

Financial Fraud Detection Using Machine Learning Techniques

Yordan Ivanov

1. Introduction

What we provide here is an extensive study of machine learning methods on different both real-world and simulated datasets in an attempt to provide better guidelines for fraud detection.

- specify the methods that are going to be used in the study

2. Literature Review

3. Methodology

Classification is one of the most widely used model framework used for the application of machine learning techniques in terms of FFD (Ngai et al. 2011). Some of the most common classification techniques include logistic regression, neural networks, support vector machine and decision trees.

3.1. Preliminaries

In the current section, we will give a description of the machine learning techniques that we apply to predict fraudulent transactions.

Let us define our binary classification dataset as $\{(x_1, y_1), (x_2, y_2), \dots, (x_p, y_p)\}$, where $x_i \in \mathbb{R}^n$ represents an n-dimensional data point and $y_i \in \{-1, 1\}$ represents the label of the class of that data point, $i = 1, \dots, p$. Let X represent the vector of features and Y the vector of dependent variables.

3.2. Machine Learning Algorithms

3.2.1. Logistic Regression

The logistic regression framework falls under the category of generalized linear models and allows the prediction of discrete outcomes. Then, by defining the probability of a transaction being fraudulent by $p(X) = Pr(Y = 1|X)$, we can portray the relationship between the dependent and independent variables as follows:

(???)

$$p(X) = \frac{e^{\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p}}{1 + e^{\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p}}$$

The number of independent variables is indexed by p . After manipulating (???), we can also see that

$$\log\left(\frac{p(X)}{1 - p(X)}\right) = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p$$

with the LHS being called the logit. Using equation (), we will predict the probabilities of a transaction being fraudulent i.e. $p(Y = 1)$. The fitting of a logistic regression is done by the method of maximum likelihood (see Appendix). The logistic regression has been among the most widely used framework in fraud detection (Ngai et al. 2011) due to simplicity of ease of implementation, but it does have its shortcomings - it tends to underperform when there are multiple or non-linear decision boundaries (*SEARCH FOR SOME PAPER OR BOOK ON IT?*)

3.2.2. Neural Networks

3.2.3. Support Vector Machines

Support Vector Machines, developed by Vapnik et. al. (Cortes and Vapnik 1995), have become a popular machine learning method that has seen its implementation rise in various domains that require the use of classification models (Batuwita and Palade 2013). Among the factors for its success is the fact that the SVMs are linear classifiers, which work in a high-dimensional feature that represents a non-linear mapping of the input space of the problem being dealt with (Bhattacharyya et al. 2011). Working in a high-dimensional

feature space has its benefits - often, the problem of non-linear classification in the original input space is transformed to a linear classification task in the high-dimensional feature space.

The goal of the SVM classifier consists of finding the optimal separating hyperplane, which manages to effectively separate the observations from the data into two classes. As mentioned above, the observations are initially transformed by a nonlinear mapping function Φ . Thus, we can write a possible separating hyperplane that resides in the transformed higher dimensional feature space by:

$$w \cdot \Phi(x) + b = 0$$

with w the weight vector normal to the hyperplane.

We will further use two variations of the SVM soft margin optimization problem - one that assigns the same cost for missclassification of the different classes and one that penalizes more the missclassification of the minority class.

Non-cost sensitive learning

For the same missclassification cost case, we can write the soft optimization problem as follows:

$$\begin{aligned} \min & \left(\frac{1}{2} w \cdot w + C \sum_{i=1}^p \xi_i \right) \\ \text{s.t. } & y_i(w \cdot \Phi(x_i) + b) \geq 1 - \xi_i \\ & \xi_i \geq 0, i = 1, \dots, p \end{aligned}$$

The slack variables $\xi_i > 0$ hold for missclassified examples. Thus, the penalty term $\sum_{i=1}^p \xi_i$ can be perceived as the total number of missclassified observations of the model. Thusm from (), we can see that there are two goals - maximizing the margin the minimizing the number of missclassifications. The cost parameter C controls the trade-off between them, i.e. assigned misclassification cost. The quadratic optimization problem in () can be represented by a dual Lagrange problem and then solved:

$$\max_{\alpha_i} \left\{ \sum_{i=1}^p \alpha_i - \frac{1}{2} \sum_{i=1}^p \sum_{j=1}^p \alpha_i \alpha_j \Phi(x_i) \cdot \Phi(x_j) \right\} \text{ s.t. } \sum_{i=1}^p y_i \alpha_i = 0, \quad 0 \leq \alpha_i \leq C, \quad i = 1, \dots, p$$

α_i are the Lagrange multipliers also satisfying the Karush-Kuhn-Tucker (KKT) conditions (see Appendix). Thanks to another one of the strenghts of SVM -

kernel representation - we don't need to explicitly know the mapping function $\Phi(x)$, but by applying a kernel function (i.e. $K(x_i, x_j) = \Phi(x_i) \cdot \Phi(x_j)$), we can rewrite () as:

$$\max_{\alpha_i} \left\{ \sum_{i=1}^p \alpha_i - \frac{1}{2} \sum_{i=1}^p \sum_{j=1}^p \alpha_i \alpha_j K(x_i, x_j) \right\} \text{ s.t. } \sum_{i=1}^p y_i \alpha_i = 0, \quad 0 \leq \alpha_i \leq C, \quad i = 1, \dots, p$$

The solution then gives us $w = \sum_{i=1}^p \alpha_i y_i \phi(x_i)$ for the optimal values of α_i and w , while b is determined from KKT. The data points that have α_i different than zero are called the support vectors. Thus, the decision function can be written was:

$$f(x) = \text{sign}(w \cdot \Phi(x) + b) = \text{sign}\left(\sum_{i=1}^p \alpha_i y_i K(x_i, x_j) + b\right)$$

Cost Sensitive learning

The regular SVM model has been effectively implemented when the dataset used has balanced classes, however it fails to produce good results when applied on imbalanced data (Batuwita and Palade 2013). When trained on a dataset with extremely imbalanced classes, the SVM framework could produce so skewed hyperplanes that all observations are recognized as the majority class (Akbani, Kwek, and Japkowicz 2004, @veropoulos1999controlling). This is due to the fact that when we take the soft margin optimization problem, we try to maximize the margin and minimize the penalty for the misclassifications. As we consider a constant C for all training examples, the minimization of the penalty is achieved through the minimization of all misclassifications. However, when the used dataset suffers from imbalanced classes, the majority class density would be higher than the minority class density, even when considering the class boundary region (through which the ideal hyperplane would pass).

Thus, we consider here the application of the Different Error Costs (DEC) variation of the SVM algorithm proposed by Veropoulos et al. (1999). The DEC method introduces different misclassification costs - C^+ for the minority and C^- for the majority class. With the inclusion of the higher misclassification cost for the minority class observations, the imbalanced class effect could be brought down. The soft margin optimization problem then has the following form:

$$\begin{aligned} \min \left(\frac{1}{2} w \cdot w + C^+ \sum_{i|y_i=+1}^p \xi_i + C^- \sum_{i|y_i=-1}^p \xi_i \right) \\ \text{s.t. } y_i(w \cdot \Phi(x_i) + b) \geq 1 - \xi_i \end{aligned}$$

$$\xi_i \geq 0, i = 1, \dots, p$$

The Dual Lagrange optimization form is the same as before, with the exception of replacing $0 \leq \alpha_i \leq C$ with $0 \leq \alpha_i^+ \leq C^+$, $0 \leq \alpha_i^- \leq C^-$ for $i = 1, \dots, p$. The α_i^+ and α_i^- are the Lagrange multipliers. The solution of the DEC dual Lagrangian problem follows the same outline as in the normal form.

Kernels

As mentioned before, we used a kernel function ($K(x_i, x_j) = \Phi(x_i) \cdot \Phi(x_j)$) in order to transform the dual Lagrange problem. The advantages of using kernel functions are computational or in some cases it allows for computations that otherwise would be impossible (James et al. (2013)). For the purpose of this study, we are using the linear and radial kernels. The radial kernel shows good performance on non-linear class separation. They have the following representations:

$$Linear : \quad K(x_i, x_j) = \sum_{k=1}^l x_{ik}x_{jk}$$

$$Radial : \quad K(x_i, x_j) = \exp(-\gamma \sum_{k=1}^l (x_{ik} - x_{jk})^2)$$

3.2.4. Tree-Based Methods

Tree-based methods involve segmentation of the feature space into a set of regions and then fitting a simple model to each one. Even though are not too complex conceptually, they are still a very powerful method (J. Friedman, Hastie, and Tibshirani 2001). Firstly, we will give a short overview of a standard classification decision tree (DT) before moving on the methods used in this study.

Let us call the set of non-overlapping regions R_1, R_2, \dots, R_M that are used to divide the feature space. The forms of those regions are high-dimensional rectangles - for simplicity and interpretability. The aim for DT would be to find the boxes that minimize the error term, which in the case of classification can be represented in several ways - misclassification error, Gini index or cross-entropy. However, it is very computationally taxing to consider every feasible partition of the feature space into M boxes. Thus, the recursive binary splitting method is used, which is a top-down, greedy approach - it starts at the top of the tree

and then it splits the feature space, making the best split possible, without looking forward. The algorithm is further described in the Appendix.

However, the beforementioned algorithm can sometimes lead to over-fitting and producing very ineffecient prediction results. Thus, the tree pruning technique is applied - first a large tree is grown, then it is “pruned” and a smaller version is obtain. The procedure leads to reduction in variance at the cost of some bias. Usually the cost complexity pruning algorithm is used in practice - a more thorough description is included in the Appendix.

The regular classification DT has high interpretability, but it sometimes lacks sufficient prediction power - they are often unstable and can be too sensitive to training data (Bhattacharyya et al. 2011). This leads us to variations of the classic classification DT.

Random Forests

A random forest algorithm is an ensemble of classification trees (Breiman 2001). The model starts with growing each tree on separate bootstrapped dataset. Moreover, only a randomly selected feature subset, typically $s \simeq \sqrt{p}$ (Khoshgoftaar, Golawala, and Van Hulse 2007), is used at each individual node. As a result, the entire algorithm is based on two basic, yet powerful, concepts - bagging and random subspace method. To better illustrate the random forest method, the pseudo-algorithm is given below:

Algorithm Random Forest

1. For $b = 1$ to B (B is representing the number of bootstrapped datasets):
 - Draw bootstrapped sample Z^* of size N from the dataset used for training
 - Grow a RF tree T_b on the previously bootstrapped dataset by recursively looping the below given steps for each tree terminal node, up until the n_{min} node (minimum size node) is reached:
 - Select s features randomly from the entire feature subset, p .
 - Choose the best available variable/split-point among the s .
 - Split node into two children nodes.
2. Output the tree ensemble $\{T_b\}_1^B$.
3. Making the prediction for a new point x : If $\hat{C}_b(x)$ is the prediction of the class for the b th RF tree, then $\hat{C}_{rf}^B(x) = \text{majority vote}\{\hat{C}_b(x)\}_1^B$.

The Random Forest algorithm has been quite popular lately, due to its simplicity and performance. It has only two parameters that can be changed - the

number of trees grown and the size of the feature subset used (Breiman 2001). Furthermore, it has often shown superior performance to one of its rival statistic learning method - the SVM (Meyer, Leisch, and Hornik 2003). In the area of financial fraud detection, RF has also shown to be a promising framework (Whitrow et al. 2009, @ngai2011application).

Stochastic Gradient Boosting Machines and eXtreme Gradient Boosting Machines

The ideas introduced by the boosting methodology have been among the most influential in the last twenty years (J. Friedman, Hastie, and Tibshirani 2001). Among the most used and influential boosting algorithms is “AdaBoost.M1” (Freund and Schapire 1997). It was the first developed “adaptive” boosting algorithm, due to its incorporated function to directly adjust its parameters to the used training dataset. This is due to the fact, that the performance of the model was re-evaluated at each iteration - the parameters weights and the final aggregated weights, were re-calculated and improved at each iteration, thus leading to an overall better performance. Moreover, it was found that not only AdaBoost, but boosting algorithms in general, managed to not only reduce variance, but also bias - an improvement over bagging, where only the variance could be decreased.

The following success led to the formulation and development of the gradient-descent based methods, which received the name gradient boosting machines or simply GBM (Freund and Schapire 1997, @friedman2000additive, @friedman2001greedy). The general idea of GBMs is to fit new models iteratively, constructing the base-learners to be maximally correlated with the loss function’s negative gradient, in order to get a better estimate of the response variable (Natekin and Knoll 2013). Moreover, the algorithm offers flexibility in terms of choosing the form of the loss function. In this study, we use the Bernoulli distribution, as we are dealing with a two-class classification problem.

Due to its flexibility and ease of implementation, the GBM algorithm has proven to be a successful model, that offers high predictive power when dealing with machine learning problems (Whiting et al. 2012, @johnson2014learning). Furthermore, when including a random element in the algorithm, as in the Stochastic GBM (SGBM), results tend to improve (Friedman 2002). We describe the pseudo-code of the SGBM for a two-class problem below.

Stochastic Gradient Tree Boosting Machines (Stochastic GBM)

Two-Class Classification as in the **gbm** R package.

1. Initialize $f_0(x) = \arg \min_{\rho} \sum_{i=1}^N \Psi(y_i, \rho)$ and set $g = \{1, 2\}$ for the two classes.
2. For $m = 1$ to M do:
 - Compute negative gradient:

$$z_{im} = -\frac{\partial \Psi(y_i, f(\mathbf{x}_i))}{\partial f(\mathbf{x}_i)} \Big|_{f(\mathbf{x}_i) = \hat{f}(\mathbf{x}_i)}, \quad i = 1, \dots, N$$

- Randomly select $p \times N$ subset.
- Fit regression tree with K number of terminal nodes on the previously selected subset, $g(\mathbf{x}) = E(z|\mathbf{x})$.
- Compute the optimal terminal node predictions, $\rho_{1m}, \dots, \rho_{Km}$, as:

$$\rho_{km} = \arg \min_{\rho} \sum_{\mathbf{x}_i \in S_k} \Psi(y_i, \hat{f}_g(\mathbf{x}_i) + \rho)$$

, where S_k is the set of \mathbf{x} 's that define terminal node k .

- Compute $\hat{f}_m(x) = \sum_{k=1}^K \rho_{km} I(x_i \in \hat{R}_{km})$, as $\{\hat{R}_{km}\}_{k=1}^K$ represents the tree structure.
- Update $\hat{f}^m(\mathbf{x})$ as $\hat{f}^m(\mathbf{x}) \leftarrow \hat{f}_{m-1}(\mathbf{x}) + \lambda \hat{f}_m(x)$.

3. Output $\hat{f}(x) = f_M(x)$

The eXtreme Gradient Boosting algorithm, or shortly XGBoost, was developed as an attempt to improve on the performance of GBM, both in terms of speed and prediction power (Chen and Guestrin 2016). It has not only succeeded in doing so, but it has also become one of the most used models (Nielsen 2016). Among the differences that XGBoost introduces in comparison to the GBM model is that the former uses clever penalization of the individual trees - the leaf weights are not all decreased at the same rate, instead the weights which are estimated without much evidence from the training set will be shrunk more heavily relative to others. Furthermore, the XGBoost employs another type of boosting structure when compared to GBM - Newton boosting. Due to this, the algorithm manages to better learn tree structures, leading to learning better neighbourhoods. XGBoost also incorporates a randomization parameter, which contributes to the individual tree decorrelation and reducing variance. Some technical detail references regarding the differences between XGBoost and GBM are included in the Appendix - such as the gradient and Newton boosting and leaf weight learning. The general pseudo-codes are similar.

3.3. Cross-Validation

4. Data

4.1. Datasets

4.1.1. Real-World Datasets

UCSD-FICO Competition

Université Libre de Bruxelles

4.1.2. Simulated Datasets

PaySim

BankSim

4.2. Problems of Imbalanced Data and Data Sampling Techniques

4.2.1. Problem of Imbalanced Data

One of the biggest challenges faced in detecting fraudulent transactions is the one of unbalanced class sizes, with legitimate class outnumbering vastly the fraudulent one (Bhattacharyya et al. 2011). The application of data-sampling techniques has been widely used in the literature with various results when combined with different algorithms, as when such a problem occurs, it could hinder the model performances (Van Hulse, Khoshgoftaar, and Napolitano 2007). Moreover, in our particular case, the cost of missclassifying the minority class could prove to be a lot more costly than predicting wrongly the majority one.

4.2.2. Data Sampling Techniques

Random Oversampling (ROS) and Random UnderSampling (RUS)

The two techniques are the simplest and most common (Van Hulse, Khoshgof-taar, and Napolitano 2007). In minority oversampling (ROS), the observations from the minority group are randomly duplicated in order to balance the dataset. In majority undersampling (RUS), the aim is the same, but it is achieved by randomly removing observations of the majority class.

SMOTE

The Synthetic Minority Oversampling Technique (SMOTE), proposed by Chawla et al. (Chawla et al. 2002), artificial minority instances are created not simply through duplication, but rather with the extrapolation between preexisting observations. The technique starts by taking into account the k nearest neighbourhoods to a minority observation for every instance from that class. Then, the artificial observation are created, taking into account just a part of the nearest neighbours or all of them (with respect to the desired oversampling specification).

5. Results

6. Further Improvements

7. Conclusion

8. Appendix

GBM vs XGBoost

Structure learning and the gradient vs Newton algorithms

On each iteration, the optimization criterias, concerning tree structure learning, of the Newton tree boosting and the Gradient tree boosting differ. When talking

about Gradient tree boosting, we are interested in the learning of the tree, that exhibits the highest correlation with the negative gradient of the current empirical risk. The tree model is fit using:

$$\{\rho_{kg}, R_{kg}\}_{k=1}^K = \arg \min_{\{\rho_{kg}, R_{kg}\}_{k=1}^K} \sum_{i=1}^n \frac{1}{2} [z_{ig} - \sum]$$

9. References

- Akbani, Rehan, Stephen Kwek, and Nathalie Japkowicz. 2004. “Applying Support Vector Machines to Imbalanced Datasets.” In *European Conference on Machine Learning*, 39–50. Springer.
- Batuwita, Rukshan, and Vasile Palade. 2013. “Class Imbalance Learning Methods for Support Vector Machines.”
- Bhattacharyya, Siddhartha, Sanjeev Jha, Kurian Tharakunnel, and J Christopher Westland. 2011. “Data Mining for Credit Card Fraud: A Comparative Study.” *Decision Support Systems* 50 (3). Elsevier:602–13.
- Breiman, Leo. 2001. “Random Forests.” *Machine Learning* 45 (1). Springer:5–32.
- Chawla, Nitesh V, Kevin W Bowyer, Lawrence O Hall, and W Philip Kegelmeyer. 2002. “SMOTE: Synthetic Minority over-Sampling Technique.” *Journal of Artificial Intelligence Research* 16:321–57.
- Chen, Tianqi, and Carlos Guestrin. 2016. “Xgboost: A Scalable Tree Boosting System.” In *Proceedings of the 22nd Acm Sigkdd International Conference on Knowledge Discovery and Data Mining*, 785–94. ACM.
- Cortes, Corinna, and Vladimir Vapnik. 1995. “Support-Vector Networks.” *Machine Learning* 20 (3). Springer:273–97.
- Freund, Yoav, and Robert E Schapire. 1997. “A Decision-Theoretic Generalization of on-Line Learning and an Application to Boosting.” *Journal of Computer and System Sciences* 55 (1). Elsevier:119–39.
- Friedman, Jerome H. 2001. “Greedy Function Approximation: A Gradient Boosting Machine.” *Annals of Statistics*. JSTOR, 1189–1232.
- . 2002. “Stochastic Gradient Boosting.” *Computational Statistics & Data Analysis* 38 (4). Elsevier:367–78.

- Friedman, Jerome, Trevor Hastie, and Robert Tibshirani. 2001. *The Elements of Statistical Learning*. Vol. 1. Springer series in statistics New York.
- Friedman, Jerome, Trevor Hastie, Robert Tibshirani, and others. 2000. "Additive Logistic Regression: A Statistical View of Boosting (with Discussion and a Rejoinder by the Authors)." *The Annals of Statistics* 28 (2). Institute of Mathematical Statistics:337–407.
- James, Gareth, Daniela Witten, Trevor Hastie, and Robert Tibshirani. 2013. *An Introduction to Statistical Learning*. Vol. 112. Springer.
- Johnson, Rie, and Tong Zhang. 2014. "Learning Nonlinear Functions Using Regularized Greedy Forest." *IEEE Transactions on Pattern Analysis and Machine Intelligence* 36 (5). IEEE:942–54.
- Khoshgoftaar, Taghi M, Moiz Golawala, and Jason Van Hulse. 2007. "An Empirical Study of Learning from Imbalanced Data Using Random Forest." In *Tools with Artificial Intelligence, 2007. ICTAI 2007. 19th Ieee International Conference on*. Vol. 2. IEEE.
- Meyer, David, Friedrich Leisch, and Kurt Hornik. 2003. "The Support Vector Machine Under Test." *Neurocomputing* 55 (1-2). Elsevier:169–86.
- Natekin, Alexey, and Alois Knoll. 2013. "Gradient Boosting Machines, a Tutorial." *Frontiers in Neurorobotics* 7. Frontiers:21.
- Ngai, EWT, Yong Hu, YH Wong, Yijun Chen, and Xin Sun. 2011. "The Application of Data Mining Techniques in Financial Fraud Detection: A Classification Framework and an Academic Review of Literature." *Decision Support Systems* 50 (3). Elsevier:559–69.
- Nielsen, Didrik. 2016. "Tree Boosting with Xgboost-Why Does Xgboost Win" Every" Machine Learning Competition?" Master's thesis, NTNU.
- Van Hulse, Jason, Taghi M Khoshgoftaar, and Amri Napolitano. 2007. "Experimental Perspectives on Learning from Imbalanced Data." In *Proceedings of the 24th International Conference on Machine Learning*, 935–42. ACM.
- Veropoulos, Konstantinos, Colin Campbell, Nello Cristianini, and others. 1999. "Controlling the Sensitivity of Support Vector Machines." In *Proceedings of the International Joint Conference on Ai*, 55:60.
- Whiting, David G, James V Hansen, James B McDonald, Conan Albrecht, and W Steve Albrecht. 2012. "Machine Learning Methods for Detecting Patterns

of Management Fraud.” *Computational Intelligence* 28 (4). Wiley Online Library:505–27.

Whitrow, Christopher, David J Hand, Piotr Juszczak, D Weston, and Niall M Adams. 2009. “Transaction Aggregation as a Strategy for Credit Card Fraud Detection.” *Data Mining and Knowledge Discovery* 18 (1). Springer:30–55.