

Adaptive random forests for data stream regression

Heitor Murilo Gomes¹, Jean Paul Barddal^{1,2}, Luis Eduardo Boiko², Albert Bifet¹

1- Department of Computer Science and Networks (INFRES), Télécom ParisTech
Université Paris-Saclay, Paris, France

2- Programa de Pós-Graduação em Informática (PPGIa)
Pontifícia Universidade Católica do Paraná, Curitiba, Brazil

Abstract. Data stream mining is a hot topic in the machine learning community that tackles the problem of learning and updating predictive models as new data becomes available over time. Even though several new methods are proposed every year, most focus on the classification task and overlook the regression task. In this paper, we propose an adaptation to the Adaptive Random Forest so that it can handle regression tasks, namely *ARF-Reg*. *ARF-Reg* is empirically evaluated and compared to the state-of-the-art data stream regression algorithms, thus highlighting its applicability in different data stream scenarios.

1 Introduction

Data stream mining is an important topic in the machine learning community. It tackles the problem of learning and updating learning models as new data becomes available over time. Even though several new methods are proposed every year, most focus on the classification task and overlook the regression task. Important examples of regression include, for instance, temperature and precipitation forecasts, stock market and household price predictions. Furthermore, the data distribution of the examples aforementioned may be *ephemeral* in the sense that it can change over time. For instance, the temperature and precipitation rates of a region may change due to unexpected environmental accidents, or the prices of stocks may vertiginously decrease if a company is found to be amidst corruption schemes, and so forth.

In this paper, we adapt the Adaptive Random Forest (ARF) learner presented in [1] to the regression task, hereafter referred to as *ARF-Reg*. *ARF-Reg* was implemented in the Massive Online Analysis (MOA) framework and it will be made publicly available for further studies on the area. The remainder of this paper is divided as follows. Section 2 describes the data stream regression task and its challenges. Section 3 overviews related works. Section 4 describes the proposed method, which is later evaluated in Section 5. Finally, Section 6 concludes this paper and reports envisioned future works.

2 Problem Definition

Despite the impressive amount of effort put on data stream mining, most of the works focus on classification and overlooked both regression and clustering

tasks. In this paper, we focus on the regression task, which aims at predicting a continuous value. Examples of regression include, for instance, temperature and precipitation forecasts, stock market value and household price predictions.

Formally, we assume S to be a data stream providing instances (\vec{x}^t, y^t) , where \vec{x}^t is feature vector, $y^t \in \mathbb{R}$ is the target meta-attribute, and t is the arrival timestamp. In regression, the goal is to iteratively learn a predictive model $h : \vec{x} \rightarrow y$ as new data becomes available. In this work, we assume a test-then-train scheme and such that y^t becomes available right after \vec{x}^t arrives. Even though this assumption might not hold in a variety of scenarios, it is by far the most used in the area.

Finally, one of the most important challenges in data streams is tackling concept drifts [2]. A drift occur when the data distribution changes, i.e., the mapping between features in \vec{x} and the target values y change over time. An important trait of concept drift relates to the rate at which it happens. The rate at which drifts happen can be abrupt, incremental, gradual or reoccurring. Notice that noise or outliers ought not be confused with drift. The difference between noise/outliers and drifts is persistence.

3 Related Works

Similarly to batch learning, the number of techniques developed for classification greatly outnumber those tailored for regression. In this section we report important contributions to the field and that are used in the empirical evaluation of the proposed method.

Regression trees are similar to decision trees as they iteratively perform splits over attributes with the goal of maximizing some goodness-of-fit criterion. **Fast and Incremental Model Trees (FIMT-DD)**, initially presented in [3] are the main example of regression trees for data streams. Similarly to standard Hoeffding Trees [4], FIMT-DD starts with an empty tree that keeps statistics from arriving data until a grace period is reached. At this point, features are ranked according to their variance, and if the two best-ranked differ by at least the Hoeffding Bound [5], the tree branches and the process is repeated. FIMT-DD also encompasses a change detection scheme that periodically flags and adapts subbranches of the tree where significant variance increases are observed. Similarly, **ORTO** also grows trees incrementally with the arrival of instances, yet, it also introduces ‘option’ nodes, which allow an instance to follow all the branches available in a tree node [3].

Regression rules are another relevant representatives of data stream regression. By far, the most used algorithm is **Adaptive Model Rules (AMRules)** [6]. AMRules learns both ordered and unordered rule set from data streams. To detect and adapt to concept drifts, each rule is associated with a Page-Hinkley drift detector [7], which prunes the rule set given changes in the incoming data.

It is also common to combine several learning models by ‘ensembling’ them. One important example for regression is the **Scale-free Network Regression (SFNR)**. SFNR is a light-weighted network-based regression ensemble for data

streams [8]. It arranges its learners in a probabilistic scale-free network such that the most accurate models tend to become more prominent ('hubs'), and thus, have higher weights during the prediction step. SFNR also uses drift detectors to eliminate inaccurate estimators according to drifts in data. Finally, it is important to mention that a previous approach for regression using Random Forests has been introduced in [9], yet, it does not include any methods to handle concept drifts.

4 Adaptive Random Forest for Regression

To describe our proposed method, ARF-Reg, we use the taxonomy presented in [10]. Precisely, we describe ARF-Reg in terms of its **voting** strategy, **diversity** induction, **base learner** characteristics and **update dynamics**.

- **voting** averages the individual predictions to obtain the final prediction;
- **diversity** is induced into the forest by training the trees on different subsets of data and by limiting the split decisions to an m randomly selected subset of features from the original input features. This follows the same methodology applied in [1], which was inspired by [11];
- the **base learner** is a regression tree, namely a FIMT-DD [12]. FIMT-DD is an incremental learner featuring an efficient attribute split and selection method;
- the **update dynamics** in ARF-Reg relies on both internal and external drift detectors for each tree and by growing trees in the background when a warning is detected.

The original ARF classifier achieved best results when its drift and warning methods were set to use the ADaptive WINdow (ADWIN) [13] algorithm. Therefore, in this work, we experiment while using the external drift detection method using ADWIN and its 'moderate' configuration as described in [1]. In ARF-Reg, besides using this external drift detection method we also experiment using the original Page-Hinkley test [7] internally to each FIMT-DD to detect and adapt to them.

In [9] authors present an Online Random Forest version that also uses FIMT-DD as its base learner, despite the aforementioned approach to deal with concept drifts another difference is how we perform resampling in ARF-Reg. Following the same strategy presented in [1], we simulate leveraging bagging ($\lambda = 6$) instead of the standard online bagging ($\lambda = 1$). The practical implications of this decision is that trees are trained with more data, which makes them more likely to split faster, thus adapting faster to drifts and rapidly building deeper trees.

5 Experiments

In these experiments, we analyze how ARF-Reg performs in terms of Root Mean Square Error (RMSE) in different scenarios including both stationary and non-stationary data. To benchmark the results we compare ARF-Reg against several state-of-the-art regression algorithms, including its base learner FIMT-DD. The experiments follow a test-then-train approach, such that each instance is first used for testing and immediately used for training. The configuration of the algorithms follows their original publications. We present experiments with 4 variations of ARF-Reg, all of them with 10 learners, namely:

- ARF-Reg: the default parametrization using ADWIN for both external drift and warning detection. The number of features in the subspace is also set to $m = \sqrt{M} + 1$, where M is the total number of features.
- ARF-Reg-inv: Similar to ARF-Reg, but uses $m = M - \sqrt{M}$, thus trying more features per split.
- ARF-Reg-int: Similar to ARF-Reg, but it disables the external drift detection and warning methods. The adaptation to drifts relies on the internal PHT provided by FIMT-DD. This version closely resembles the ensemble method presented in [9].
- ARF-Reg-int-inv: It is a combination of ARF-Reg-int and ARF-Reg-inv.

Table 1 presents the datasets used during the experiments, which are further discussed below. We used four synthetically generated datasets to represent streams that exhibit incremental, abrupt, gradual and no drifts at all. The first is derived from the Hyperplane generator often used in classification tasks, but that can be adapted to regression problems as well. The **HyperplaneReg** [14] (i.e. Hyperplane for Regression) generator creates a hyperplane, which is a flat, $(M - 1)$ dimensional subset of a M space that divides it into two disjoint parts. Instances values are generated following a uniform distribution, and three different functions that map instances to their outcomes are used: (i) the Euclidian distance between the instance in the feature space and the hyperplane, (ii) the square of this distance, or (iii) the cube of the same distance. Incremental drifts are simulated by slowly changing the hyperplane location. Setting a new hyperplane by varying the random seed one can simulate an abrupt or gradual drift. The final synthetic data is the **FRIED** [15] dataset. FRIED is a classical artificial regression dataset where each instance is represented by 10 features whose values are independently and uniformly distributed over $[0,1]$. The outcome value y is given by an equation that takes as input only 5 of the 10 features.

We use three real-world datasets in our experimentation framework. The first two (HOUSE16N and Ailerons) are not representations of streaming data, but these can be interpreted as stationary streams in our analysis. The goal of the **HOUSE16H**. dataset is to estimate the median house price in a given region according to 16 features representing demographic and house market data. The

Table 1: Datasets overview. *i*, *a* and *g* stands for incremental, abrupt and gradual drifts, respectively.

Experiment	# Instances	# Features	# drifts
HYPER(i)	100,000	10	-
HYPER(a)	100,000	10	2
HYPER(g)	100,000	10	2
FRIED	40,768	10	-
AILERONS	13750	40	-
BIKE	17,389	12	-
HOUSE16H	22,784	16	-

Table 2: RMSE obtained during the experiments.

Experiment	FIMT-DD	ORTO	AMRules	SFNR	ARF-Reg	ARF-Reg-inv	ARF-Reg-int	ARF-Reg-int-inv
HYPER(i)	0.3499	0.8042	0.0333	0.3403	0.5050	0.3915	0.5050	0.3915
HYPER(a)	0.2496	0.3214	0.0790	0.2109	0.2128	0.2001	0.2164	0.2036
HYPER(g)	0.2552	0.3262	0.1060	0.2246	0.2234	0.2117	0.2260	0.2142
FRIED	2.9516	3.4416	2.4802	2.7659	3.1874	2.8676	3.1874	2.8676
AILERONS	0.0003	0.0005	0.0020	0.0003	0.0003	0.0003	0.0003	0.0003
BIKE	114.5981	101.3098	135.0353	108.3146	106.7113	86.5060	100.0952	84.1412
HOUSE16H	43236.1339	71470.0178	46033.5364	42739.1849	41612.9812	39956.4878	41629.3639	40124.0958
Avg. Rank	6.00	7.20	3.60	4.40	4.70	2.30	5.10	2.70

AILERONS dataset contains information about an F16 aircraft, and the goal is to predict the control action to be applied to the ailerons of the aircraft. The **BIKE** [16] dataset includes 2 years (2011 and 2012) worth of a bike-sharing service from Washington D.C., USA. The goal is to predict how many bikes in total will be rented in the next hour using weather and temporal data (e.g., time of the day, the day of the week and so forth). Some features from the original data were removed to avoid data leakage (e.g., **registered** and **casual**), as their sum converges to the target variable (i.e., **cnt**).

5.1 Discussion

Comparing the results obtained across all the experiments, a combination of Friedman and Nemenyi [17] statistical tests show that ARF-Reg-inv is the best performing one. Nevertheless, no significant differences were observed across all learners, except ORTO, which was significantly worse than the others. In this comparison, *ARF-Reg* obtained compelling error rates, mainly on real-world scenarios, thus showing its efficacy. The results from ARF-Reg-inv and ARF-Reg-int-inv in comparison to ARF-Reg can be explained by the small amount of features per dataset. In these cases, using more features per splits tends to achieve better results. When we compare the ARF-Reg variations to AMRules, we can observe that AMRules perform very well on synthetic datasets. However, it could not achieve reasonable results in the real-world datasets. Finally, comparing ARF-Reg variations against its base learner (FIMT-DD), we can see that improvements are consistent, with the exception of HYPER(i).

6 Conclusion

In this paper, we introduced an adaptation of the Adaptive Random Forest method for data stream regression, called *ARF-Reg*. The proposed method was empirically assessed and compared to existing works of the area. In future works, we intend to investigate other combinations techniques, analyze other drift detection techniques and thoroughly evaluate the computational resources.

References

- [1] Heitor M. Gomes, Albert Bifet, Jesse Read, Jean Paul Barddal, Fabrício Enembreck, Bernhard Pfahringer, Geoff Holmes, and Talel Abdesslem. Adaptive random forests for evolving data stream classification. *Machine Learning*, 2017.
- [2] João Gama, Indre Zliobaite, Albert Bifet, Mykole Pechenizkiy, and Abderlhamid Bouchachia. A survey on concept drift adaptation. *ACM Computing Surveys*, 46(4):44:1–44:37, March 2014.
- [3] Elena Ikonomovska, João Gama, and Sašo Džeroski. Learning model trees from evolving data streams. *Data mining and knowledge discovery*, 23(1):128–168, 2011.
- [4] Pedro Domingos and Geoff Hulten. Mining high-speed data streams. In *Proceedings of the Sixth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '00, pages 71–80, New York, NY, USA, 2000. ACM.
- [5] Wassily Hoeffding. Probability inequalities for sums of bounded random variables. *Journal of the American statistical association*, 58(301):13–30, 1963.
- [6] Ezilda Almeida, Carlos Ferreira, and Joao Gama. Adaptive model rules from data streams. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 480–492. Springer, 2013.
- [7] H. Mouss, D. Mouss, N. Mouss, and L. Sefouhi. Test of page-hinckley, an approach for fault detection in an agro-alimentary production system. In *Control Conference, 2004. 5th Asian*, volume 2, pages 815–818 Vol.2, 2004.
- [8] Jean Paul Barddal, Heitor Murilo Gomes, and Fabrício Enembreck. Advances on concept drift detection in regression tasks using social networks theory. *International Journal of Natural Computing Research (IJNCR)*, 5(1):26–41, 2015.
- [9] Elena Ikonomovska, João Gama, and Sašo Džeroski. Online tree-based ensembles and option trees for regression on evolving data streams. *Neurocomputing*, 150:458–470, 2015.
- [10] Heitor Murilo Gomes, Jean Paul Barddal, Fabrício Enembreck, and Albert Bifet. A survey on ensemble learning for data stream classification. *ACM Comput. Surv.*, 50(2):23:1–23:36, 2017.
- [11] Leo Breiman. Random forests. *Machine learning*, 45(1):5–32, 2001.
- [12] Elena Ikonomovska, João Gama, and Sašo Džeroski. Learning model trees from evolving data streams. *Data mining and knowledge discovery*, 23(1):128–168, 2011.
- [13] Albert Bifet and Ricard Gavaldà. Learning from time-changing data with adaptive windowing. In *SIAM*, 2007.
- [14] Geoff Hulten, Laurie Spencer, and Pedro Domingos. Mining time-changing data streams. In *Proceedings of the seventh ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 97–106. ACM, 2001.
- [15] Leo Breiman. Bagging predictors. *Machine learning*, 24(2):123–140, 1996.
- [16] Hadi Fanaee-T and Joao Gama. Event labeling combining ensemble detectors and background knowledge. *Progress in Artificial Intelligence*, pages 1–15, 2013.
- [17] Janez Demšar. Statistical comparisons of classifiers over multiple data sets. *Journal of Machine Learning Research*, 7:1–30, December 2006.