

A systematic comparison between different Base Learners in AdaBoosting model

Mateus Maia ¹, Anderson Ara ²

Resumo: Os métodos de boosting estão se tornando cada vez mais populares devido ao seu excelente desempenho quando comparados com outras técnicas de aprendizado estatístico. O Adaptive Boosting, ou simplesmente AdaBoost, foi uma das primeiras técnicas de aprimoramento desenvolvidas e consiste, geralmente, em uma combinação linear de modelos fracos (modelos que se saem um pouco melhor do que um palpite aleatório) para construir um classificador forte. Este trabalho propõe uma comparação sistemática entre os possíveis modelos que podem ser utilizados como classificadores base, além da árvore de decisão que é o método padrão para o AdaBoosting, utilizando os outros métodos comumente usados na área do aprendizado estatístico de máquina, sendo estes: K Vizinhos Mais Próximos, Análise Linear Discriminante, Regressão Logística, Redes Neurais e Máquinas de Vetores de Suporte. Todas possíveis formas do AdaBoosting foram aplicadas à dez bases de dados diferentes, e foram comparadas através da acurácia utilizando a técnica de Holdout Repetido, com trinta repetições. Como resultado, apesar da performance média geral das árvores de decisões ainda se manter superior, esta foi muito próxima às outras técnicas, como o MVS, além de que em metade das bases de dados a Árvore de Decisão teve performance inferior quando comparada à outros métodos. Dessa forma, concluiu-se que, apesar de estabelecido como modelo de base padrão do AdaBoost é importante avaliar a aplicação de outros classificadores além da árvore de decisão a depender do problema investigado.

Palavras-chave: AdaBoosting; Ensemble Models; Statistical Learning; Machine Learning.

Abstract: *Boosting methods are becoming popular due their outstanding performance when compared with some others statistical learning techniques. The Adaptive Boosting consists in a linear combination of weak models to built a strong classifier. This work proposes a systematic comparison between the possible models that can be used as base classifiers.*

Keywords: AdaBoosting; Ensemble Models; Statistical Learning; Machine Learning.

Introduction

The Adaptive Boosting algorithm (AdaBoost), developed by Freund and Schapire (1997), has been showed to be a great statistical model that can outperform a lot of others statistical learning algorithms. Like all others ensemble methods, the AdaBoosting is built by the combination of several models that vote to classify and predict an observation. In AdaBoosting, these classifiers are modeled sequentially and each new model it's weighted

¹Universidade Federal da Bahia. e-mail: mateusmaia11@gmail.com

²Universidade Federal da Bahia. e-mail: anderson.ara@ufba.br

considering the capacity to predict correctly the previous missclassified observations. Generally, the base models are the decision-tree algorithm (C4.5) (Quilan,1993), however any other weak learner can be used in AdaBoosting. In order to explore the capacity of use a variety of base learners, this work present a complete comparisson between the most commons models used in the statistical learning tasks in their simplest form. The definition of a weak learner can be set as the model that classifies observations a little better than random guessing. Also, in the most cases, those weak models are associated with linear classification rules or decision boundaries. Many researches focus on enhancing the performance of AdaBoost, by choosing more discriminative classifier (Ratsch, 2001; Schwenk and Bengio, 2000; Li et al., 2008), so change the base learner it's way to emphasize this aspect.

The AdaBoost algorithm was modified to use the following models: K Nearest Neighbors, Linear Discriminant Analysis, Logistic Regression, Neural Networks, Support Vector Machines. All of them were applied in 10 datasets, and their accuracys were evaluated using a repeated holdout validation technique.

Methodology

Essentially, boosting consists of repeatedly using a base weak learning algorithm, on differently weighted versions of the training data, yielding a sequence of weak classifiers that are combined in a addition function. The weighting of each model depends on the accuracy of the previous, in order to increase the importance of classify correctly wrong predicted observations from the last model. The ensemble prediction function of AdaBoost $H : X \rightarrow \{-1, 1\}$ is given by

$$H(\mathbf{x}) = \text{sign} \left(\sum_{m=1}^M \alpha_m h_m(\mathbf{x}) \right) \quad (1)$$

where $\alpha_1, \dots, \alpha_M$ is a set of weights from respective h_1, \dots, h_M models.

To build this model, we followed the pseudo-code below, varying the base models h_i by those mentioned before

- Given (x_1, y_1) , where $x_i \in X, y_i \in \{-1, 1\}$
- Initialize: $D_1 = \frac{1}{n}$ for $i = 1, \dots, n$
- For $m = 1, \dots, M$
 - Train the weak learner using distribution D_m
 - Get the hypothesis $h_m : X \rightarrow \{-1, 1\}$
 - Aim: Select h_m with lower weighted error.

$$\epsilon_m = \text{Pr}_i \sim D_t[h_m(x_i) \neq y_i]$$

- Choose $\alpha_m = \frac{1}{2} \ln \left(\frac{1-\epsilon_m}{\epsilon_m} \right)$
- Update for $i = 1, \dots, n$

$$D_{m+1} = \frac{D_m(i) \exp(-\alpha_m y_i h_m)}{Z_m}$$

Where Z_m is a normalization factor.

Then the output is given by the Equation (1).

Were choosen six type of statistical models to use as base learners in AdaBoosting, which follows:

- **K Nearest Neighbors (KNN)**, with the parameter k defined by tuning.
- **Linear Discriminant Analysis**.
- **Logistic Regression** in canonical form.
- **Neural Networks** with one layer and one perceptron.
- **Support Vector Machines** with the linear kernel.
- **Decision Trees** with just one split node (Stump Models).

Each model was applied to differents datasets (Table 1), that can be acessed in *UCI ML Data Repository*, to evaluate empiracally the peformance from each method. They were all a binary classification task, where $y_i \in \{-1, 1\}$. The validation technique used was the repeated holdout, with 30 repetitions, and the performace metric obtained was the accuracy, once all datasets were balanced.

Table 1: All Datasets

Data Set	# Observations	# Covariates
<i>hepatitis</i>	80	20
<i>parkinsons</i>	195	23
<i>sonar</i>	208	61
<i>heart-statlog</i>	270	14
<i>haberman</i>	306	4
<i>liver-disorder</i>	345	7
<i>ionosphere</i>	351	35
<i>vertebral-column</i>	310	7
<i>heart-h</i>	261	11
<i>magic04</i>	19020	11

Results and Discussion

The main result is represented by the Figure 1, where it's possible to see a boxplot from the accuracy results from each round of the holdout split. To each AdaBoost model where generated 100 models of each classifier type.

As we can observe from Table 2, the best AdaBoost isn't always that which refers to the standard Stump Models, and specifically in the half of the cases he isn't the one with greatest accuracy. The Linear Discriminant Model as well the Logistic Regression performs relatively well in some databases with low dimensionality. However, the LDA for example, can't perform in some datasets where some covariates have an strong colinearity. The base learner that appears the most, together with the Decision Tree Stumpo is the SVM, appearing 5 out of 10 times, followed by the neural networks, LDA, and logistic that appears 4 times each.

Those results confirms the necessity of an investigation when choose the base leanear for an specific dataset. Despite the flebility from Decision Tree Models, which can deal with continous and categorical features, and data with high dimensionality, once a dataset has some simpler characteristics and behaviours, others base classifiers can perform better then the stump models.

Table 2: Table of Accuracy different AdaBoosting to all datasets

Datasets	KNN	LDA	Logistic	NN	SVM	Tree Stump
<i>hepatitis</i>	0.58 \pm 0.11	-	0.64 \pm 0.09	0.66\pm0.08	0.66\pm0.07	0.65 \pm 0.09
<i>parkinsons</i>	0.80 \pm 0.11	-	0.81 \pm 0.06	0.65 \pm 0.24	0.89\pm0.04	0.89\pm0.04
<i>sonar</i>	0.68 \pm 0.16	0.73 \pm 0.04	0.71 \pm 0.05	0.83\pm0.04	0.78 \pm 0.05	0.83\pm0.04
<i>heart-statlog</i>	0.65 \pm 0.16	0.84\pm0.04	0.84 \pm0.04	0.79 \pm 0.05	0.83 \pm 0.04	0.83 \pm 0.04
<i>haberman</i>	0.73 \pm 0.04	0.75\pm0.03	0.75 \pm0.03	0.75\pm0.03	0.75\pm0.03	0.73 \pm 0.03
<i>liver-disorder</i>	0.66 \pm 0.05	0.70 \pm 0.05	0.69 \pm 0.04	0.71 \pm 0.04	0.68 \pm 0.04	0.73\pm0.04
<i>ionosphere</i>	0.88 \pm 0.04	-	0.86 \pm 0.05	0.90 \pm 0.03	0.87 \pm 0.03	0.93\pm0.02
<i>vertebral-column</i>	0.80 \pm 0.04	0.84 \pm0.04	0.84\pm0.04	0.81 \pm 0.04	0.84\pm0.04	0.82 \pm 0.03
<i>heart-h</i>	0.64 \pm 0.05	0.80 \pm0.04	0.80\pm0.04	0.76 \pm 0.04	0.80\pm0.03	0.79 \pm 0.03
<i>magic04</i>	0.81 \pm 0.01	0.79 \pm 0.01	0.79 \pm 0.01	0.83 \pm 0.01	0.79 \pm 0.01	0.83\pm0.01
All Datasets	0.72 \pm 0.10	0.78 \pm 0.06	0.78 \pm 0.08	0.77 \pm 0.11	0.79 \pm 0.08	0.81\pm0.09

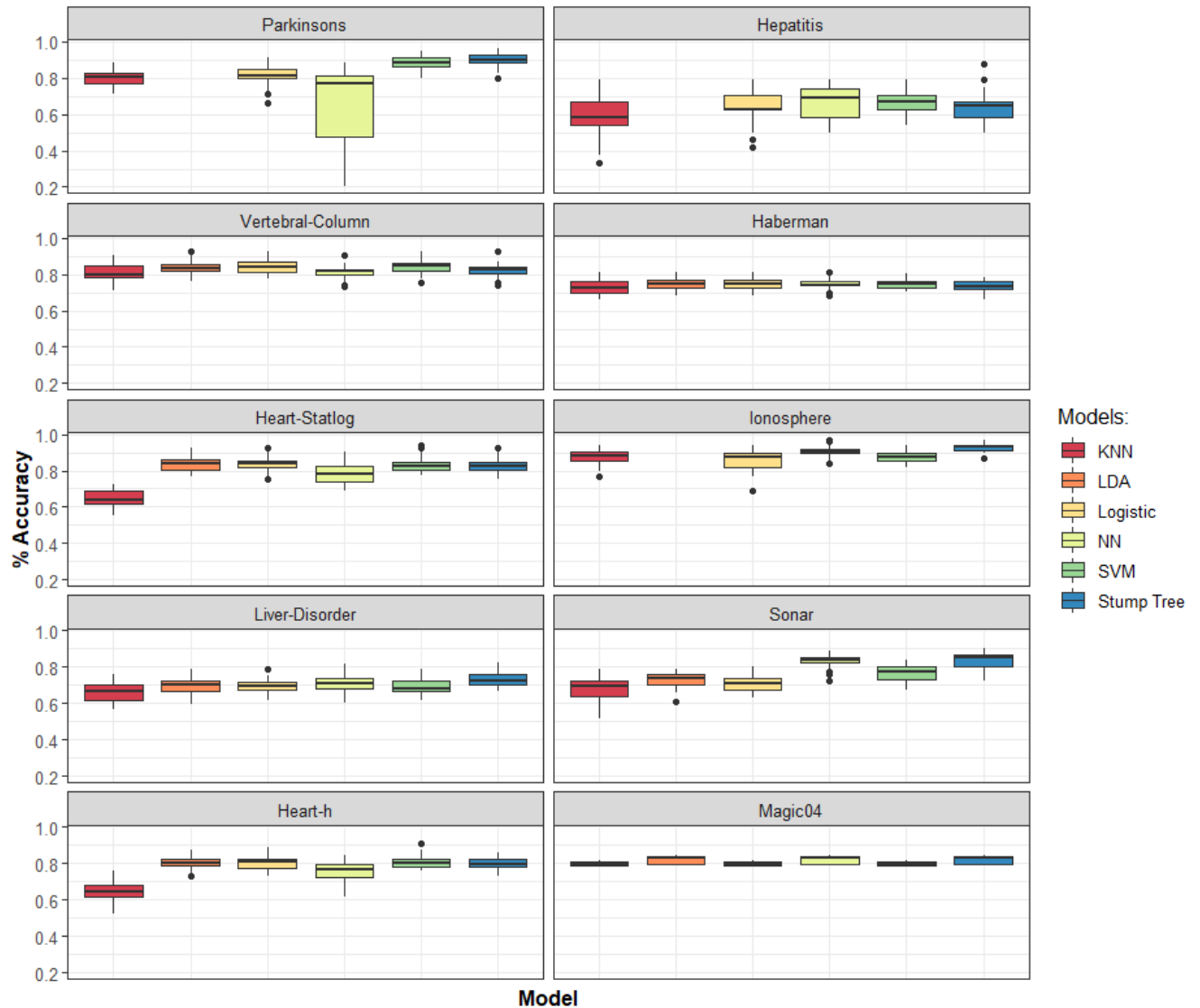


Figure 1: AdaBoosting accuracy to all datasets and base learners.

Conclusion

The AdaBost can be defined as a powerful ensemble classifier formed by successively refitting a weak classifier to different weighted realizations of a data set. In this work we proposed a comparisson between different base leaners models used in AdaBoost, instead the standard Decision Tree Stumps in order to study the efficiency of the others methods to predict correctly new observtions, and create more discriminant classifiers to compose the AdaBoosting classifiers. We could observe that several methods improved or equated the standard AdaBoosting suggesting that's interessting to analyze in each situation what could be the better base learner to use. To futures works is important to try to combine multiples learners in a single AdaBoost model, and maybe change the hiperparameters iterativily in each model.

Acknowledgment

I would like to thank for the support given by the CNPq/CAPES during this work.

References

- FREUND, Yoav; SCHAPIRE, Robert E. A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of computer and system sciences*, v. 55, n. 1, p. 119-139, 1997.
- SALZBERG, Steven L. C4. 5: Programs for machine learning by j. ross quinlan. morgan kaufmann publishers, inc., 1993. *Machine Learning*, v. 16, n. 3, p. 235-240, 199.
- Rätsch, G., Onoda, T., & Müller, K. R. (2001). Soft margins for AdaBoost. *Machine learning*, 42(3), 287-320.
- Schwenk, Holger, and Yoshua Bengio. "Boosting neural networks." *Neural computation* 12.8 (2000): 1869-1887.
- Li, X., Wang, L., & Sung, E. (2008). AdaBoost with SVM-based component classifiers. *Engineering Applications of Artificial Intelligence*, 21(5), 785-795.
- Hadley Wickham (2017). tidyverse: Easily Install and Load the 'Tidyverse'. R package version 1.2.1. <https://CRAN.R-project.org/package=tidyverse>
- Terry Therneau and Beth Atkinson (2018). rpart: Recursive Partitioning and Regression Trees. R package version 4.1-13. <https://CRAN.R-project.org/package=rpart>
- Venables, W. N. & Ripley, B. D. (2002) *Modern Applied Statistics with S*. Fourth Edition. Springer, New York. ISBN 0-387-95457-0
- Alexandros Karatzoglou, Alex Smola, Kurt Hornik, Achim Zeileis (2004). kernlab - An S4 Package for Kernel Methods in R. *Journal of Statistical Software* 11(9), 1-20. URL <http://www.jstatsoft.org/v11/i09/>