

Multi-class Boosting for Imbalanced Data

Antonio Fernández-Baldera¹, José M. Buenaposada², and Luis Baumela¹

¹ Dept. de Inteligencia Artificial, Universidad Politécnica de Madrid
antonio.fbaldera@upm.es, lbaumela@fi.upm.es

² ETSII, Universidad Rey Juan Carlos
josemiguel.buenaposada@urjc.es

Abstract. We consider the problem of multi-class classification with imbalanced data-sets. To this end, we introduce a cost-sensitive multi-class Boosting algorithm (*BAdaCost*) based on a generalization of the Boosting margin, termed *multi-class cost-sensitive margin*. To address the class imbalance, we also present a procedure to estimate the cost matrix from the 0|1-loss confusion matrix. Finally, we evaluate the performance of the approach with synthetic and real data-sets and compare our results with the AdaC2M1 algorithm.

Keywords: Boosting, multi-class classification, imbalanced data

1 Introduction

Imbalanced classification problems are characterized for having large differences in the number of samples in each class. This frequently occurs in complex data-sets, such as those involving class overlapping, small sample size, or within-class imbalance. In this situation, standard classifiers perform poorly since they minimize the number of misclassified training samples disregarding minority classes [1,2]. Solutions to the class imbalance problem may be coarsely organized into data-based, that re-sample the data space to balance the classes, and algorithm-based approaches, that introduce new algorithms that bias the learning towards the minority class [2,1]. Boosting methods have been extensively used to address the problem of classification with imbalanced data-sets [1,2,3] and cost-sensitive classification for two-class problems [4,5,6]. However, with the exception of AdaC2M1 [3], no previous work has addressed the problem of multi-class Boosting in presence of imbalanced data.

In our proposal we merge both multi-class and cost-sensitive perspectives into a new Boosting algorithm, BAdaCost, that stands for *Boosting Adapted for Cost-matrix*. We introduce the concept of *Multi-class Cost-sensitive Margin*, which serves as link between multi-class margins and the values of the cost matrix, both of them considered as argument of a loss function. We also present a procedure to estimate this matrix from the confusion matrix of a standard 0|1-loss classifier. We justify BAdaCost's good properties in a set of experiments.

2 Background

In this section we briefly review some Boosting results related to our proposal. We start by introducing AdaBoost [7]. Given N training data instances $\{(\mathbf{x}_i, l_i)\}$, where $\mathbf{x}_i \in X$ encodes the object to be classified and $l_i \in L = \{+1, -1\}$ is the class label, the goal of AdaBoost is learning a strong classifier $\text{sign}(\mathbf{H}(\mathbf{x})) = \text{sign}(\sum_{m=1}^M \beta_m G_m(\mathbf{x}))$ based on a linear combination of weak classifiers, $G_m : X \rightarrow L$. At each round m , a *direction* for classification, $G_m(\mathbf{x}) = \pm 1$, and a *step size*, β_m are added to an additive model whose goal is to minimize the empirical risk of the *Exponential Loss Function* [8], $\mathcal{L}(l, G_m(\mathbf{x})) = \exp(-l G_m(\mathbf{x}))$, defined over $z = l G_m(\mathbf{x})$, usually known as the *margin* [9].

2.1 Multi-class Boosting with vectorial encoding

A successful way to generalize the symmetry of class-label representation in the binary case to the multi-class case is using a set of vector-valued codes that represent the correspondence between the multi-class label set $L = \{1, \dots, K\}$ and a collection of vectors $Y = \{\mathbf{y}_1, \dots, \mathbf{y}_K\}$, where \mathbf{y}_k has a value 1 in the k -th coordinate and $\frac{-1}{K-1}$ elsewhere. It is immediate to see the equivalence between classifiers G defined over L and classifiers \mathbf{g} defined over Y , $G(\mathbf{x}) = l \in L \Leftrightarrow \mathbf{g}(\mathbf{x}) = \mathbf{y}_l \in Y$. Zou, Zhu and Hastie [9] used this codification to generalize the concept of binary margin to the multi-class case using a related vectorial codification in which a K -vector $\mathbf{y} = (y_1, \dots, y_K)^\top$ is said to be a *margin vector* if it satisfies the *sum-to-zero* condition, $\sum_{i=1}^K y_i = 0$. The SAMME algorithm generalizes the binary AdaBoost to the multi-class case [10]. It uses the above codification and an exponential loss whose risk is minimized using a stage-wise additive gradient descent approach. In this loss function the binary margin, $z = lG(\mathbf{x})$, is replaced by the multi-class vectorial margin defined with a scalar product, $z = \mathbf{y}^\top \mathbf{g}(\mathbf{x})$, producing the *Multi-class Exponential Loss Function* (MELF),

$$\mathcal{L}(\mathbf{y}, \mathbf{g}(\mathbf{x})) = \exp\left(-\frac{\mathbf{y}^\top \mathbf{g}(\mathbf{x})}{K}\right). \quad (1)$$

In this paper we generalize the class-label representation here described so that our Boosting algorithm can model the asymmetries arising when training on an unbalanced data set.

2.2 Cost-sensitive binary boosting

Classifiers that weigh certain types of errors more heavily than others are called cost-sensitive. They are used for example in medical diagnosis and object detection problems. Optimal cost-sensitive Boosting in two-class problems has been already studied in the literature [4,5,6]. The solution in [4] is based on minimizing the *Cost-sensitive Binary Exponential Loss Function* (CBELF)

$$\mathcal{L}(l, f(\mathbf{x})) = I(l = 1) \exp(-lC_1 f(\mathbf{x})) + I(l = -1) \exp(-lC_2 f(\mathbf{x})), \quad (2)$$

where $I(\cdot)$ is the indicator function and C_j are the costs of the two possible errors. Classification with imbalanced data-sets is a typical cost-sensitive problem. To counter balance the bias in the data we want the classifier’s loss function to under-weight errors from the majority class. In this paper we are guided by the cost-sensitive Boosting approach [4] to introduce an extension to the multiple-class case and use it to solve imbalanced problems.

3 Multi-class cost-sensitive margin

In this section we introduce the multi-class cost sensitive margin, based on which we derive the BAdaCost algorithm. Let us suppose the misclassification costs for our multi-class problem are encoded using a $K \times K$ -matrix \mathbf{C} , where each value $C(i, j)$ represents the cost of misclassifying an instance with real label i as j . We can assume without loss of generality [11] that $C(i, i) = 0, \forall i \in L$, i.e. the cost of correct classifications is null. We introduce an essential change in the MELF to handle this kind of problems. Firstly, let \mathbf{C}^* be a $K \times K$ -matrix defined in the following way

$$C^*(i, j) = \begin{cases} C(i, j) & \text{if } i \neq j \\ -\sum_{h=1}^K C(j, h) & \text{if } i = j \end{cases}, \quad \forall i, j \in L. \quad (3)$$

For our cost-sensitive classification problem each value $C^*(j, j)$ will represent a “negative cost” associated to a correct classification, i.e. a “reward”.

The j -th row in \mathbf{C}^* , denoted as $\mathbf{C}(j, -)$, is a margin vector that encodes the cost structure associated to the j -th label. By using it we can define the *multi-class cost-sensitive margin value* for an instance (\mathbf{x}, l) with respect to the multi-class vectorial classifier $\mathbf{g}(\cdot)$ as $z_C := \mathbf{C}^*(l, -) \cdot \mathbf{g}(\mathbf{x})$. It is easy to verify that if $\mathbf{g}(\mathbf{x}) = \mathbf{y}_i \in Y$, for a certain $i \in L$, then $\mathbf{C}^*(l, -) \cdot \mathbf{g}(\mathbf{x}) = \frac{K}{K-1} \mathbf{C}^*(l, i)$. Hence, multi-class cost-sensitive margins obtained from a discrete classifier $\mathbf{g} : \mathbf{x} \rightarrow Y$ can be computed using the “label valued” analogous of \mathbf{g} , $G : \mathbf{x} \rightarrow L$,

$$z_C = \mathbf{C}^*(l, -) \cdot \mathbf{g}(\mathbf{x}) = \frac{K}{K-1} \mathbf{C}^*(l, G(\mathbf{x})). \quad (4)$$

We use this generalized margin as argument for the MELF in order to obtain the *Cost-sensitive Multi-Class Exponential Loss Function (CMELF)*,

$$\mathcal{L}_C(l, \mathbf{g}(\mathbf{x})) := \exp(z_C) = \exp(\mathbf{C}^*(l, -) \cdot \mathbf{g}(\mathbf{x})) \quad (5)$$

as the loss function for our problem. Although any other margin-based loss functions could have been used, we use the exponential loss to maintain the similarity with the original AdaBoost algorithm. The new margin, z_C , yields negative values when classifications are correct under the cost-sensitive point of view, and positive values for from costly (wrong) assignments. Moreover, the range of margin values of z_C in (4) is much broader than the $z = \pm 1$ values of AdaBoost.

The CMELF is a generalization of the MELF (1) and CBELF (2). Let $\mathbf{C}_{0|1}$ be the cost matrix with zeros in the diagonal and ones elsewhere. This matrix encodes a multi-class problem free of costs. Further, it is well known that any matrix $\lambda \mathbf{C}_{0|1}$, with $\lambda > 0$, represents the same problem [11]. If we take into account that $\mathbf{y}^\top \mathbf{g}(\mathbf{x})$ has two possible values ($\frac{K}{K-1}$ when correct and $\frac{-K}{(K-1)^2}$ for errors) it is straightforward to prove that $\frac{1}{K(K-1)} \mathbf{C}_{0|1}$ will lead exactly to the same values of (1) when applied over the CMELF. In other words, the MELF is a particular case of the CMELF. On the other hand, it is also immediate to see that for a binary classification problem the values of \mathbf{C}^* lead to the CBELF (2). Hence, it is a special case of CMELF as well.

Vectorial classifiers, $\mathbf{f}(\mathbf{x}) = (f_1(\mathbf{x}), \dots, f_K(\mathbf{x}))^\top$, provides us with a *degree of confidence* for classifying sample \mathbf{x} into every class. Hence, they use the max rule, $\arg \max_{k=1, \dots, K} f_k(\mathbf{x})$, for label assignment [10,12]. It is immediate to see that this criterion is equivalent to assigning the label that maximizes the multi-class margin, $\arg \max_{k=1, \dots, K} \mathbf{y}_k^\top \mathbf{f}(\mathbf{x}) = \arg \min_{k=1, \dots, K} -\mathbf{y}_k^\top \mathbf{f}(\mathbf{x})$. Since $-\mathbf{y}_k^\top \mathbf{f}(\mathbf{x})$ is proportional to $\mathbf{C}_{0|1}^*(k, -)^\top \mathbf{f}(\mathbf{x})$, we can extend the decision rule to the cost-sensitive field just by assigning $\arg \min_{k=1, \dots, K} \mathbf{C}^*(k, -) \mathbf{f}(\mathbf{x})$.

4 BAdaCost: Boosting Adapted for Cost-matrix

In this section we present the BAdaCost, a multi-class cost-sensitive boosting algorithm. As we have defined the CMELF and given a training sample $\{(\mathbf{x}_i, l_i)\}$ we minimize the empirical expected loss, $\sum_{n=1}^N \mathcal{L}_C(l_n, \mathbf{f}(\mathbf{x}_n))$. The minimization is carried out by fitting an additive model, $\mathbf{f}(\mathbf{x}) = \sum_{m=1}^M \beta_m \mathbf{g}_m(\mathbf{x})$. The weak learner selected at each iteration m will consists of an optimal step of size β_m along the direction \mathbf{g}_m of the largest descent of the expected CMELF. In *Lemma 1* we show how to compute them.

Lemma 1 *Under the above assumptions, both β_m and \mathbf{g}_m are given by minimizing:*

$$(\beta_m, \mathbf{g}_m(\mathbf{x})) = \arg \min_{\beta, \mathbf{g}(\cdot)} \sum_{j=1}^K \left(S_j \exp(\beta C^*(j, j)) + \sum_{k \neq j} E_{j,k} \exp(\beta C^*(j, k)) \right), \quad (6)$$

where $S_j = \sum_{\{n: G(\mathbf{x}_n)=l_n=j\}} \omega_n$, $E_{j,k} = \sum_{\{n: l_n=j, G(\mathbf{x}_n)=k\}} \omega_n$. Given a known a direction \mathbf{g} , the optimal step $\beta(\mathbf{g})$ can be obtained as the solution to

$$\sum_{j=1}^K \sum_{k \neq j} E_{j,k} C^*(j, k) A(j, k)^\beta = - \sum_{j=1}^K S_j C^*(j, j) A(j, j)^\beta, \quad (7)$$

being $A(j, k) = \exp(C^*(j, k))$, $\forall i, j$. Finally, given a value of β , the optimal descent direction \mathbf{g} , equivalently $G(\cdot)$, is

$$\arg \min_{G(\cdot)} \sum_{n=1}^N \omega_n \left(A(l_n, l_n)^\beta I[G(\mathbf{x}_n) = l_n] + \sum_{k \neq l_n} A(l_n, k)^\beta I[G(\mathbf{x}_n) = k] \right). \quad (8)$$

The BAdaCost pseudo-code is shown in Algorithm 1. At each iteration, we add a new multi-class weak learner $\mathbf{g}_m : X \rightarrow Y$ to the additive model weighted by β_m , a measure of the confidence in the prediction of \mathbf{g}_m . The optimal weak learner that minimizes (8) is a cost-sensitive multi-class classifier trained using the data weights, w_i , and a modified cost matrix, \mathbf{C}_{wl} , with $C_{wl}(i, j) = A(i, j)^\beta, \forall i, j$.

Algorithm 1 : BAdaCost

- 1: **Input:** Cost matrix \mathbf{C} , N labeled training instances (\mathbf{X}, \mathbf{Y}) and number of iterations M
 - 2: **Output:** The trained weak learners and weights $(G_m, \beta_m), m = 1, \dots, M$
 - 3:
 - 4: Initialize weight vector $\mathbf{w}, \omega_i = 1/N; \forall i = 1, \dots, N$.
 - 5: $\mathbf{C}^* := \text{computeFullCostMatrix}(\mathbf{C})$ { *Using equation (3)* }
 - 6: **for** $m = 1, \dots, M$ **do**
 - 7: $\beta := 1; c := \infty; \Delta c := \infty$.
 - 8: **while** $\Delta c \geq \gamma$ **do**
 - 9: $\mathbf{C}_{wl} := \text{computeWLCostMatrix}(\mathbf{C}^*, \beta)$.
 - 10: $G := \text{trainMulticlassCostSensitiveWL}(\mathbf{X}, \mathbf{Y}, \mathbf{w}, \mathbf{C}_{wl})$.
 - 11: $\beta := \text{computeBeta}(\mathbf{C}^*, G, \mathbf{w}, \mathbf{Y})$ { *Solving equation (7)* }
 - 12: $c_{new} := \text{computeCost}(\mathbf{C}^*, G, \mathbf{w}, \mathbf{Y}, \beta)$ { *Using β and G in equation (6)* }
 - 13: $\Delta c := c - c_{new}; c := c_{new}$.
 - 14: **end while**
 - 15: $G_m := G; \beta_m := \beta$.
 - 16: Translate G_m into $\mathbf{g}_m : X \rightarrow Y$.
 - 17: Update weights $\omega_i = \omega_i \exp(\beta_m \mathbf{C}^*(l_i, -)\mathbf{g}_m(\mathbf{x}_i))$ for $i = 1, \dots, N$, and re-normalize vector \mathbf{w} .
 - 18: **end for**
 - 19: Output Classifier: $\mathbb{H}(\mathbf{x}) = \arg \min_k \mathbf{C}^*(k, -)\mathbf{f}(\mathbf{x})$, where $\mathbf{f}(\mathbf{x}) = \sum_{m=1}^M \beta_m \mathbf{g}_m(\mathbf{x})$.
-

5 Experiments

In this section we experimentally evaluate BAdaCost’s accuracy on imbalanced data-sets. In our experiments we use CART weak-learners and regularize our Boosting algorithm using shrinkage and re-sampling.

5.1 Cost matrix construction

A preliminary issue when using a cost-sensitive algorithm for solving an imbalance problem is establishing the cost matrix, \mathbf{C} . A straightforward solution would be to set the costs inversely proportional to the class imbalance ratios [1]. However, this solution does not take into account the complexity of the classification problem. i.e. the amount of class overlapping, within-class imbalance, etc. Here we introduce an alternative solution that considers the problem complexity.

Let \mathbf{F} be the confusion matrix and \mathbf{F}^* the matrix obtained when dividing each row i , $\mathbf{F}(i, -)$, by $\mathbf{F}(i, \cdot) = \sum_j F(i, j)$, i.e. the number of samples in class i . Then $F^*(i, j)$ is the proportion of data in class i classified as j . In a complex and imbalanced data-set, a 0|1-loss classifier (e.g. BAdaCost with 0|1-losses) will tend to over-fit the majority classes. So, off-diagonal elements in rows $\mathbf{F}^*(i, -)$ for majority (alt. minority) classes will have low (high) scores. Hence, the resulting matrix after setting $F^*(i, i) = 0, \forall i = 1 \dots K$ is already a cost matrix. Finally, to improve numerical conditioning, we set $\mathbf{C} = \lambda \mathbf{F}^*$, for a small $\lambda > 0$.

5.2 Synthetic datasets

The aim of this synthetic experiment is to visually analyze the performance of BAdaCost. We sample data from 2D Mixtures of Gaussian (MoG) probability distributions (see Fig. 1a) that represent a typical multi-class computer vision object detection problem. Two minority classes represent the target objects (colored areas in Fig. 1a) and a majority class that represents the background (black area in Fig. 1a). For training we sample 500 data for each class 1 (red) and 2 (green), and 5000 samples for class 3 (black). We also sample 1000 data from each class for testing. We run BAdaCost twice, firstly using a 0|1 cost matrix (see results in Fig. 1b) and in second place using a cost matrix built as described in Section 5.1 (see results in Fig. 1d).

In Fig. 1c we can see that, for 0|1 costs, training error evolves close to zero, whereas the testing error rate in a balanced data set levels-off above 0.4. This is an expected behavior, since this classifier optimizes the number of misclassified training samples, which come mostly from the background class, thus overfitting it. When using a cost matrix, see Fig. 1e, we get a much better testing error rate of 0.2. Note that the training error rate in this case is higher than that for the 0|1 cost matrix. This is also as expected, since the cost matrix has effectively moved the class boundary towards the majority class.

To visually appreciate the effect produced when training with an unbalanced data-set and the benefits of BAdaCost, in Fig. 1 b and d we show respectively the result of classifying all points on a grid in the feature space of this problem with the 0|1 and the imbalanced cost matrix. We can see a much better reconstruction using the imbalanced cost matrix.

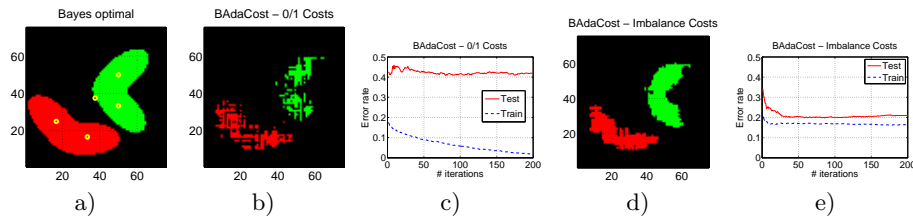


Fig. 1: Synthetic experiment 1. Simulated computer vision object detection problem with majority background class (black).

5.3 Real data-set: Synapse and Mitochondria segmentation

In the last years we have seen advances in the automated acquisition of large series of images of brain tissue. The complexity of these images and the high number of neurons in a small section of the brain, makes the automated analysis of these images the only practical solution. Mitochondria and synapses are two interesting cell structures that will be the object of detection. Unfortunately, the proportion of them w.r.t. the background is quite small, which makes the problem highly skewed. We collected a training set composed of 10000 background, 4000 mitochondria and 1000 synapse data and a testing set with 20000 data per class.

In this section we use the BAdaCost algorithm to label pixels in these images as mitochondria, synapse and background, and compare the results with those achieved by the AdaC2M1 algorithm. Following [13], we apply to each image in the stack a set of linear Gaussian filters at different scales to compute zero, first and second order derivatives. For each pixel we get a vector of responses $S = (s_{00}, s_{10}, s_{01}, s_{02}, s_{11}, s_{02})$ that are respectively obtained applying the filters $G_{\sigma^*}, \sigma \cdot G_{\sigma^*} * \frac{\partial}{\partial x}, \sigma \cdot G_{\sigma^*} * \frac{\partial}{\partial y}, \sigma^2 \cdot G_{\sigma^*} * \frac{\partial^2}{\partial x^2}, \sigma^2 \cdot G_{\sigma^*} * \frac{\partial^2}{\partial xy}, \sigma^2 \cdot G_{\sigma^*} * \frac{\partial^2}{\partial y^2}$ where G_{σ} is a zero mean Gaussian with σ standard deviation. For a given σ the pixel feature vector is given by $f(\sigma) = (s_{00}, \sqrt{s_{10}^2 + s_{01}^2}, \lambda_1, \lambda_2)$ where λ_1 and λ_2 are the eigenvalues of the Hessian matrix of the pixel, that depend on s_{20}, s_{02} and s_{11} . The final 16 dimensional feature vector for each pixel is given by the concatenation of the $f(\sigma)$ vector at 4 scales (values of σ).

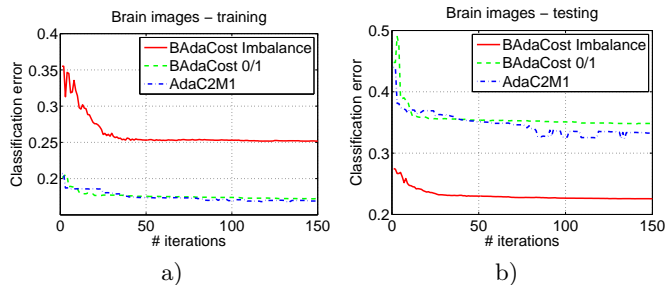


Fig. 2: Brain images experiment with a heavily imbalanced data-set.

In this experiment we compare BAdaCost (with 0|1 costs), AdaC2M1 and BAdaCost (with the imbalanced cost matrix described in Section 5.1). In Fig. 2 we show the training and testing classification errors of the three algorithms. The behavior of BAdaCost in the real experiment is similar to those obtained with synthetic data. The 0|1-cost classifier has lower training error and higher testing error, whereas the classifier with imbalanced cost matrix achieves the best generalization on the test set. The AdaC2M1 classifier with imbalanced matrix achieve marginally better results than the 0|1-cost BAdaCost, but clearly worse than the imbalanced BAdaCost.

6 Conclusions

In this paper we have addressed the problem of multi-class classification with imbalanced data-sets. By extending the notion of multi-class margin to the cost-sensitive margin we introduced the *BAdaCost* algorithm. We have shown experimentally that BAdaCost performs as expected from a cost-sensitive algorithm and outperforms the AdaC2M1 when dealing with imbalanced data.

Additionally, to address the class imbalance, we developed a simple procedure to estimate the cost matrix from the 0|1-loss confusion matrix.

References

1. Galar, M., Fernandez, A., Barrenechea, E., Bustince, H., Herrera, F.: A review on ensembles for the class imbalance problem: Bagging-, boosting-, and hybrid-based approaches. *IEEE Transactions on Systems, Man, and Cybernetics, Part C: Applications and Reviews* **42**(4) (July 2012) 463–484
2. He, H., Garcia, E.A.: Learning from imbalanced data. *IEEE Transactions on Knowledge and Data Engineering* **21**(9) (sept. 2009) 1263–1284
3. Sun, Y., Kamel, M.S., Wang, Y.: Boosting for learning multiple classes with imbalanced class distribution. In: *Proc. of the Int. Conference on Data Mining. ICDM '06* (2006) 592–602
4. Masnadi-Shirazi, Hamed, Vasconcelos, N.: Cost-sensitive boosting. *Transactions on Pattern Analysis and Machine Intelligence* **33** (2011) 294–309
5. Fan, W., Stolfo, S.J., Zhang, J., Chan, P.K.: Adacost: Misclassification cost-sensitive boosting. In: *Proceedings of the 16th International Conference on Machine Learning*. (1999) 97–105
6. Landesa-Vazquez, I., Alba-Castro, J.L.: Double-base asymmetric adaboost. *Neurocomputing* **118** (2013) 101–114
7. Freund, Y., Schapire, R.E.: A decision theoretic generalization of on-line learning and an application to boosting. *Journal of Computer and System Sciences* **55** (1997) 199–139
8. Friedman, J., Hastie, T., Tibshirani, R.: Additive logistic regression: a statistical view of boosting. *The Annals of Statistics* **28**(2) (2000) 337–407
9. Zou, H., Zhu, J., Hastie, T.: New multiclass boosting algorithms based on multiclass fisher-consistent losses. *Annals of Applied Statistics* **2** (2008) 1290–1306
10. Zhu, J., Zou, H., Rosset, S., Hastie, T.: Multi-class AdaBoost. *Statistics and Its Interface* **2** (2009) 349–360
11. O’Brien, D.B., Gupta, M.R., Gray, R.M.: Cost-sensitive multi-class classification from probability estimates. In: *Proceedings of the 25th International Conference on Machine Learning*. (2008) 712–719
12. Fernandez-Baldera, A., Baumela, L.: Multi-class boosting with asymmetric weak-learners. *Pattern Recognition* **47**(5) (May 2014) 2080–2090
13. Cetina, K., Márquez-Neila, P., Baumela, L.: A comparative study of feature descriptors for mitochondria and synapse segmentation. In: *Proc. of the Int. Conference on Pattern Recognition*. (2014)