



Example Dependent cost-sensitive logistic regression

Fraud Analytics: Assignment 3

Jatin Sharma
CS17BTECH11020

Anjani Kumar
CS17BTECH11002

Dhananjay Raut
CS17BTECH11014

Vijay Tadikamalla
CS17BTECH11040

(Teamsize of 4 with special permission from sir)

Objective

Several real-world classification problems are example-dependent cost-sensitive in nature, where the costs due to misclassification vary between examples – e.g. Credit Scoring. The objective of this assignment is to incorporate this requirement in real-world cost agnostic ML models and analyze the trade-off between classification accuracy and financial cost.

Introduction

Standard cost-insensitive binary classification algorithms, like Logistic Regression, Decision Trees, are often used in practice for real-world classification problems like Credit Scoring. The objective in credit Scoring is to classify when the target customer is likely to default a financial contract based on past financial experience. However, in the financial world, the cost associated with approving a potential defaulter varies differently and is quite different from falsely denying a good customer. Some authors have proposed methods that include the miss-classification cost [1]-[2], but assuming a constant misclassification cost is a major drawback.

We implement a framework where misclassification cost varies across examples, i.e. cost-sensitive example dependent classification. We then implement cost-sensitive logistic regression, by changing the objective function of the model to one that is cost-sensitive. We then evaluate the model

with vanilla logistic regression and analyze the trade-off between classification accuracy and incurred financial cost. The results will show that the enhanced model will out-perform the base model and reduces the financial cost by a huge factor.

Logistic regression

Logistic regression is a classification model that, in the specific context of binary classification, estimates the posterior probability of the positive class, as the logistic sigmoid of a linear function of the feature vector [Bishop, 2006]. The probability for the true class is estimated as

$$\hat{p} = P(y = 1|x) = h_{\theta}(x) = g(\theta^T x)$$

where $g(\cdot)$ is the sigmoid function defined as

$$g(z) = \frac{1}{(1 + \exp(-z))}$$

Logistic regression tries to solve the classification problem by minimizing the cost

function $J(\theta)$ which is the negative logarithmic of the likelihood, such that

$$J(\theta) = \frac{1}{N} \sum_{i=1}^N J_i(\theta)$$

where

$$J_i(\theta) = -y_i \log(h_\theta(x_i)) - (1 - y_i) \log(1 - h_\theta(x_i))$$

There are several methods to estimate the logistic regression, in particular, MLE and coordinate descent used in practice, utilizes the convexity of the logistic loss.

However, this cost function assigns the same weight to different errors (namely false positives and false negatives), which is not the case of real-world classification problems.

More specifically, the above cost function assigns 0 costs for true positives and true negatives and inf. for false positives and false negatives.

COST-SENSITIVE LOSS

We design a new cost function that is cost-sensitive and example dependent, by merging the various costs as follows

$$J^c(\theta) = \frac{1}{N} \sum_{i=1}^N \left(y_i (h_\theta(x_i) C_{TP_i} + (1 - h_\theta(x_i)) C_{FN_i}) + (1 - y_i) (h_\theta(x_i) C_{FP_i} + (1 - h_\theta(x_i)) C_{TN_i}) \right).$$

Here we take into account example dependent costs associated, mainly C_{TP} , C_{TN} , C_{FP} , C_{FN} .

Experimentation

Dataset Summary

The dataset consists of customers' records such as illegal/legal transactions, whether the customer has filed within time limits, average tax per month, etc, and costs associated with each example.

Statistics

The dataset contains about 140000 samples, each one with 12 features and the class label. The proportion of default or positive examples is 29.85%

Data partitioning

About 1/3 of the dataset was used for testing and the rest was used for training purposes, sampled based on class label distribution, to avoid undersampling in training examples.

Results

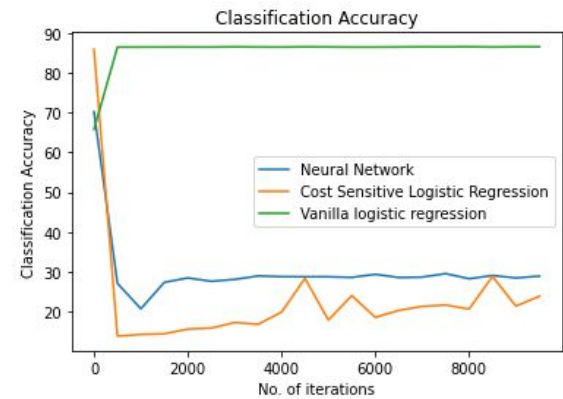
We compare our results with the base logistic regression model which uses cost-insensitive loss function and with cost-sensitive example dependent neural network and analyze the trade-offs.

Cost independent logistic loss



As we can see from the above plot, cost-sensitive models perform worse than the cost agnostic logistic regression. This is because of the fact, that logistic regression is minimizing the negative log-likelihood loss, while other models are minimizing the cost dependent loss function.

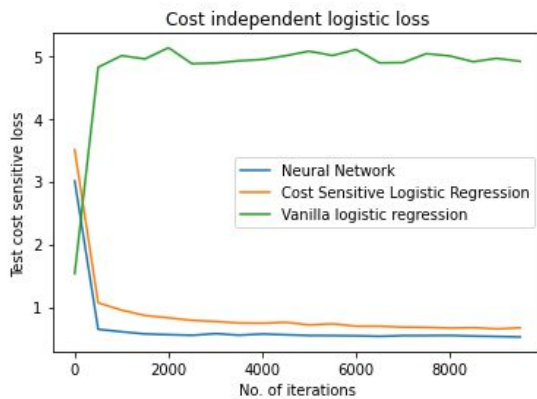
Classification Accuracy



We will see that even though cost-sensitive models compromises the classification accuracy, the cost dependent loss

minimization far exceeds the base logistic regression.

Cost-Sensitive Loss



The above plot shows the tradeoff between classification accuracy and financial loss. Even though the cost-sensitive models compromise the classification accuracy, they outperform the base model by a huge factor (~10) in minimizing the cost sensitive loss, which is the actual goal.

Conclusion

The cost-sensitive example dependent models, clearly, outperforms other models when the goal is to minimize the real-world loss. We also saw the tradeoff between the classification accuracy and loss, i.e. even though a model's classification accuracy is very high, it can still lead to a huge loss in real-world scenarios.

References

- <https://towardsdatascience.com/fraud-detection-with-cost-sensitive-machine-learning-24b8760d35d9>
- A. C. Bahnsen, D. Aouada and B. Ottersten, "Example-Dependent Cost-Sensitive Logistic Regression for Credit Scoring," 2014 13th International Conference on Machine Learning and Applications, Detroit, MI, 2014, pp. 263-269, doi: 10.1109/ICMLA.2014.48.
- C. M. Bishop, Pattern Recognition and Machine Learning, ser. Information science and statistics. Springer, 2006, vol. 4, no. 4.
- <https://towardsdatascience.com/logistic-regression-detailed-overview-46c4da4303bc>