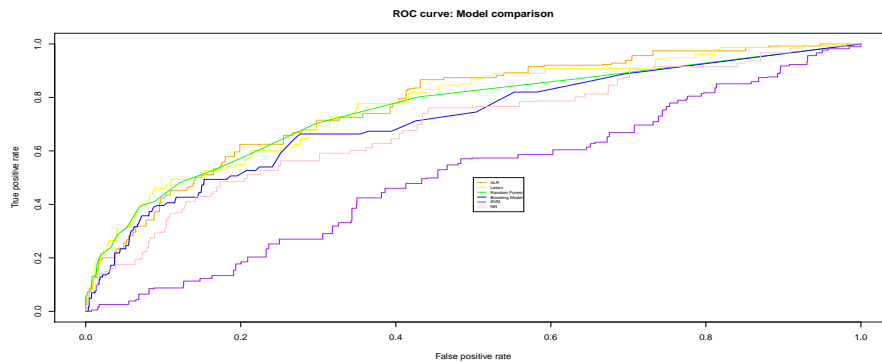


Figure 12: This figure allows comparing the quality of the models adjusted in terms of Gini drawing the ROC curves of each model parallel to each other.

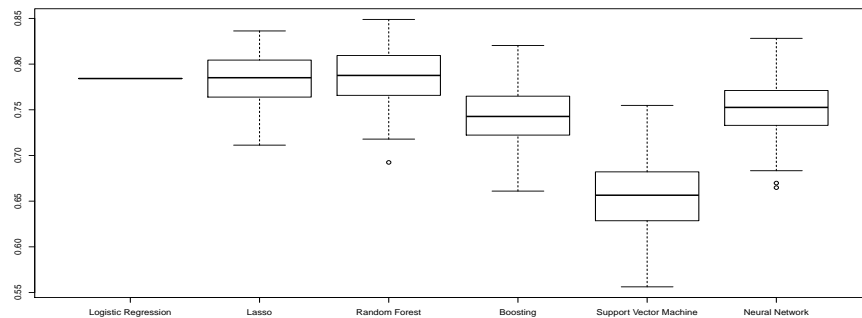


Figure 13: This figure allows comparing the quality of the models adjusted in terms of Gini. The boxes provides a measure of the volatility of the gini.

Considering the results presented in the previous section, it is interesting to note that though all the models could be used to achieve the same tasks, the way to perform the computations are not identical and do not philosophically imply the same thing. Indeed, though for the logistic regression, the data has been selected such that no correlations are present, the other models do not necessarily require that. The random forest strategies have the capability to capture non-linear dependencies, while neural networks are entirely based on building and capturing relationships between variables.

Furthermore, these results have been obtained on a static data set; therefore the evolution of the

data set may lead to a different ranking of the models. Indeed, here though random forests seems in average better than the others, this is not necessarily the same if the data set were to evolve. Besides, we do not know yet how the integration of other data sources (such as integration of social medias) would impact this ranking, however, considering the nature of the models, we are inclined to think that neural networks or random forests would be more appropriate than the others.

Finally, it is interesting to note that Figure 13 has been obtained considering a single sample based on 80% of the original data, and we observed that if another subsample of 80% of the original data were to be used, the Gini index of the logistic regression may decrease by up to 17%, while for other approaches, the results would still be contained in the confidence intervals of the box plots. This observation may greatly impact the original ranking disclosed in this paper. We have also observed that the selection of the activation function, the algorithms embedded within the machine learning approaches presented, and the index used to evaluate the quality of the models may lead to different results too. Therefore, this point should be carefully handled.

6 Conclusion

While the regulation and the legislation related to big data is still embryonic, the implementation of machine learning models has to be done carefully as the dynamic philosophy implied cannot yet properly be handled by the current static regulation and internal governance.

As we have shown in Section 5, the choice of the models is not naive, nor was the choice of the indicator we chose to have as an answer to our objective. This exercise points to a number of questions that will be the purpose of companion papers. If it appears necessary to compare different models to a given data set, this exercise would be even more relevant considering a larger data set and considering data coming from alternative sources. Indeed, with a larger number of data points, naturally containing more information, more dependence, some non linearity and specific features like trends, seasonalities and so on, should lead to more relevant outcomes, as well as potentially controversial questions. The choice of the indicator to retain the "best " model is also very difficult when we base our modelling on machine learning: it is important to know

exactly what are the underlying processes for automation to understand the specific features we retain, the question of the number of iterations is also important to get a form of convergence, and what convergence. This point is also important to discuss.

Concerning the existence of specific features, we cannot analyse strong dependence if all the models are based on linear regression. The question also is to exactly know the error we measure with this kind of modelling, are we close to a type 1 or a type 2 error? The application we provide gives a picture at time t of the capability of companies to obtain a loan or not based on their creditworthiness, therefore an important question is how to integrate new pieces of information in these procedures.

Finally, this thoughtful process leads to questions concerning regulators. There is a strong possibility of risk materialisation depending on the chosen models; the attitude of the regulator - historically when dealing with the risk topic - is to be conservative or close to a form of uniformity. This posture seems complicated with Big Data as the emergence of continuous flows mechanically prevents imposing the same model to all institutions. Dynamics need to be introduced, and it is not simple. We point out the fact that, even if all the models used are very well known in the non-parametric statistics, these dynamics need to be clearly understood for the use they provide. Furthermore, target specifications and indicators need to be put in place in order to avoid large approximations and bad results.

References

- Ben-Hur, A., Horn, D., Siegelmann, H. and Vapnik, V. (2001), 'Support vector clustering', *Journal of Machine Learning Research* **2**, 125–137.
- Breiman, L. (1996), 'Bagging predictors.', *Machine Learning* **24(2)**, 123–140.
- Cortes, C. and Vapnik, V. (1995), 'Support-vector networks', *Machine Learning* **20(3)**, 273–297.
- deVille, B. (2006), *Decision Trees for Business Intelligence and Data Mining: Using SAS Enterprise Miner*, SAS Press, Cary (NC).

Everitt, B., Landau, S., Leese, M. and Stahl, D. (2011), *Cluster Analysis*, 5th edition, John Wiley and Sons, New York.

Gini, C. (1936), ‘On the measure of concentration with special reference to income and statistics’, *Colorado College Publication, General Series* **208**, 73–79.

Hanley, J. A. and McNeil, B. J. (1982), ‘The meaning and use of the area under a receiver operating characteristic (roc) curve.’, *Radiology* **143**(1), 29–36.

Hastie, T., Tibshirani, R. and Friedman, J. (2009), *The Elements of Statistical Learning: Data mining, Inference, and Prediction.*, Springer, New York.

Mohri, M., Rostamizadeh, A. and Talwalkar, A. (2012), *Foundations of Machine Learning*, The MIT Press, Cambridge (MA).

Palace, W. (1996), ‘Data mining: What is data mining?’.

URL: http://www.anderson.ucla.edu/faculty_pages/jason.frand/teacher/technologies/palace/datamini

Russell, S. and Norvig, P. (2009), *Artificial Intelligence: A Modern Approach*, (3rd ed.), Pearson, London.

Sutton, R. and Barto, A. (1998), *Reinforcement Learning: An Introduction*, A Bradford Book, The MIT Press, Cambridge (MA).

Tibshirani, R. (1996), ‘Regression shrinkage and selection via the lasso’, *Journal of the Royal Statistical Society. Series B (Methodological)* pp. 267–288.

Van der Sloot, B. and Van Schendel, S. (2016), ‘Ten questions for future regulation of big data: A comparative and empirical legal study’, *JIPITEC* **7**(2).

URL: <http://nbn-resolving.de/urn:nbn:de:0009-29-44386>

Zou, H. and Hastie, T. (2005), ‘Regularization and variable selection via the elastic net’, *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **67**(2), 301–320.