

Accelerating AdaBoost using UCB

Balázs Kégl

LAL/LRI, University of Paris-Sud, CNRS
Orsay, 91898, France

BALAZS.KEGL@GMAIL.COM

Róbert Busa-Fekete

LAL, University of Paris-Sud, CNRS
Orsay, 91898, France
Research Group on Artificial Intelligence of the
Hungarian Academy of Sciences and University of Szeged
Aradi vértanúk tere 1., H-6720 Szeged, Hungary

BUSAROBIGMAIL.COM

Editor: Gideon Dror, Marc Boullé, Isabelle Guyon, Vincent Lemaire, David Vogel

Abstract

This paper explores how multi-armed bandits (MABs) can be applied to accelerate AdaBoost. AdaBoost constructs a strong classifier in a stepwise fashion by adding simple base classifiers to a pool and using their weighted “vote” to determine the final classification. We model this stepwise base classifier selection as a sequential decision problem, and optimize it with MABs. Each arm represents a subset of the base classifier set. The MAB gradually learns the “utility” of the subsets, and selects one of the subsets in each iteration. ADABOOST then searches only this subset instead of optimizing the base classifier over the whole space. The reward is defined as a function of the accuracy of the base classifier. We investigate how the well-known UCB algorithm can be applied in the case of boosted stumps, trees, and products of base classifiers. The Fast Track of *KDD Cup 2009* was a large-scale learning task with a limited training time, thus this challenge offered us a good opportunity to test the utility of our approach. Our best results came in the *Up-selling* task where our model was within 2% of the best AUC rate.

Keywords: AdaBoost, Multi-Armed Bandit Problem, Upper Confidence Bound

1. Introduction

ADABOOST (Freund and Schapire, 1997) is one of the best off-the-shelf learning methods developed in the last decade. It constructs a classifier in a stepwise fashion by adding simple classifiers (called *base classifiers*) to a pool, and using their weighted “vote” to determine the final classification. The simplest base learner used in practice is the *decision stump*, a one-decision two-leaf decision tree. Learning a decision stump means to select a feature and a threshold, so the running time of ADABOOST with stumps is proportional to the number of data points n , the number of attributes d , and the number of boosting iterations T . When *trees* (Quinlan, 1993) or *products* (Kégl and Busa-Fekete, 2009) are constructed over the set of stumps, the computational cost is multiplied by an additional factor of the number of tree levels N or the number of terms m . Although the running time is linear in each of these factors, the algorithm can be prohibitively slow if the data size n and/or the number of features d is large.

There are essentially two ways to accelerate ADABOOST in this setting: one can either limit the number of data points n used to train the base learners, or one can cut the search space by using only

a subset of the d features. Although both approaches increase the number of iterations T needed for convergence, the net computational time can still be significantly decreased. The former approach has a basic version when the base learner is not trained on the whole weighted sample, rather on a small subset selected randomly using the weights as a discrete probability distribution (Freund and Schapire, 1997). A recently proposed algorithm of the same kind is FILTERBOOST (Bradley and Schapire, 2008), which assumes that an oracle can produce an unlimited number of labeled examples, one at a time. In each boosting iteration, the oracle generates sample points that the base learner can either accept or reject, and then the base learner is trained on a small set of accepted points. The latter approach was proposed by (Escudero et al., 2000) which introduces several feature selection and ranking methods used to accelerate ADABOOST. In particular, the LAZYBOOST algorithm chooses a fixed-size random subset of the features in each boosting iteration, and trains the base learner using only this subset. This technique was successfully applied to face recognition where the number of features can be extremely large (Viola and Jones, 2004).

In this paper we aim to improve the latter approach by “aiding” the random feature selection. It is intuitively clear that certain features are more important than others for classification. In specific applications the utility of features can be assessed a-priori (e.g., on images of characters, we know that background pixels close to the image borders are less informative than pixels in the middle of the images), however, our aim here is to *learn* the importance of features by evaluating their empirical performance during the boosting iterations. At the same time, in boosting it is important to keep a high level of base learner diversity, which is the reason why we opted for using *multi-armed bandits* (MAB) which are known to manage the exploration-exploitation trade-off very well.

MAB techniques have recently gained great visibility due to their successful applications in real life, for example, in the game of GO (Gelly and Silver, 2008). In the classical bandit problem the decision maker can select an arm at each discrete time step (Auer et al., 2002b). Selecting an arm results in a random reward, and the goal of the decision maker is to maximize the expected sum of the rewards received. Our basic idea is to partition the base classifier space into subsets and use MABs to learn the utility of the subsets. In each iteration, the bandit algorithm selects an optimal subset, then the base learner finds the best base classifier in the subset and returns a reward based on the accuracy of this optimal base classifier. By reducing the search space of the base learner, we can expect a significant decrease of the complete running time of ADABOOST. We use the UCB algorithm (Auer et al., 2002a) by assigning each feature to a subset. In the case of trees and products we use UCB by considering each tree or product as a sequence of decisions, and using the same partitioning as with decision stumps at each inner node.

The paper is organized as follows. First we describe the ADABOOST.MH algorithm and the necessary notations in Section 2. Section 3 contains our main contribution of using MABs for accelerating the selection of base classifiers. In Section 4 we present experiments conducted during the development period of the competition, our competition results, and some post-challenge analysis. Closing discussions are in Section 5.

2. ADABOOST.MH

For the formal description let $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$ be the $n \times d$ *observation matrix*, where $x_i^{(j)}$ are the elements of the d -dimensional observation vectors $\mathbf{x}_i \in \mathbb{R}^d$. We are also given a *label matrix* $\mathbf{Y} = (\mathbf{y}_1, \dots, \mathbf{y}_n)$ of dimension $n \times K$ where $\mathbf{y}_i \in \{+1, -1\}^K$. In *multi-class* classification one and only one of the elements of \mathbf{y}_i is +1, whereas in *multi-label* (or *multi-task*) classification \mathbf{y}_i

is arbitrary, meaning that the observation \mathbf{x}_i can belong to several classes at the same time. In the former case we will denote the index of the correct class by $\ell(\mathbf{x}_i)$.

```

ADABOOST.MH( $\mathbf{X}, \mathbf{Y}, \mathbf{W}^{(1)}, \text{BASE}(\cdot, \cdot, \cdot), T$ )
1   for  $t \leftarrow 1$  to  $T$ 
2        $\mathbf{h}^{(t)}(\cdot) \leftarrow \alpha^{(t)} \mathbf{v}^{(t)} \varphi^{(t)}(\cdot) \leftarrow \text{BASE}(\mathbf{X}, \mathbf{Y}, \mathbf{W}^{(t)})$ 
3       for  $i \leftarrow 1$  to  $n$  for  $\ell \leftarrow 1$  to  $K$ 
4            $w_{i,\ell}^{(t+1)} \leftarrow w_{i,\ell}^{(t)} \frac{\exp(-h_\ell^{(t)}(\mathbf{x}_i)y_{i,\ell})}{\sum_{i'=1}^n \sum_{\ell'=1}^K w_{i',\ell'}^{(t)} \exp(-h_{\ell'}^{(t)}(\mathbf{x}_{i'})y_{i',\ell'})}$ 
5   return  $\mathbf{f}^{(T)}(\cdot) = \sum_{t=1}^T \mathbf{h}^{(t)}(\cdot)$ 

```

Figure 1: The pseudocode of the ADABOOST.MH algorithm. \mathbf{X} is the observation matrix, \mathbf{Y} is the label matrix, $\mathbf{W}^{(1)}$ is the initial weight matrix, $\text{BASE}(\cdot, \cdot, \cdot)$ is the base learner algorithm, and T is the number of iterations. $\alpha^{(t)}$ is the base coefficient, $\mathbf{v}^{(t)}$ is the vote vector, $\varphi^{(t)}(\cdot)$ is the scalar base classifier, $\mathbf{h}^{(t)}(\cdot)$ is the vector-valued base classifier, and $\mathbf{f}^{(T)}(\cdot)$ is the final (strong) classifier.

The goal of the ADABOOST.MH algorithm ((Schapire and Singer, 1999), Figure 1) is to return a vector-valued classifier $\mathbf{f} : \mathcal{X} \rightarrow \mathbb{R}^K$ with a small *Hamming loss*

$$R_H(\mathbf{f}^{(T)}, \mathbf{W}^{(1)}) = \sum_{i=1}^n \sum_{\ell=1}^K w_{i,\ell}^{(1)} \mathbb{I} \left\{ \text{sign}(f_\ell^{(T)}(\mathbf{x}_i)) \neq y_{i,\ell} \right\}^1$$

by minimizing its upper bound (the exponential margin loss)

$$R_e(\mathbf{f}^{(T)}, \mathbf{W}^{(1)}) = \sum_{i=1}^n \sum_{\ell=1}^K w_{i,\ell}^{(1)} \exp(-f_\ell^{(T)}(\mathbf{x}_i)y_{i,\ell}), \quad (1)$$

where $f_\ell(\mathbf{x}_i)$ is the ℓ th element of $\mathbf{f}(\mathbf{x}_i)$. The user-defined weights $\mathbf{W}^{(1)} = [w_{i,\ell}^{(1)}]$ are usually set either uniformly to $w_{i,\ell}^{(1)} = 1/(nK)$, or, in the case of multi-class classification, to

$$w_{i,\ell}^{(1)} = \begin{cases} \frac{1}{2n} & \text{if } \ell = \ell(\mathbf{x}_i) \text{ (i.e., if } y_{i,\ell} = 1), \\ \frac{1}{2n(K-1)} & \text{otherwise (i.e., if } y_{i,\ell} = -1) \end{cases} \quad (2)$$

to create K well-balanced one-against-all classification problems. ADABOOST.MH builds the final classifier \mathbf{f} as a sum of *base classifiers* $\mathbf{h}^{(t)} : \mathcal{X} \rightarrow \mathbb{R}^K$ returned by a *base learner* algorithm $\text{BASE}(\mathbf{X}, \mathbf{Y}, \mathbf{W}^{(t)})$ in each iteration t . In general, the base learner should seek to minimize the *base objective*

$$E(\mathbf{h}, \mathbf{W}^{(t)}) = \sum_{i=1}^n \sum_{\ell=1}^K w_{i,\ell}^{(t)} \exp(-h_\ell(\mathbf{x}_i)y_{i,\ell}). \quad (3)$$

1. The indicator function $\mathbb{I}\{A\}$ is 1 if its argument A is true and 0 otherwise.

Using the weight update formula of line 4 (Figure 1), it can be shown that

$$R_e(\mathbf{f}^{(T)}, \mathbf{W}^{(1)}) = \prod_{t=1}^T E(\mathbf{h}^{(t)}, \mathbf{W}^{(t)}), \quad (4)$$

so minimizing (3) in each iteration is equivalent to minimizing (1) in an iterative greedy fashion. By obtaining the multi-class prediction

$$\widehat{\ell}(\mathbf{x}) = \arg \max_{\ell} f_{\ell}^{(T)}(\mathbf{x}),$$

it can also be proven that the “traditional” multi-class loss (or *one-error*)

$$R(\mathbf{f}^{(T)}) = \sum_{i=1}^n \mathbb{I} \left\{ \ell(\mathbf{x}_i) \neq \widehat{\ell}(\mathbf{x}_i) \right\} \quad (5)$$

has an upper bound $K R_e(\mathbf{f}^{(T)}, \mathbf{W}^{(1)})$ if the weights are initialized uniformly, and $\sqrt{K-1} R_e(\mathbf{f}^{(T)}, \mathbf{W}^{(1)})$ with the multi-class initialization (2). This justifies the minimization of (1).

2.1 Learning the base classifier

In this paper we use *discrete* ADABOOST.MH in which the vector-valued base classifier $\mathbf{h}(\mathbf{x})$ is represented as

$$\mathbf{h}(\mathbf{x}) = \alpha \mathbf{v} \varphi(\mathbf{x}),$$

where $\alpha \in \mathbb{R}^+$ is the *base coefficient*, $\mathbf{v} \in \{+1, -1\}^K$ is the *vote vector*, and $\varphi(\mathbf{x}) : \mathbb{R}^d \rightarrow \{+1, -1\}$ is a *scalar* base classifier. It can be shown that for minimizing (3), one has to choose φ that maximizes the *edge*

$$\gamma = \sum_{i=1}^n \sum_{\ell=1}^K w_{i,\ell} v_{\ell} \varphi(\mathbf{x}_i) y_{i,\ell}, \quad (6)$$

using the votes

$$v_{\ell} = \begin{cases} 1 & \text{if } \sum_{i=1}^n w_{i,\ell} \mathbb{I} \{ \varphi(\mathbf{x}_i) = y_{i,\ell} \} > \sum_{i=1}^n w_{i,\ell} \mathbb{I} \{ \varphi(\mathbf{x}_i) \neq y_{i,\ell} \}, \\ -1 & \text{otherwise,} \end{cases} \quad \ell = 1, \dots, K. \quad (7)$$

The optimal coefficient is then

$$\alpha = \frac{1}{2} \ln \frac{1 + \gamma}{1 - \gamma}.$$

It is also well known that the base objective (3) can be expressed as

$$E(\mathbf{h}, \mathbf{W}) = \sqrt{(1 + \gamma)(1 - \gamma)} = \sqrt{1 - \gamma^2}. \quad (8)$$

The simplest scalar base learner used in practice is the *decision stump*, a one-decision two-leaf decision tree of the form

$$\varphi_{j,b}(\mathbf{x}) = \begin{cases} 1 & \text{if } x^{(j)} \geq b, \\ -1 & \text{otherwise,} \end{cases}$$

where j is the index of the selected feature and b is the decision threshold. If the features are pre-ordered before the first boosting iteration, a decision stump maximizing the edge (6) can be found very efficiently in $\Theta(ndK)$ time (making the total running time $\Theta(nd(\log n + KT))$).

Although boosting decision stumps often yields satisfactory results, state-of-the-art performance of ADABOOST is usually achieved by using decision trees as base learners. In this paper we use an “in-house” implementation that calls the decision stump optimizer as a subroutine. The learner is similar to Quinlan’s C4.5 algorithm (Quinlan, 1993), except that we use the edge improvement (instead of Quinlan’s entropy-based criterion) to select the next node to split and the threshold b . The base learner has one hyperparameter, the number of the leaves N , which also shows up as a linear factor in the running time.

We also test our approach using a recently proposed base learner that seems to outperform boosted trees (Kégl and Busa-Fekete, 2009). The goal of this learner is to optimize products of simple base learners of the form

$$\mathbf{h}(\cdot) = \alpha \prod_{j=1}^m \mathbf{v}_j \varphi_j(\cdot), \quad (9)$$

where the vote vectors \mathbf{v}_j are multiplied element-wise. The learner also calls the decision stump optimizer as a subroutine but in an iterative rather than a recursive fashion. The hyperparameter m , again, appears as a linear factor in the total running time.

3. Using multi-armed bandits to reduce the search space

In this section we will first describe the MAB framework and next we show how bandit algorithm UCB can be used to accelerate the base learning step in AdaBoost.

3.1 Multi-armed bandits

In the classical bandit problem there are M arms that the decision maker can select at discrete time steps. Selecting arm j in iteration t results in a random reward $r_j^{(t)} \in [0, 1]$ whose (unknown) distribution depends on j . The goal of the decision maker is to maximize the expected sum of the rewards received. Intuitively, the decision maker’s policy has to balance between using arms with large past rewards (exploitation) and trying arms that have not been tested enough times (exploration). The UCB algorithm (Auer et al., 2002a) manages this trade-off by choosing the arm that maximizes the sum of the average reward

$$\bar{r}_j^{(t)} = \frac{1}{T_j^{(t)}} \sum_{t'=1}^t \mathbb{I}\{\text{arm } j \text{ is selected}\} r_j^{(t')}$$

and a confidence interval term

$$c_j^{(t)} = \sqrt{\frac{2 \ln t}{T_j^{(t)}}},$$

where $T_j^{(t)}$ is the number of times when arm j has been selected up to iteration t . To avoid the singularity at $T_j^{(t)} = 0$, the algorithm starts by selecting each arm once. We use a generalized version, denoted by $\text{UCB}(k)$, in which the best k arms are selected for evaluation, and the one that maximizes the actual reward $r_j^{(t)}$ is finally chosen.

3.2 The application of $\text{UCB}(k)$ for accelerating AdaBoost

The general idea is to partition the base classifier space into (not necessarily disjoint) subsets and use MABs to learn the utility of the subsets. In each iteration, the bandit algorithm selects an optimal subset (or, in the case of $\text{UCB}(k)$, a union of subsets). The base learner then finds the best base

classifier in the subset, and returns a reward based on this optimal base learner. By reducing the search space of the base learner, we can expect a significant decrease of the complete running time of ADABOOST.

The upper bound (4) together with (8) suggest the use of $-\frac{1}{2} \log(1 - \gamma^2)$ for the reward. In practice we found that

$$r_j^{(t)} = 1 - \sqrt{1 - \gamma^2}$$

works as well as the logarithmic reward; it was not surprising since the two are almost identical in the lower range of the $[0, 1]$ interval where the majority of the edges are. The latter choice has another advantage of always being in the $[0, 1]$ interval which is a formal requirement in MABs.

The actual partitioning of the base classifier set depends on the particular base learner. In the case of decision stumps, the most natural choice for UCB is to assign each feature to a subset. In principle, we could also further partition the threshold space but that would not lead to further savings in the linear computational time since, because of the changing weights $w_{i,\ell}$, all data points and labels would have to be visited anyway. On the other hand, subsets that contain more than one feature can be efficiently handled by $UCB(k)$.

In the case of trees and products we use UCB by considering each tree or product as a sequence of decisions, and using the same partitioning as with decision stumps at each inner node. In this setup we lose the information in the dependence of the decisions on each other *within* a tree or a product.

4. Experiments

4.1 Dataset description and data preparation

In the Fast Track of KDD Cup 2009 we were given a training set containing 50000 instances with three different label sets corresponding to the three tasks of *Churn*, *Appetency*, and *Up-selling*. The training and the standard evaluation dataset had the same size ($n = 50000$). Both contained 14740 numerical and 260 nominal features for a total of $d = 15000$. About 2% of all values were missing; we treated them by using an out-of-range value (i.e., we set all missing value to ∞). Next, we investigated the utility of the features. Using the cross-entropy based feature ranker of WEKA package (Witten and Frank, 2005), we found that only a relatively small number of features have positive score for any of task. In our experiments, we used only those features which had positive score for at least one of the three tasks. In this way we selected 6546 features.

4.2 Competition Results

We set number of iterations T to 8000 using a 60% – 40% simple validation on the training set. Figure 2 shows the AUC curves for the three tasks. Due to the time limit in the Fast Track we did not validate the number of tree leaves N and number of product terms m . We set the $N = 8$ and $m = 3$ based on the former experiments using our program package (Kégl and Busa-Fekete, 2009). The only remaining hyperparameter we had to choose was the number of best arms to be evaluated in the case of $UCB(k)$. We set k to 50. We also carried out some experiments with a lower value ($k = 20$), but we found that this only slightly influenced the results. Table 1 shows our official results. The *Churn* and *Up-selling* tasks were evaluated by ROC analysis because they were ranking problems, whereas the *Appetency* task, which were a classification task, was evaluated in terms of Balanced Accuracy².

2. <http://www.kddcup-orange.com/>

We also calculated the time requirements for the methods we employed on the large dataset. Table 2 shows the time in minutes needed for training and testing for 8000 iterations. This means that each model can be trained and tested in less than 10 hours on this large-scale database.

learner \ data set	Churn (AUC)		Appetency (B. Acc.)		Upselling (AUC)	
	Validation	Evaluation	Validation	Evaluation	Validation	Evaluation
STUMP	0.7424	0.6633	0.7051	0.6617	0.8938	0.8840
PRODC	0.6702	0.6314	0.6999	0.6807	0.8888	0.8666
TREE	0.7088	0.6819	0.7424	0.7343	0.8956	0.8883
BEST	–	0.7611	–	0.883	–	0.9038

Table 1: The validation and evaluation results in Fast Track.

Base learner	STUMP	PRODUCT	TREE
Time requirements	274	456	384

Table 2: Training and testing running times (in minutes) on the large dataset. The number of iterations T was equal to 8000.

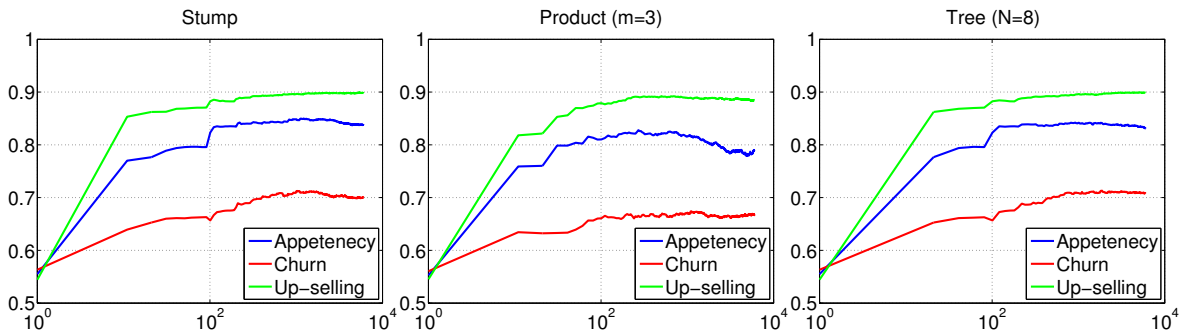


Figure 2: The ROC values vs. number of iterations on the validation set using discrete ADABOOST.MH with Stump, Product and Tree base learners.

5. Discussion and Conclusions

Our goal in the paper was to accelerate ADABOOST using multi-armed bandits. Recently, in machine learning, many applications have become the center of interest in which millions of training examples and thousands of features are not uncommon. In this scenario, fast optimization becomes more important than the asymptotic statistical optimality (Bottou and Bousquet, 2008). From this point of view, our approach has solved the tasks well because it reduced greatly the computational complexity of the learning phase. In the official competition our approach can achieve competitive results only on the *Up-selling* task. It seems that on small data set of the Slow Track we could also have results close to the best on the *Appetency* task, but since there was a confusion during the

competition (people merged the small and large data sets), we decided not to optimize our results. Once the submission site is reopened, we will provide quantitative results.

There are several possibilities to improve our boosting approach. First, handling missing values and nominal features can be done better than in our implementation. Our code was not yet prepared for these technical subtleties, but the algorithms are well known so it is only a question of programming effort to use them. We will also have to investigate the relatively large difference between the validation and test results. The early overfitting of ADABOOST indicates the existence of significant label noise in many learning tasks (Mease and Wyner, 2007); although standard solutions do not yet exist for this problem, we foresee to investigate the possibility to use a recent proposition of Freund (2009).

References

- P. Auer, N. Cesa-Bianchi, and P. Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine Learning*, 47:235–256, 2002a.
- P. Auer, N. Cesa-Bianchi, Y. Freund, and R.E. Schapire. The non-stochastic multi-armed bandit problem. *SIAM Journal on Computing*, 32(1):48–77, 2002b.
- L. Bottou and O. Bousquet. The tradeoffs of large scale learning. In *Advances in Neural Information Processing Systems*, volume 20, pages 161–168, 2008.
- J.K. Bradley and R.E. Schapire. FilterBoost: Regression and classification on large datasets. In *Advances in Neural Information Processing Systems*, volume 20. The MIT Press, 2008.
- G. Escudero, L. Márquez, and G. Rigau. Boosting applied to word sense disambiguation. In *Proceedings of the 11th European Conference on Machine Learning*, pages 129–141, 2000.
- Y. Freund. A more robust boosting algorithm. Technical report, Arxiv/0905.2138, 2009.
- Y. Freund and R. E. Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of Computer and System Sciences*, 55:119–139, 1997.
- S. Gelly and D. Silver. Achieving master level play in 9 x 9 computer go. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 1537–1540, 2008.
- B. Kégl and R. Busa-Fekete. Boosting products of base classifiers. In *Proceedings of the 26th International Conference on Machine Learning*, 2009.
- D. Mease and A. Wyner. Evidence contrary to the statistical view of boosting. *Journal of Machine Learning Research*, 9:131–156, 2007.
- J. Quinlan. *C4.5: Programs for Machine Learning*. Morgan Kaufmann, 1993.
- R. E. Schapire and Y. Singer. Improved boosting algorithms using confidence-rated predictions. *Machine Learning*, 37(3):297–336, 1999.
- P. Viola and M. Jones. Robust real-time face detection. *International Journal of Computer Vision*, 57:137–154, 2004.
- I.H. Witten and E. Frank. *Data Mining: Practical Machine Learning Tools and Techniques*. Morgan Kaufmann, 2005.