

Approximating the optimal threshold for an  
abstaining classifier based on a reward function  
with regression

*BACHELOR THESIS*

Jonas Fassbender

jonas@fassbender.dev  
11117674

In the course of studies  
*COMPUTER SCIENCE*

For the degree of  
*BACHELOR OF SCIENCE*

Technical University of Cologne  
Faculty of Computer Science and Engineering

First supervisor: Prof. Dr. Heinrich Klocke  
Technical University of Cologne

Second supervisor: Prof. Dr. Fotios Giannakopoulos  
Technical University of Cologne

Overath, July 2019

## 1. Introduction

An abstaining classifier (see e.g. Vanderlooy et al., 2009)—also called a classifier with reject option (see e.g. Fischer et al., 2016)—is a kind of confidence predictor. It can refuse from making a prediction if its confidence in the prediction is not high enough. High enough, in this context, means that the confidence is greater than a certain—hopefully optimal—threshold. Optimality is dependent on a performance metric set beforehand.

This thesis introduces a new kind of method for approximating the optimal threshold based on a reward function—better known from reinforcement learning than from the supervised learning setting (see e.g. Sutton and Barto, 2018). The method treats the reward function as unknown, making it a very general approach and giving quite the amount of freedom in designing the reward function.

In supervised learning the concept that is closest to a reward function is a cost function and many abstract types of cost in supervised learning are known (see Turney, 2002).

Probably today’s most used methods for obtaining the optimal threshold for reducing the expected cost of an abstaining classifier are based on the receiver operating characteristic (ROC) rule (see Tortorella, 2000; Pietraszek, 2005; Vanderlooy et al., 2009; Guan et al., 2018).

The method presented in this thesis is more flexible than the methods based on the ROC rule and can—depending on the context of the classification problem—produce results better interpretable than results from a cost setting (see Chapter 2). Also it is more natural with multi-class classification problems than the methods based on the ROC rule, all assuming binary classification problems, wherefore the classifiers generated by these methods must be transformed to multi-class classifiers for non-binary problems.

On the other hand the presented method can suffer from its very general approach and only produces approximations. This can result in non-optimal and unstable thresholds.

This thesis first presents a motivational example. In Chapter 3 the proposed method is presented. After that experiments on data sets from the UCI machine learning repository (see Dua and Graff, 2017) are discussed. At last further research ideas are listed and a conclusion is drawn.

## 2. Motivational example

This chapter will point out the usefulness of abstaining classifiers in real world application domains where reliability is key. It will show an example why the reward setting can improve readability in some domains. First another example, for which the cost setting—more commonly used in supervised learning—comes more natural is given and the differences are discussed.

Abstaining classifiers—compared to typical classifiers, which classify every prediction, maybe even without a confidence value in it (then called a bare prediction)—can be easily integrated into and enhance processes where they partially replace some of the decision

making, since they can delegate the abstained predictions back to the underlying process. The use of abstaining classifiers in domains where reliability—in regard to prediction errors—is important, has an interesting aspect in giving reliability while still being able to decrease work, cost, etc. to some degree. This is a valuable property if there does not exist a typical classifier good enough to fully replace the underlying process.

Many real world application domains for abstaining classifiers can express a cost function associated to the decisions about predicting and abstaining of the classifier—which then chooses the threshold with which it produces the least amount of cost, therefore minimizing the cost of introducing the abstaining classifier to the process.

For example, the real world application domain could be a facial recognition system at a company which regulates which employee can enter a trust zone and which can not. The process which should be enhanced with the facial recognition system is a manual process where the employee has to fill out a form in order to receive a key which opens the trust zone.

In this example, the costs of miss-classifying an unauthorized person as authorized can be huge for the company while abstaining or classifying an authorized employee as unauthorized produces quite low costs—the authorized employee just has to start the manual process, which should be replaced by the facial recognition system.

On the other hand, for some real world application domains a reward function based on which the abstaining classifier chooses the threshold by maximizing the reward—rather than minimizing the cost—comes more natural.

Such a domain would be the finance industry, where we often can associate a certain amount of money an abstaining classifier can produce or save by supporting the decision making of an underlying process.

An example for such a process would be the process of a bank for granting a consumer credit. The bank requests information about the consumer from a credit bureau in order to assess the consumer’s credit default risk. Now the bank wants to predict the consumer’s credit default risk based on information the bank has about the consumer. If the credit default risk is very high or very low the bank can save money not making a request to the credit bureau for this consumer. The optimal threshold for the abstaining classifier making the prediction about the credit default risk can easily be expressed by a reward function. Every correct decision saves the bank the money the request to the credit bureau costs. Every miss-classification costs the bank either the amount of money it would gain by granting the credit, or the money it loses by giving a credit to somebody that does not pay the rates. Abstention cost is the cost of making a request to the credit bureau.

Using a reward function—like in the example above—instead of a cost function has an advantage in readability. One can easily assess the gain of introducing the abstaining classifier to the process. Is the reward generated by the abstaining classifier higher than zero, the process is enhanced by the abstaining classifier. Otherwise the abstaining classifier would produce more cost than it would save and it is not valuable for the bank to introduce it to its process of assessing a consumer’s credit default risk.

### 3. Proposed method based on reward

Let  $\mathbf{X}$  be our observation space and  $\mathbf{Y}$  our label space.  $|\mathbf{Y}| < \infty$  since only classification is discussed. Let  $\mathbf{Z}$  be the cartesian product of  $\mathbf{X}$  and  $\mathbf{Y}$ :  $\mathbf{Z} := \mathbf{X} \times \mathbf{Y}$ .  $\mathbf{Z}$  is called the example space. Let an example  $z_i$  from  $\mathbf{Z}$  be:  $z_i := (x_i, y_i); z_i \in \mathbf{Z}$ . A data set<sup>1</sup> containing examples  $z_1, \dots, z_n$  is annotated as  $\{z_1, \dots, z_n\}$ .

A classical machine learning predictor—in the previous chapter called a typical classifier—can be represented as a function:

$$D : \mathbf{Z}^* \times \mathbf{X} \rightarrow \mathbf{Y}. \quad (1)$$

Its first argument being a data set with an arbitrary length the classifier is trained on, while the second is an observation which should be predicted (mapped to a label from  $\mathbf{Y}$ ).

Let  $D_{\{z_1, \dots, z_n\}}$  be a classical machine learning predictor trained on the data set  $\{z_1, \dots, z_n\}$  and let  $D_{\{z_1, \dots, z_n\}}(x)$  be equivalent to (1).

The proposed method relies on scoring classifiers. A scoring classifier does not return just a label but instead returns some score for each label from our label space. The only constraint on the scores is that higher scores are better than lower. A score could be a probability or just an uncalibrated confidence value (see Vanderlooy et al., 2009).

Let  $S$  be a scoring classifier:

$$S : \mathbf{Z}^* \times \mathbf{X} \rightarrow (\mathbf{Y} \rightarrow \mathbb{R}).$$

$S$  takes the same arguments as (1) but instead of producing bare predictions it returns a function which maps every label from the label space to a score determined by  $S$ .

The method proposed is only interested in the highest score and the associated label. For that two functions  $k$  and  $v$  are defined:

$$\begin{aligned} k(S_{\{z_1, \dots, z_n\}}, x) &= \arg \max_{y \in \mathbf{Y}} S_{\{z_1, \dots, z_n\}}(x)(y) \\ v(S_{\{z_1, \dots, z_n\}}, x) &= \max_{y \in \mathbf{Y}} S_{\{z_1, \dots, z_n\}}(x)(y). \end{aligned}$$

The composition  $kv$  of  $k$  and  $v$  returns the tuple with the label mapped to the highest score:

$$kv(S_{\{z_1, \dots, z_n\}}, x) = (k(S_{\{z_1, \dots, z_n\}}, x), v(S_{\{z_1, \dots, z_n\}}, x)).$$

The proposed method relies on an meta architecture comparable to the one described in Smirnov et al. (2009).

---

1. not an actual set but a multi-set since it can contain the same element more often than one time.

## 4. Experiments

## 5. Further research

## 6. Conclusion

## Appendix

### A. Plots

### References

- Dheeru Dua and Casey Graff. UCI machine learning repository, 2017. URL <http://archive.ics.uci.edu/ml>.
- Lydia Fischer, Barbara Hammer, and Heiko Wersing. Optimal local rejection for classifiers. *Neurocomputing*, 214:445 – 457, 2016. ISSN 0925-2312. doi: <https://doi.org/10.1016/j.neucom.2016.06.038>. URL <http://www.sciencedirect.com/science/article/pii/S0925231216306762>.
- Hongjiao Guan, Yingtao Zhang, Heng-Da Cheng, and Xianglong Tang. Abstaining classification when error costs are unequal and unknown. *CoRR*, abs/1806.03445, 2018. URL <http://arxiv.org/abs/1806.03445>.
- Tadeusz Pietraszek. Optimizing abstaining classifiers using roc analysis. In *Proceedings of the 22Nd International Conference on Machine Learning, ICML '05*, pages 665–672, New York, NY, USA, 2005. ACM. ISBN 1-59593-180-5. doi: 10.1145/1102351.1102435. URL <http://doi.acm.org/10.1145/1102351.1102435>.
- Evgueni Smirnov, Georgi Nalbantovi, and A. M. Kaptein. Meta-conformity approach to reliable classification. *Intelligent Data Analysis*, 13, 01 2009. doi: 10.3233/IDA-2009-0400.
- Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. MIT Press, Cambridge, MA, USA, 2nd edition, 2018.
- Francesco Tortorella. An optimal reject rule for binary classifiers. pages 611–620, 08 2000. doi: 10.1007/3-540-44522-6\_63.
- Peter D. Turney. Types of cost in inductive concept learning. *CoRR*, cs.LG/0212034, 2002. URL <http://arxiv.org/abs/cs.LG/0212034>.
- Stijn Vanderlooy, Ida G. Sprinkhuizen-Kuyper, Evgueni N. Smirnov, and H. Jaap van den Herik. The roc isometrics approach to construct reliable classifiers. *Intell. Data Anal.*, 13(1):3–37, January 2009. ISSN 1088-467X. URL <http://dl.acm.org/citation.cfm?id=1551758.1551760>.

## **Erklärung**

Ich versichere, die von mir vorgelegte Arbeit selbstständig verfasst zu haben. Alle Stellen, die wörtlich oder sinngemäß aus veröffentlichten oder nicht veröffentlichten Arbeiten anderer oder der Verfasserin/des Verfassers selbst entnommen sind, habe ich als entnommen kenntlich gemacht. Sämtliche Quellen und Hilfsmittel, die ich für die Arbeit benutzt habe, sind angegeben. Die Arbeit hat mit gleichem Inhalt bzw. in wesentlichen Teilen noch keiner anderen Prüfungsbehörde vorgelegen.

---

Ort, Datum

---

Rechtsverbindliche Unterschrift