



This document is shared under a **CC-BY-NC-ND license**.

You are free to share (that is, copy and redistribute the material in any medium or format) if you follow these licence terms:



**Attribution (by)**

You must give appropriate credit, provide a link to the license.



**Non Commercial (nc)**

You can copy, distribute, display, perform, and use this material for any purpose other than commercially (unless you get permission first). Non Commercial means not primarily intended for or directed towards commercial advantage or monetary compensation.



**No Derivatives (nd)**

If you remix, transform, or build upon the material, you may **not** distribute the modified material. But note that simply changing the format does not create a derivative.

**No additional restrictions** — You may not apply legal terms or technological measures that legally restrict others from doing anything the licence permits.

#### **Notices**

You do not have to comply with the license for elements of the material in the public domain or where your use is permitted by an applicable exception or limitation.

For the full text of this Creative Commons licence go to:

<https://creativecommons.org/licenses/by-nc-nd/4.0/legalcode>

# An effective learning strategy for cascaded object detection

A. Bria<sup>a</sup>, C. Marrocco<sup>a</sup>, M. Molinara<sup>a</sup>, F. Tortorella<sup>a,\*</sup>

<sup>a</sup>*Department of Electrical and Information Engineering, University of Cassino and L.M.,  
Via Di Biasio 43, 03043 Cassino (FR), Italy*

---

## Abstract

To distinguish objects from non-objects in images in real-time, a suitable solution is to employ a cascade detector that consists of a sequence of node classifiers with increasing discriminative power. However, among the millions of image patches generated from an input image, only very few contain the searched object. When trained on these highly unbalanced data sets, the node classifiers tend to have poor performance on the minority class. Thus, we propose a learning strategy aimed at maximizing the node classifiers ranking capability rather than their accuracy. We also provide an efficient implementation yielding the same time complexity of the original Viola-Jones cascade training. Experimental results on highly unbalanced real problems show that our approach is both efficient and effective when compared to other node training strategies for skewed classes.

*Keywords:* pattern recognition, image processing, object detection, cascade, unbalanced data, ranking

---

## 1. Introduction

2 Detecting objects in images and videos is a crucial task in many real-world  
3 problems, ranging from face detection in images to pedestrian detection on

---

\*Corresponding author. Address: Department of Electrical and Information Engineering, University of Cassino and L.M., Via Di Biasio 43, 03043 Cassino (FR), Italy. Tel.: +39 07762993605

*Email addresses:* a.bria@unicas.it (A. Bria), c.marrocco@unicas.it (C. Marrocco), m.molinara@unicas.it (M. Molinara), tortorella@unicas.it (F. Tortorella)

4 roads, from biomedical image analysis to videosurveillance applications. From  
5 the pattern recognition point of view, object detection can be cast as a classi-  
6 fication problem in which one has to distinguish between the object class and  
7 the “non-object” class. The difficulty is that the “non-object” class actually  
8 contains all the patches in the images that are not instances of the searched  
9 object and this makes very challenging this apparently simple problem.

10 First, whereas the object class is quite precisely defined, the non-object class  
11 can be very inhomogeneous, since it includes several subclasses of other objects  
12 very different each other. In real scenes, the searched objects are embedded in  
13 a cluttered background containing various kinds of non-objects, some of which  
14 are very different from what we are searching for, whereas others can be very  
15 similar. As an example, when detecting faces in an office scene, it would be  
16 simple to distinguish between a face and a chair, but it could be more complex  
17 to distinguish between a face and a cartoon.

18 The second issue is that almost every real-world detection task requires  
19 processing a huge number of pixels and this could involve a computational load  
20 not easy to bear, specially if a complex classifier architecture is used.

21 The last point is related to the skewed class priors: in a typical detection  
22 problem, the vast majority of image locations do not contain the searched ob-  
23 ject and this makes detection a severely unbalanced classification problem. As  
24 a consequence, learning an effective classifier is very difficult since classifiers  
25 trained on highly unbalanced data sets tend to have poor performance on the  
26 minority class.

27 In this regard, a commonly adopted solution, originally proposed by Viola  
28 and Jones [22], is to employ an ensemble of classifiers structured as a cascade  
29 of dichotomizers with increasing complexity. Such an approach allows each di-  
30 chotomizer in the cascade to deal with only a part of the non-object class, thus  
31 parting the complexity of the whole problem among the classifiers. In particular,  
32 the first stages of the cascade are built to reject the most distinguishable back-  
33 ground regions, while the last stages are specialized to discriminate between  
34 actual object and the most confusing background patches (see Fig. 1). This

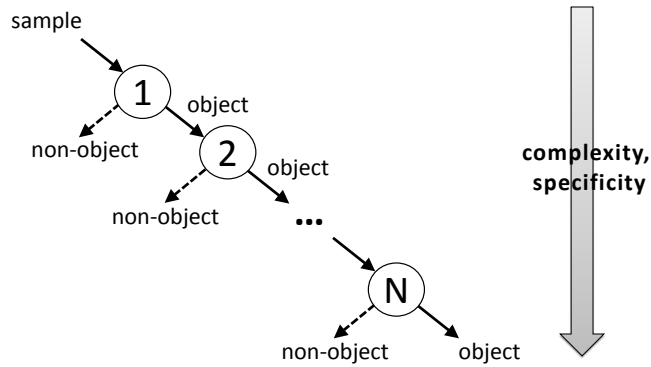


Figure 1: Scheme of the cascade structure with  $N$  nodes

35 aims at reducing the number of false positives produced by the detector and  
 36 concentrate the computational complexity of the system on the last classifiers  
 37 of the cascade.

38 Such approach is a valid solution to the first two issues. Indeed, whereas  
 39 a monolithic classifier would hardly ensure both a good sensitivity and a good  
 40 specificity and, in any case, typically requires an excessive computational bur-  
 41 den, the cascade of dichotomizers provides a high constant sensitivity and a  
 42 growing specificity through the stages at a reasonable computational cost. This  
 43 is obtained by connecting classifiers that provide high sensitivity and sufficient  
 44 specificity on the subproblem they face and are sufficiently simple to ensure a  
 45 real-time response. To this end, a commonly used learning algorithm is *Ada-*  
 46 *Boost* [9] that allows both effective feature selection and proficient classifica-  
 47 tion.

48 However, the original Viola-Jones approach does not deal well with the class  
 49 imbalance problem. In fact, even though each dichotomizer in the cascade is  
 50 trained to discriminate between the object class and a part of the non-object  
 51 class, the problem is asymmetric anyway since the object class is still much less  
 52 numerous. In this regard, *AdaBoost* is not able to effectively face the asymmetry  
 53 between the classes since it minimizes a quantity related to classification error,  
 54 that is significantly biased by skewed class priors. To address this problem, some  
 55 new approaches have been developed. In particular, Viola and Jones proposed

56 a new algorithm, *AsymBoost* [23], that handles the unbalanced classes through  
57 an asymmetric weight updating mechanism of the samples in the training set.  
58 Starting from this approach, Visentini et al. [24] proposed *AsymBoost\** that  
59 extended the *AsymBoost* cascade algorithm by introducing a reactive control  
60 of the asymmetry at both cascade and node learning level. Both these ap-  
61 proaches modify the learning strategy of the dichotomizer without abandoning  
62 the beneficial *AdaBoost* scheme that simultaneously selects features and builds  
63 the classifier.

64 Other approaches alter the Viola-Jones standard framework: for example, a  
65 *multiexit cascade* is proposed in [13], where the  $i$ th node combines the scores  
66 from the first  $i - 1$  dichotomizers, while [28] decouples the feature selection  
67 process and the node classifier, which is designed to explicitly address the class  
68 asymmetry. In [25] the two previous approaches are combined. All such solu-  
69 tions, however, do not guarantee that the benefits of the original cascade frame-  
70 work are still provided: in fact, the multiexit cascade increases the complexity  
71 and the computational load of the single node whereas the separated feature  
72 selection process could not make the optimal feature choice for the subsequent  
73 classifier learning.

74 In this paper we address the problem of class imbalance within the Viola-  
75 Jones standard cascade framework where we introduce a new learning algorithm  
76 for the node classifiers aimed at maximizing their ranking capability rather than  
77 their accuracy. Although employed by all the techniques described before, accu-  
78 racy is not a correct measure for handling rare cases that have less impact than  
79 common cases [27]. As a consequence, learning algorithms based on accuracy  
80 lead to poor minority-class performance [26]. A common way to alleviate this  
81 problem is to adopt a cost-sensitive learning method that assigns a higher cost  
82 to the errors made on the minority class (as in *AsymBoost* and *AsymBoost\**).  
83 The difficulty of this approach lies in correctly defining the cost matrix: a com-  
84 mon choice is to set the ratio of the error costs equal to the ratio of the prior  
85 probabilities of the two classes, but it is not always possible to reliably estimate  
86 such value. Another possible approach is to attempt to balance the class distri-

87 butions by undersampling the majority class or by oversampling the minority  
88 class, but both solutions present serious drawbacks. Whereas undersampling  
89 can potentially remove important examples of the majority class, the samples  
90 artificially generated and added to the minority class in oversampling can pro-  
91 duce an unfaithful training set [12].

92 Our approach is to use rank metrics instead of accuracy metrics for train-  
93 ing the single dichotomizer in the cascade. More specifically, rank metrics are  
94 more concerned with the relative ordering of cases than with making absolute  
95 predictions for cases and thus place more emphasis on learning to distinguish  
96 classes than on learning the internal structure of classes [3]. To this end, the  
97 training of the node dichotomizers is based on a reformulation of *RankBoost* for  
98 bipartite ranking problems [8], suitably modified to be embedded in a cascade  
99 structure. This solution allows the detector to take advantage of the benefits  
100 of the standard cascade framework and to effectively face the asymmetry still  
101 present in each node classifier training. A partial and preliminary presentation  
102 of this approach has been made in [2], whereas in this paper we extend the theo-  
103 retical framework and consider other significant detection problems in the exper-  
104 iments. In particular, the cascade trained with rank-based node dichotomizers  
105 is compared with the original Viola-Jones method, with *AsymBoost* and with  
106 *AsymBoost\**. Moreover, we have also considered a cascade of dichotomizers  
107 trained with *RUSBoost* [18], a boosting-based learning algorithm that faces the  
108 class imbalance by undersampling the majority class. As real-world applica-  
109 tions, we have examined in our experiments a face detection problem and two  
110 medical imaging problems: the detection of microcalcifications on digital mam-  
111 mograms (to extend the preliminary comparison in [2]) and the detection of  
112 microaneurysms in digital fundus images. All the examined applications are  
113 characterized by a considerable class imbalance and the obtained results show  
114 that the proposed approach provides performance similar or better than the  
115 other algorithms thus confirming its effectiveness.

116 The rest of this paper is organized as follows: Section 2 recalls some previ-  
117 ously proposed learning strategies for the node classifier, while in Section 3 we

118 describe the rank-based learning algorithm for the node dichotomizer. Section 4  
119 provides an efficient implementation of the proposed approach. In Section 5 the  
120 results of an extensive experimental comparison are reported and discussed.  
121 Finally, Section 6 draws some conclusions and outlines directions for future  
122 research.

## 123 2. Accuracy-based node learning strategies for skewed classes

124 To point out the drawbacks of accuracy-based node training in the original  
125 Viola-Jones framework, let us briefly recall the *AdaBoost* approach. A weak  
126 learner  $h_\tau(\cdot)$  is selected in each of a series of rounds  $\tau = 1, 2, \dots$ , so as to  
127 minimize the *weighted exponential loss*  $\sum_j D_\tau(j) \exp(-y_j h_\tau(\mathbf{x}_j))$ , where  $D_\tau(j)$   
128 is the weight on the  $j$ th sample  $\mathbf{x}_j$  and  $y_j \in \{-1, 1\}$  is the class label of  $\mathbf{x}_j$ . It is  
129 easy to see that, in the case of highly skewed classes, such sum is dominated by  
130 the error produced on the majority class and thus the choice of the weak learner  
131 is not optimal for predicting the minority class. The approaches proposed so far  
132 to alleviate such problem can be attributed to two different groups: *cost-based*  
133 *strategies* and *sampling-based strategies*.

### 134 2.1. Cost-based strategies

135 This category includes the methods that employ the *Adaboost* learning pro-  
136 cedure, with a cost assigned to false negatives greater than to false positives.  
137 Both *AsymBoost* and *AsymBoost\** follow this way. In *AsymBoost*, each sample  
138  $\mathbf{x}_j$  is pre-weighted at each round with an asymmetric cost  $\exp\left(\frac{1}{T} y_j \log \sqrt{k}\right)$   
139 where  $T$  is the total number of rounds of the dichotomizer. The parameter  $k$   
140 estimates how much more false negatives cost than false positives and its choice  
141 should bias the classifier to perform well on the minority object class. The dif-  
142 ference between *AsymBoost* and *AsymBoost\** is in the value of  $k$  that is fixed  
143 in *AsymBoost* whereas it is automatically tuned at each round in *AsymBoost\**.

144 Other cost-sensitive boosting algorithms have been proposed in the literature  
145 [11, 7, 21, 19]. The problem common to all these methods is that specific cost  
146 information is rarely available or hard to estimate, as well as the ratio between  
147 the costs.

148 *2.2. Sampling-based strategies*

149 Another popular solution is to modify the original distribution of the data so  
150 as to obtain an artificially balanced training set. This approach can be accom-  
151 plished by two means: undersampling the majority class by eliminating data  
152 points until it reaches the size of the minority class or oversampling the minor-  
153 ity class by adding artificially generated data points until the desired balance  
154 is achieved. An algorithm introducing data undersampling into the *AdaBoost*  
155 procedure is *RUSBoost* [18] that applies a random undersampling of the major-  
156 ity class for building a balanced training set at each iteration of the boosting  
157 algorithm. On the other hand, *SMOTEBoost* [5] applies oversampling to *Ad-*  
158 *aBoost* according to *SMOTE* [4], an algorithm that oversamples the minority  
159 class by introducing new, non-replicated minority-class examples. Even though  
160 very popular, such approach is hardly applicable when dealing with high dimen-  
161 sional feature spaces, as in the case of the Viola-Jones framework. Moreover, in  
162 particular applications (e.g. medical imaging) it could be unsafe to add artifi-  
163 cially generated samples that do not correspond to real situations.

164 **3. Ranking-based node learning**

165 In this section we draw a ranking-based learning strategy to train the cascade  
166 dichotomizers. The dichotomizer  $H_i(\mathbf{x})$  added at the  $i$ th stage is based on a  
167 reformulation of *RankBoost* for bipartite ranking problems [8], suitably modified  
168 to be embedded in a cascade structure. It consists of a linear combination of  
169 weak learners  $h_{i,\tau}(\mathbf{x}) \in \{0, 1\}$  (0 for the negative class, 1 for the positive one)  
170 weighted by  $\alpha_{i,\tau} \in \mathbb{R}$  and added in subsequent rounds  $\tau = 1, 2, \dots$  so that after  
171  $t$  rounds we obtain

$$H_{i,t}(\mathbf{x}) = \sum_{\tau=1}^t \alpha_{i,\tau} h_{i,\tau}(\mathbf{x}). \quad (1)$$



172 To simplify the notation, let us omit the subscript  $i$  throughout this section. A  
 173 weak learner  $h_\tau(\mathbf{x})$  consists of a simple *decision stump* given by

$$h_\tau(\mathbf{x}) = \begin{cases} 1 & \text{if } \varphi_\tau(\mathbf{x}) > \theta_\tau \\ 0 & \text{if } \varphi_\tau(\mathbf{x}) \leq \theta_\tau \end{cases} \quad (2)$$

174 where  $\varphi_\tau(\mathbf{x}) \in \mathcal{F}$  is the feature selected at round  $\tau$  from the feature set  $\mathcal{F}$  and  
 175  $\theta_\tau$  is the threshold selected at round  $\tau$  in the range  $\left(\min_{\mathbf{x} \in \mathcal{T}} \varphi_\tau(\mathbf{x}), \max_{\mathbf{x} \in \mathcal{T}} \varphi_\tau(\mathbf{x})\right)$ ,  
 176 with  $\mathcal{T}$  being the training set. Differently from the original *RankBoost*, such  
 177 weak learners do not abstain and are  $\{0, 1\}$ -valued so as to preserve the or-  
 178 dering information provided by the feature while ignoring its specific scoring  
 179 information.

180 To maximize the ranking capability of the dichotomizer, let us consider the  
 181 *crucial pairs*  $(\mathbf{p}, \mathbf{n})$ , defined as all the sample pairs made by a positive sample  $\mathbf{p}$   
 182 and a negative sample  $\mathbf{n}$ . A crucial pair is correctly ranked when  $H_t(\mathbf{n}) < H_t(\mathbf{p})$ .  
 183 Therefore, our goal is to minimize the *ranking loss*  $R_t$ , defined as the number  
 184 of misranked crucial pairs and given by<sup>1</sup>

$$R_t = \sum_{(\mathbf{p}, \mathbf{n})} [H_t(\mathbf{n}) \geq H_t(\mathbf{p})] \quad (3)$$

185 A weight distribution  $w_t(\mathbf{p}, \mathbf{n})$  is maintained over the crucial pairs so that the  
 186 misranked crucial pairs will be more influential in the following rounds. The  
 187 weight update rule is given by

$$w_{t+1}(\mathbf{p}, \mathbf{n}) = \frac{w_t(\mathbf{p}, \mathbf{n}) \exp(\alpha_t (h_t(\mathbf{n}) - h_t(\mathbf{p})))}{W_t} \quad (4)$$

188 where  $W_t$  is a normalization factor so that  $\sum_{\mathbf{p}, \mathbf{n}} w_{t+1}(\mathbf{p}, \mathbf{n}) = 1$ . Assuming  
 189  $\alpha_t > 0$ , the update rule decreases the weight of crucial pairs in case of correct  
 190 ranking (i.e.,  $h_t(\mathbf{p}) = 1$  and  $h_t(\mathbf{n}) = 0$ ) and increases the weight otherwise.  
 191 Combining Eq. 3 with Eq. 4, the goal becomes to minimize the weighted number

---

<sup>1</sup>the notation  $[pr]$  (Iverson bracket) is defined to be 1 if predicate  $pr$  holds and 0 otherwise.

192 of misranked crucial pairs  $R_t^w$  (*weighted ranking loss*) given by

$$R_t^w = \sum_{(\mathbf{p}, \mathbf{n})} w_t(\mathbf{p}, \mathbf{n}) [H_t(\mathbf{n}) \geq H_t(\mathbf{p})] \quad (5)$$

193 from which  $h_t(\cdot)$  (i.e.,  $\varphi_t(\cdot)$  and  $\theta_t$ ) and  $\alpha_t$  can be chosen by

$$(h_t, \alpha_t) = \arg \min_{h, \alpha} \sum_{(\mathbf{p}, \mathbf{n})} w_t(\mathbf{p}, \mathbf{n}) [H_{t-1}(\mathbf{n}) + \alpha h(\mathbf{n}) \geq H_{t-1}(\mathbf{p}) + \alpha h(\mathbf{p})] \quad (6)$$

194 In fact,  $\alpha_t$  can be found directly by minimizing an upperbound of  $R_t^w$  as in [8],  
195 which yields

$$R_t^w \leq \sqrt{1 - r_t^2} \quad (7)$$

196 by choosing

$$\alpha_{h_t} = \frac{1}{2} \ln \left( \frac{1 + r_t}{1 - r_t} \right) \quad (8)$$

197 where

$$r_t = \sum_{\mathbf{p}, \mathbf{n}} w_t(\mathbf{p}, \mathbf{n}) (h_t(\mathbf{p}) - h_t(\mathbf{n})) \quad (9)$$

198 Then, combining Eq. 6 with Eq. 7, the choice of  $h_t(\cdot)$  at round  $t$  becomes

$$h_t = \arg \max_h |r_t| \quad (10)$$

199 In summary, at each round  $t$ , the choice of the weak learner  $h_t(\cdot)$  and sub-  
200 sequently the selection of the feature  $\varphi_t(\cdot)$  is made to maximize the correct  
201 pairwise ranking. An example of the training process described so far is shown  
202 in Fig. 2.

#### 203 4. Efficient implementation of ranking-based training

204 The training procedure presented in the previous section is embedded in a  
205 Viola-Jones standard cascade structure which is constructed as follows. Let  $\mathcal{P}$   
206 and  $\mathcal{N}$  denote the sets of positive and negative samples. The node classifiers  
207  $H_i(\mathbf{x})$  are added in subsequent stages  $i = 1, 2, \dots$  until the desired false positive  
208 rate  $F_{target}$  is reached or  $\mathcal{N} = \{\emptyset\}$ . To construct the  $i$ th node classifier  $H_i(\mathbf{x})$ ,  
209  $\mathcal{P}$  and  $\mathcal{N}$  (hereafter referred to as *pool*) are used to form a training set  $\mathcal{T}_i$  and

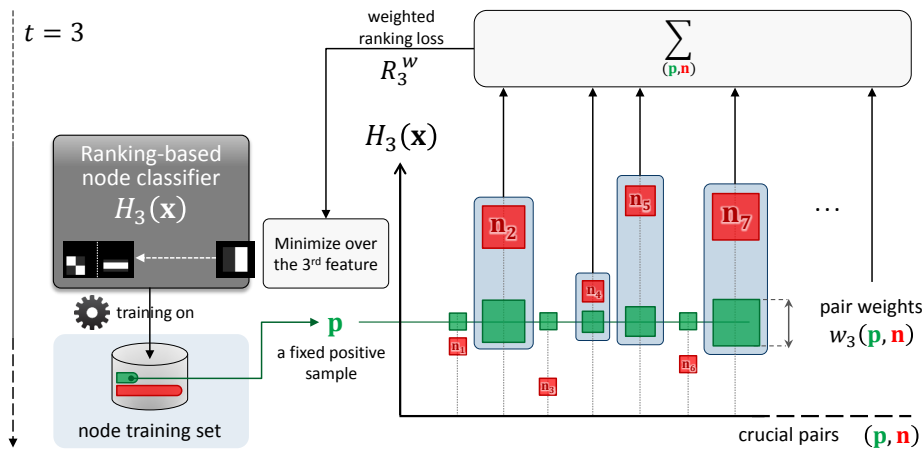


Figure 2: An example of the training process for a node classifier. At each round  $t$  ( $t = 3$  in the depicted graph), a new weak learner which is restricted to use a single feature is added. The resulting classifier function is then used to rank the samples in the node training set. The sum of weights of misranked crucial pairs (highlighted with a gray box) is evaluated and the feature minimizing this quantity is selected and added to the current node.

210 a validation set  $\mathcal{V}_i$ , respectively to train and to set the decision threshold  $\Theta_i$  of  
 211 the  $i$ th node classifier  $H_i(\mathbf{x})$  so as to meet the node learning goals  $d$  (detection  
 212 rate) and  $f$  (false positive rate) (see Algorithm 1, lines 13-14).

213 A naive implementation of this approach may result in a computational  
 214 workload not easy to sustain in the training phase, especially when dealing with  
 215 millions of negative samples, as in the case of many object detection problems.  
 216 Indeed, there are three major issues that make it challenging to efficiently im-  
 217 plement the proposed learning strategy: (i) the high computational complexity  
 218 of weak learner selection in Eq. 10; (ii) the computation of features throughout  
 219 training; and (iii) the need of a criterion to control the number of training rounds  
 220 in a node. Each of these problems is discussed in the following subsections. The  
 221 pseudocode of the overall training procedure is provided in Algorithm 1 along  
 222 with the list of used symbols.

#### 223 4.1. Reducing the time complexity of weak learner selection

224 From Eq. 10, the time-per-round requirements are  $O(|\mathcal{T}_i^{\mathbf{P}}| \cdot |\mathcal{T}_i^{\mathbf{n}}| \cdot |\mathcal{H}|)$ , where  
 225  $\mathcal{T}_i^{\mathbf{P}}$  is the subset of  $\mathcal{T}_i$  containing only positive samples  $\mathbf{p}$ ,  $\mathcal{T}_i^{\mathbf{n}}$  is the subset of  $\mathcal{T}_i$   
 226 containing only negative samples  $\mathbf{n}$  and  $\mathcal{H}$  is the set of candidate weak learners

227  $h(\cdot)$ . To improve this, as suggested in [8] we maintain only a one-argument  
 228 weight distribution  $v_{i,t}(\mathbf{x})$  under the constraint

$$w_{i,t}(\mathbf{p}, \mathbf{n}) = v_{i,t}(\mathbf{p})v_{i,t}(\mathbf{n}) \quad (11)$$

229 Then, omitting  $i, t$  subscripts we obtain from Eq. 9

$$\begin{aligned} r &= \sum_{\mathbf{p}, \mathbf{n}} w(\mathbf{p}, \mathbf{n}) (h(\mathbf{p}) - h(\mathbf{n})) \\ &= \sum_{\mathbf{p}} \sum_{\mathbf{n}} v(\mathbf{p})v(\mathbf{n}) (h(\mathbf{p})s(\mathbf{p}) + h(\mathbf{n})s(\mathbf{n})) \\ &= \sum_{\mathbf{p}} (v(\mathbf{p}) \sum_{\mathbf{n}} v(\mathbf{n})) h(\mathbf{p})s(\mathbf{p}) + \sum_{\mathbf{n}} (v(\mathbf{n}) \sum_{\mathbf{p}} v(\mathbf{p})) h(\mathbf{n})s(\mathbf{n}) \\ &= \sum_{\mathbf{x}} \pi(\mathbf{x})h(\mathbf{x}) \end{aligned} \quad (12)$$

230 where

$$s(\mathbf{x}) = \begin{cases} +1, & \text{if } \mathbf{x} \in \mathcal{T}_i^{\mathbf{p}} \\ -1, & \text{if } \mathbf{x} \in \mathcal{T}_i^{\mathbf{n}} \end{cases} \quad (13)$$

231 and

$$\pi(\mathbf{x}) = s(\mathbf{x})v(\mathbf{x}) \sum_{\mathbf{x}': s(\mathbf{x}') \neq s(\mathbf{x})} v(\mathbf{x}') \quad (14)$$

232 is referred to as the *potential*  $\pi(\mathbf{x})$  and can be precomputed at the beginning of  
 233 each round in only  $O(|\mathcal{T}_i^{\mathbf{p}}| + |\mathcal{T}_i^{\mathbf{n}}|) = O(|\mathcal{T}_i|)$  time. Combining Eqs. 2, 10 and  
 234 12 we obtain the final solution

$$h_{i,t} = \arg \max_h \left| \sum_{\mathbf{x}: \varphi(\mathbf{x}) > \theta} \pi_{i,t}(\mathbf{x})h(\mathbf{x}) \right| \quad (15)$$

235 which has the reduced time complexity  $O(|\mathcal{T}_i| \cdot |\mathcal{H}|)$ . Remarkably, this is the  
 236 same time complexity of the weak learner selection in the Viola-Jones standard  
 237 cascade framework.

#### 238 4.2. On-line efficient features computation

239 From Eq. 2, it follows that each candidate weak learner  $h(\mathbf{x}) \in \mathcal{H}$  relies  
 240 on the evaluation of a feature  $\varphi(\mathbf{x}) \in \mathcal{F}$  on a sample  $\mathbf{x}$  with respect to a

---

**Algorithm 1** Ranking-based cascade training
 

---

$F_{target}, F$ : desired and actual false positive rate of the cascade classifier  
 $f, f_{i,t}$ : desired and actual false positive rate of the  $i$ th node at round  $t$   
 $d$ : desired detection rate of each node  
 $\mathcal{P}, \mathcal{N}$ : sets of positive samples  $\mathbf{p}$  and negative samples  $\mathbf{n}$   
 $\mathcal{T}_i, \mathcal{V}_i$ : training and validation sets of the  $i$ th node  
 $\mathcal{T}_i^{\mathbf{x}}, \mathcal{V}_i^{\mathbf{x}}$ : subsets of  $\mathcal{T}_i, \mathcal{V}_i$  containing only positive ( $\mathbf{x}=\mathbf{p}$ ) or negative ( $\mathbf{x}=\mathbf{n}$ ) samples  
 $v_{i,t}(\mathbf{x})$ : weight of sample  $\mathbf{x}$  at round  $t$  of the  $i$ th node  
 $[pr]$ : defined to be 1 if  $pr$  holds and 0 otherwise (Iverson bracket)  
 $s(\mathbf{x})$ : function equal to +1 if  $\mathbf{x} = \mathbf{p}$  or to -1 if  $\mathbf{x} = \mathbf{n}$   
 $\varphi(\mathbf{x})$ : a feature evaluated on a sample  $\mathbf{x}$   
 $H_i(\mathbf{x})$ :  $i$ th node classifier  
 $\Theta_i$ : decision threshold of the  $i$ th node classifier  
 $h_{i,t}(\mathbf{x})$ : weak learner of the  $i$ th node at round  $t$  evaluated on the sample  $\mathbf{x}$   
 $\alpha_{i,t}$ : the weight of  $h_{i,t}(\mathbf{x})$

1:  $i \leftarrow 0$ ;  $F \leftarrow 1.0$ ;  
 2: **while**  $F > F_{target} \wedge \mathcal{N} \neq \{\emptyset\}$  **do** /\* add a new node \*/  
 3:    $i \leftarrow i + 1$ ;  $t \leftarrow 0$ ;  $f_{i,0} \leftarrow 1.0$   
 4:   precompute all candidate weak learners  $h(\cdot) \in \mathcal{H}$  on  $\mathcal{T}_i$   
 5:    $v_{i,0}(\mathbf{p}) = 1/|\mathcal{T}_i^{\mathbf{p}}|$ ;  $v_{i,0}(\mathbf{n}) = 1/|\mathcal{T}_i^{\mathbf{n}}|$  /\* initialize sample weights \*/  
 6:   **while**  $f_{i,t} > f \wedge \neg varstop$  **do** /\* add a new weak learner \*/  
 7:      $t \leftarrow t + 1$   
 8:      $\pi_{i,t}(\mathbf{x}) = s(\mathbf{x})v_{i,t}(\mathbf{x}) \sum_{\mathbf{x}':s(\mathbf{x}') \neq s(\mathbf{x})} v_{i,t}(\mathbf{x}')$  /\* precompute potentials \*/  
 9:      $h_{i,t} = \arg \max_h |r_{i,t}|$  /\* find best weak learner \*/  
     $\alpha_{i,t} = \frac{1}{2} \ln \left( \frac{1+r_{i,t}}{1-r_{i,t}} \right) \Big|_{h=h_{i,t}}$   
    where  $r_{i,t} = \sum_{\mathbf{x}:\varphi(\mathbf{x}) > \theta} \pi_{i,t}(\mathbf{x})h(\mathbf{x})$   
 10:      $v_{i,t+1}(\mathbf{x})|_{\mathbf{x}=\mathbf{p},\mathbf{n}} = \frac{v_{i,t}(\mathbf{x})e^{-s(\mathbf{x})\alpha_{i,t}h_{i,t}(\mathbf{x})}}{\sum_{\mathbf{x}} v_{i,t}(\mathbf{x})e^{-s(\mathbf{x})\alpha_{i,t}h_{i,t}(\mathbf{x})}}$  /\* sample weights update \*/  
 11:     precompute weak learner  $h_{i,t}$  on  $\mathcal{V}_i$  and  $\mathcal{N}$   
 12:      $H_{i,t}(\mathbf{x})|_{\mathbf{x} \in \mathcal{V}_i} = \sum_{\tau=1}^t \alpha_{i,\tau} h_{i,\tau}(\mathbf{x})$  /\* evaluate  $H_{i,t}(\cdot)$  on  $\mathcal{V}_i$  \*/  
 13:      $\Theta_i : |\{\mathbf{p} \in \mathcal{V}_i : H_{i,t}(\mathbf{p}) \geq \Theta_i\}| \geq d|\mathcal{V}_i^{\mathbf{p}}|$  /\* meet node learning goal  $d$  \*/  
 14:      $f_{i,t} \leftarrow f_{i,t-1} |\{\mathbf{n} \in \mathcal{V}_i : H_{i,t}(\mathbf{n}) \geq \Theta_i\}|$  /\* update false positive rate  $f_{i,t}$  \*/  
 15:      $varstop \leftarrow [Var(\{f_{i,\tau}\}_{\tau=t-\Delta, \dots, t}) \leq \epsilon]$  /\* stop criterion \*/  
 16:   **end while**  
 17:   output  $H_i(\mathbf{x}) = \left[ \sum_{\tau=1}^t \alpha_{i,\tau} h_{i,\tau}(\mathbf{x}) \geq \Theta_i \right]$  /\* the final  $i$ th node classifier \*/  
 18:    $F \leftarrow F \times f_{i,t}$  /\* update overall false positive rate \*/  
 19:   **if**  $F > F_{target}$  **then** /\* prepare  $\mathcal{T}_{i+1}$  and  $\mathcal{V}_{i+1}$  for a new node \*/  
 20:      $\mathcal{T}_{i+1} = \mathcal{T}_i \setminus \{\mathbf{n} \in \mathcal{T}_i : H_i(\mathbf{n}) = 0\}$ ;  $\mathcal{V}_{i+1} = \mathcal{V}_i \setminus \{\mathbf{n} \in \mathcal{V}_i : H_i(\mathbf{n}) = 0\}$   
 21:      $\mathcal{N} = \{\mathbf{n} \in \mathcal{N} : H_i(\mathbf{n}) = 1\}$   
 22:     refill adequately both  $\mathcal{T}_{i+1}$  and  $\mathcal{V}_{i+1}$  by random sampling from  $\mathcal{N}$   
 23:   **end if**  
 24: **end while**

---

241 threshold  $\theta$ . Being  $O(1)$  the per-feature evaluation time requirement [22], pre-  
 242 calculating the features for all the samples used in the training phase would  
 243 cost  $O(|\mathcal{F}|(|\mathcal{P}| + |\mathcal{N}|))$ . Moreover, additional costs have to be considered when  
 244 there is no available system memory to keep all these feature values at the same  
 245 time, and a caching strategy is used instead. To cope with these problems, we  
 246 draw an *on-line* features computation strategy aimed at minimizing the num-  
 247 ber of needed features calculations and keeping feature values into memory until  
 248 they are no longer needed. In fact, only three contributions are required to train  
 249 the  $i$ th node: (i)  $|\mathcal{T}_i| |\mathcal{F}|$  feature values for finding the best weak learner on  $\mathcal{T}_i$   
 250 (see Algorithm 1, line 9); (ii)  $t |\mathcal{V}_i|$  feature values for evaluating the node clas-  
 251 sifier  $H_{i,t}(\mathbf{x})$  on  $\mathcal{V}_i$  (see Algorithm 1, lines 12 and 20); and (iii)  $t |\mathcal{N}|$  feature  
 252 values for evaluating the  $i$ th node classifier on  $\mathcal{N}$  (see Algorithm 1, line 21).  
 253 Therefore, we can pre-compute and keep into memory the  $|\mathcal{T}_i| |\mathcal{F}|$  feature values  
 254 before training the  $i$ th node (see Algorithm 1, line 4) and do the same for each  
 255 selected feature evaluated on  $\mathcal{V}_i$  and  $\mathcal{N}$  at each round (see Algorithm 1, line 11).  
 256 When a node has been trained, we release the memory used for storing the fea-  
 257 ture values calculated on the samples discarded from  $\mathcal{T}_i$ ,  $\mathcal{V}_i$  and  $\mathcal{N}$ . In this way,  
 258 the time-per-node requirements for feature computation are only  $O(|\mathcal{T}_i| |\mathcal{F}|)$ ,  
 259 where we supposed  $t(|\mathcal{V}_i| + |\mathcal{N}|) \ll |\mathcal{T}_i| |\mathcal{F}|$  since  $|\mathcal{N}|$  decreases exponentially  
 260 throughout the cascade (see Algorithm 1, lines 21-22) and  $t$  can be supposed to  
 261 be a relatively small number (as it will be discussed in the next paragraph).

#### 262 4.3. Determining the number of training rounds

263 In principle, the training rounds for the  $i$ th node go on until the node learn-  
 264 ing goals  $d$  and  $f$  are met. In fact, the decision threshold  $\Theta_i$  of the  $i$ th node  
 265 classifier  $H_i(\mathbf{x})$  is always chosen at each round so as to meet  $d$ , hence the number  
 266 of rounds is governed only by whether the condition  $f_{i,t} \leq f$  is satisfied, being  
 267  $f_{i,t}$  the actual false positive rate of the  $i$ th node achieved at round  $t$ . How-  
 268 ever, when the classification task is getting more and more complex throughout  
 269 the cascade, it could be too difficult to satisfy such a condition, thus causing  
 270 several unnecessary features to be added without substantially reducing  $f_{i,t}$ .

271 To solve this problem, we add new features until a significant reduction of  $f_{i,t}$   
272 can be achieved. Let  $\psi_i(\Delta) = \{f_{i,\tau}\}_{\tau=t-\Delta,t-\Delta+1,\dots,t}$  be the latest  $\Delta$  achieved  
273 false positive rates of the  $i$ th node. Then, we define a stopping criterion being  
274  $varstop = [Var(\psi_i(\Delta)) \leq \epsilon]$  so that new features are added until the condition  
275  $f_{i,t} > f \wedge \neg varstop$  holds (see Algorithm 1, lines 6,15).

## 276 5. Experiments

277 To evaluate the performance of the proposed approach, we have considered  
278 three different real problems with highly unbalanced classes. The first is face de-  
279 tection that represents a well-known topic in scientific literature and a historical  
280 problem for the cascade approach as it was used in the seminal paper of Viola  
281 and Jones [22]. The other two problems belong to the field of medical image  
282 analysis. In particular, we examined the detection of microcalcifications (MC)  
283 on digital mammograms for the automated early detection of breast cancer, and  
284 the detection of microaneurysms (MA) on digital ocular fundus images, that is  
285 an important task in computer aided diagnosis of diabetic retinopathy.

286 In all the experiments, Haar-like features were used to describe the regions to  
287 be detected. These features are simple and computationally efficient and have  
288 been successfully used in face detection and classification problems. As in [10],  
289 we considered the set of (i) edge features; (ii) line features; (iii) center-surround  
290 features; (iv) and special diagonal line features. All features have been scaled  
291 and separately translated across all possible combinations on the subwindow,  
292 thus obtaining tens of thousands of features.

293 To verify the effectiveness of our approach (hereafter referred to as *Rank-*  
294 *ingCascade*) we also analyzed the behavior of other methods proposed in the  
295 literature. In particular, we implemented and evaluated the performance of  
296 different cascades that represent different solutions for facing class imbalance:

- 297 • *AsymBoost-based Cascade* (hereafter abbreviated as *Asym*). AsymBoost  
298 [23] is a cost-sensitive variant of AdaBoost that uses a parameter  $k$  to  
299 estimate the balance between the errors on the positive (False Negative

300 Rate, FNR) and the negative class (False Positive Rate, FPR). The value  
301 of  $k$  is usually chosen equal to the imbalance level, i.e., the ratio between  
302 the cardinality of the negative and the positive sets.

- 303 • *AsymBoost\** (hereafter abbreviated as *Asym\**). This method, proposed in  
304 [24], is a variant of *Asym* and introduces two principal changes: a dynamic  
305 optimization of the FPR in each stage of the cascade and a control of the  
306 FNR in each round so as to avoid an exponential increment of the false  
307 positives.
- 308 • *RUSBoost-based Cascade* (hereafter abbreviated as *RUS*). This approach  
309 is based on the random undersampling of the majority class. To this end,  
310 we implemented in each node of the cascade the RUSBoost classifier [18].  
311 The training data in input at each stage are undersampled so as to obtain  
312 a ratio of 35 : 65 (positive:negative).

313 For the sake of comparison, we have also considered the original Viola-Jones  
314 cascade detector (hereafter abbreviated as *Ada*), where the node stage is built  
315 with an AdaBoost dichotomizer.

316 The detectors have been evaluated in terms of Receiver Operating Charac-  
317 teristics (ROC) curve by plotting True Positive Rate (TPR) against FPR for  
318 a series of thresholds on the confidence degree associated to each sample. In  
319 this case not all the ROC curve is of interest since only low values of FPR are  
320 acceptable. For this reason, our analysis was focused on the initial portion of  
321 the curve and, subsequently, we considered as performance measure the *partial*  
322 *Area Under the ROC Curve* (pAUC) defined as  $pAUC(x) = \int_0^x TPR \, dFPR$   
323 where  $[0, x]$  is the range of interest for FPR [30].

324 To determine if *RankingCascade* performs significantly different from the  
325 other approaches we applied the bootstrap procedure [17]. The test set was  
326 sampled with replacement 2,000 times so that each new set of sampled data  
327 contained the same number of examples as the original set. We considered  
328 two different FPR ranges  $[0, 10^{-3}]$  and  $[0, 10^{-4}]$  which are close to practical  
329 application requirements for all the problems considered. For each range, the



Table 1: Performance differences between the compared detectors for a pAUC evaluated at  $10^{-3}$  on the face detection problem.

Methods	Mean Differences	p-value
<i>Ada</i>	$2.19 \cdot 10^{-5}$	< <b>0.001</b>
<i>Asym</i>	$-4.68 \cdot 10^{-6}$	0.842
<i>Asym</i> *	$6.42 \cdot 10^{-7}$	0.430
<i>RUS</i>	$2.05 \cdot 10^{-5}$	< <b>0.001</b>

Table 2: Performance differences between the compared detectors for a pAUC evaluated at  $10^{-4}$  on the face detection problem.

Methods	Mean Differences	p-value
<i>Ada</i>	$6.33 \cdot 10^{-6}$	< <b>0.001</b>
<i>Asym</i>	$3.29 \cdot 10^{-6}$	<b>0.001</b>
<i>Asym</i> *	$1.07 \cdot 10^{-6}$	0.081
<i>RUS</i>	$1.98 \cdot 10^{-5}$	< <b>0.001</b>

330 differences in pAUC between *RankingCascade* and the other detectors were  
 331 computed. Resampling 2,000 times resulted in 2,000 values for each performance  
 332 difference. To compare the performance of our approach against the other 4  
 333 methods we evaluated the *p*-value, defined as the fraction of the corresponding  
 334 pAUC differences that were negative or zero. The statistical significance level  
 335 was  $\alpha = 0.01$  but, due to the number of comparisons, we applied the Bonferroni  
 336 correction [6] and thus, performance differences were considered significant if  
 337  $p < 0.0025$ , (i.e.,  $p < \alpha/4$ ).

### 338 5.1. Face Detection

339 Face detection is the first analyzed problem. Following a common procedure  
 340 in the literature [22, 23, 28], two different data sets have been created to train  
 341 and test the cascades. The positives samples for the training phase were 1,600  
 342 subwindows extracted from the frontal faces of the well-known FERET database  
 343 [14] and scaled in a standard format of  $24 \times 24$  pixels. The negative samples  
 344 were 22,906,769 subwindows, all of size  $24 \times 24$ , taken from the 9,028 images  
 345 of a publicly available non-face dataset [28] at different scales. The positive  
 346 samples were equally parted between training and validation set. The negative  
 347 samples were divided in three data sets: training set, validation set and pool  
 348 respectively of 20,000, 60,000 and 22,826,769 samples. As test set we adopted

349 the MIT+CMU database [20, 16] that consists of 180 images containing 734  
350 faces. The negative set has been taken from 1,417 images randomly extracted  
351 from different categories (abbey, forests, greenhouse, shipyard, skyscraper) of  
352 the publicly available SUN database [29] for a total of 21,181,195 samples.

353 The cascade detectors were built using  $d = 0.999$  and  $f = 0.300$ . The  
354 training stage produced 4 nodes for all the cascades except *Asym\** that was  
355 composed by 6 stages. The total number of features considered at each stage  
356 was 116,544.

357 Results of the comparisons between detectors are reported in Tables 1 and  
358 2 for the two FPRs considered in the pAUC evaluation. In both tables, the  
359 second column shows the mean difference between the pAUC of the proposed  
360 approach and the compared detectors while for each comparison the  $p$ -value is  
361 given in the third column. Statistically significant differences are listed in bold.

362 In Table 1 the performance of the proposed method is compared to the  
363 other cascades for an FPR lower than  $10^{-3}$ . *RankingCascade* reveals to be  
364 significantly better than *Ada* and *RUS*. The differences, instead, are not statis-  
365 tically significant when compared to *Asym\** and *Asym*. When looking at Table  
366 2, i.e., for a pAUC evaluated between 0 and  $10^{-4}$ , results are even better since  
367 the performance of our method becomes statistically higher than those of *Asym*.  
368 The difference with *Asym\**, instead, is again not statistically significant. As a  
369 final remark, it is worth noting that *Asym* is a parametric approach whereas  
370 *Asym\**, using a dynamic optimization of the FPR in each stage, tends to a  
371 classification system with a higher number of stages and a higher number of  
372 features per stage, so increasing the computational complexity of the cascade.

### 373 5.2. Microcalcification Detection

374 The second experiment deals with the problem of detecting microcalcifica-  
375 tions on digital mammograms. Microcalcifications appear as bright small cir-  
376 cular spots in the image and represent a subtle sign of breast cancer in women.  
377 A private full-field digital mammographic database of 198 images has been ex-  
378 ploited to extract 8,000 microcalcifications (positive samples) and 22,760,320

Table 3: Performance differences between the compared detectors for a pAUC evaluated at  $10^{-3}$  on the microcalcifications detection problem.

Methods	Mean Differences	p-value
<i>Ada</i>	$2.82 \cdot 10^{-5}$	< <b>0.001</b>
<i>Asym</i>	$8.81 \cdot 10^{-6}$	< <b>0.001</b>
<i>Asym</i> *	$2.74 \cdot 10^{-5}$	< <b>0.001</b>
<i>RUS</i>	$9.52 \cdot 10^{-5}$	< <b>0.001</b>

Table 4: Performance differences between the compared detectors for a pAUC evaluated at  $10^{-4}$  on the microcalcifications detection problem.

Methods	Mean Differences	p-value
<i>Ada</i>	$5.61 \cdot 10^{-6}$	< <b>0.001</b>
<i>Asym</i>	$3.51 \cdot 10^{-6}$	< <b>0.001</b>
<i>Asym</i> *	$4.64 \cdot 10^{-6}$	< <b>0.001</b>
<i>RUS</i>	$1.69 \cdot 10^{-5}$	< <b>0.001</b>

379 background regions (negative class), all of size  $12 \times 12$ . Ten-fold cross vali-  
 380 dation has been performed considering 9 folds for training set, validation set  
 381 and pool and the remaining fold as test set. In each cross validation step we  
 382 used 7,200 positive samples equally parted between training and validation sets  
 383 and 20,484,351 negative samples subdivided in 20,000, 60,000 and 20,404,351  
 384 respectively for training set, validation set and pool.

385 The cascade detectors were built using  $d = 0.99$  and  $f = 0.30$ . The training  
 386 stage consists in ten different cross validation runs and thus, for each method,  
 387 ten different cascades were obtained. The number of stages for the *RankingCas-*  
 388 *cade* was equal to 5 for all the runs while for the other approaches varied among  
 389 4 and 6 according to the cross validation step considered. The total number of  
 390 features considered at each stage was 14,709.

391 The results of the comparisons for the two different FPRs used in pAUC  
 392 evaluation (i.e.,  $10^{-3}$  and  $10^{-4}$ ) are reported respectively in Tables 3 and 4. In  
 393 both cases the *RankingCascade* exhibits higher performance than all the other  
 394 approaches even with a high statistical significance as proved by the very low  
 395 *p*-values.

396 The effectiveness of the approach on this particular problem has been also  
 397 confirmed by the results obtained in [1], where it has been employed in a

Table 5: Performance differences between the compared detectors for a pAUC evaluated at  $10^{-3}$  on the microaneurysm detection problem.

Methods	Mean Differences	p-value
<i>Ada</i>	$2.82 \cdot 10^{-5}$	< <b>0.001</b>
<i>Asym</i>	$8.81 \cdot 10^{-6}$	< <b>0.001</b>
<i>Asym*</i>	$2.74 \cdot 10^{-5}$	< <b>0.001</b>
<i>RUS</i>	$9.52 \cdot 10^{-5}$	< <b>0.001</b>

Table 6: Performance differences between the compared detectors for a pAUC evaluated at  $10^{-4}$  on the microaneurysm detection problem.

Methods	Mean Differences	p-value
<i>Ada</i>	$5.61 \cdot 10^{-6}$	< <b>0.001</b>
<i>Asym</i>	$3.51 \cdot 10^{-6}$	< <b>0.001</b>
<i>Asym*</i>	$4.64 \cdot 10^{-6}$	< <b>0.001</b>
<i>RUS</i>	$1.69 \cdot 10^{-5}$	< <b>0.001</b>

398 *Computer-Aided Detection and Diagnosis* (CAD) system that revealed to be  
 399 competitive with the state-of-the-art commercial CAD systems.

### 400 5.3. Microaneurysm Detection

401 The third experiment deals with the problem of detecting microaneurysms  
 402 on digital ocular fundus images. Microaneurysms appear as a dark small cir-  
 403 cular spots in the image and represent a subtle sign of retinopathy. A pub-  
 404 lic database [15] of 50 digital ocular fundus images has been used to extract  
 405 2,890,972 patches of size  $15 \times 15$  pixels, 1,997 containing microaneurysms (pos-  
 406 itive samples) and 2,888,975 containing background tissue (negative samples).  
 407 Ten-fold cross validation has been performed considering 9 folds for training  
 408 set, validation set and pool and the remaining fold as test set. In each cross  
 409 validation step we used 1,797 positive samples equally parted between training  
 410 and validation sets and 2,620,077 negative samples subdivided in 10,000, 10,000  
 411 and 2,600,077 respectively for training set, validation set and pool.

412 The cascade detectors were built using  $d = 0.95$  and  $f = 0.30$ . The training  
 413 stage consists in ten different cross validation runs and thus, for each method,  
 414 ten different cascades were obtained. The number of stages for all the methods  
 415 considered varied among 14 and 16 according to the cross validation step con-  
 416 sidered, except for *RUSBoost* where it varied among 3 and 5. The total number

Table 7: Average execution time (in seconds) for selecting a weak learner

	Face Detection	MC Detection	MA Detection
<b><i>RankingCascade</i></b>	<b>22.6</b>	<b>3.5</b>	<b>3.7</b>
<i>Ada</i>	104.7	9.1	12.4
<i>Asym</i>	89.6	11.6	15.8
<i>Asym*</i>	92.2	10.6	15.4
<i>RUS</i>	19.1	3.1	3.3

Table 8: Average execution time (in seconds) for building a node classifier

	Face Detection	MC Detection	MA Detection
<b><i>RankingCascade</i></b>	<b>360.0</b>	<b>49.4</b>	<b>103.0</b>
<i>Ada</i>	830.7	107.0	306.2
<i>Asym</i>	1030.1	135.9	400.3
<i>Asym*</i>	1625.6	150.9	509.2
<i>RUS</i>	248.4	34.1	67.5

417 of features considered at each stage was 43,172.

418 The results of the comparisons for the two different FPRs used in pAUC  
 419 evaluation (i.e.,  $10^{-3}$  and  $10^{-4}$ ) are reported respectively in Tables 5 and 6. In  
 420 both cases the *RankingCascade* exhibits higher performance than all the other  
 421 approaches even with a high statistical significance as proved by the very low  
 422 *p*-values.

#### 423 5.4. Computational time

424 To confirm the effectiveness of the implementation proposed in Section 4, for  
 425 each performed experiment we report in Table 7 the average execution time for  
 426 selecting a weak learner and in Table 8 the average execution time for building  
 427 a node classifier. All the considered methods have been implemented in C++  
 428 with multi-threading and run on a workstation equipped with two Intel Xeon  
 429 E5520 and 96.0 GB of RAM. Remarkably, *RankingCascade* performed faster  
 430 than all the other learning strategies except *RUS*, which takes advantage of  
 431 undersampling to drastically reduce the number of calculations (at the cost of  
 432 worse detection performance, as shown before).

## 433 6. Conclusions

434 In this paper we have addressed the problem of class imbalance in the node  
435 classifier learning within the standard Viola-Jones cascade framework and pro-  
436 posed a new learning strategy aimed at maximizing the node classifiers ranking  
437 capability rather than their accuracy. Such approach revealed to be an effective  
438 solution to the class asymmetry problem and provided good performance in ex-  
439 periments when compared with other approaches. In particular, when tested on  
440 three real-world severe detection problems, our learning strategy provided simi-  
441 lar or better results than methods, such as *AsymBoost* and *AsymBoost\**, which  
442 rely on a cost-based approach to face class imbalance. In this regard, it is worth  
443 mentioning that our method (as well as *AsymBoost\**) does not introduce any  
444 further parameter to the original cascade whereas *AsymBoost* requires a cost  
445 ratio to be specified. In the same experiments, the proposed approach showed  
446 significantly better results also when compared with a cascade using *RUSBoost*,  
447 so demonstrating that the ranking-based strategy in node classifier learning  
448 performs better than undersampling. In summary, the experimental results  
449 supported the rationale on which our method is based, i.e., that ranking-based  
450 learning is effective in facing class imbalance and allows the construction of pro-  
451 ficient cascade detectors. Possible directions for future work include evaluating  
452 the effectiveness of the proposed approach also in other cascade architectures  
453 which extend the standard Viola-Jones framework.

## 454 References

- 455 [1] A. Bria, N. Karssemeijer, F. Tortorella, Learning from unbalanced data: A  
456 cascade-based approach for detecting clustered microcalcifications, *Medical*  
457 *Image Analysis* 18 (2) (2014) 241–252.
- 458 [2] A. Bria, C. Marrocco, M. Molinara, F. Tortorella, A ranking-based cascade  
459 approach for unbalanced data, in: *Pattern Recognition (ICPR)*, 2012 21st  
460 International Conference on, 2012, pp. 3439–3442.

- 461 [3] R. Caruana, Learning from Imbalanced Data: Rank Metrics and Extra  
462 Tasks, in: The AAAI Workshop on Learning from Imbalanced Data Sets,  
463 2000, pp. 51–57.
- 464 [4] N. V. Chawla, K. W. Bowyer, L. O. Hall, W. P. Kegelmeyer, SMOTE: Syn-  
465 thetic Minority Over-sampling Technique, *Journal of Artificial Intelligence*  
466 *Research* 16 (1) (2002) 321–357.
- 467 [5] N. V. Chawla, A. Lazarevic, L. O. Hall, K. W. Bowyer, SMOTEBoost:  
468 Improving Prediction of the Minority Class in Boosting, in: N. Lavrač,  
469 D. Gamberger, L. Todorovski, H. Blockeel (eds.), *Knowledge Discovery in*  
470 *Databases: PKDD 2003*, vol. 2838 of *Lecture Notes in Computer Science*,  
471 Springer Berlin Heidelberg, 2003, pp. 107–119.
- 472 [6] O. J. Dunn, Multiple Comparisons Among Means, *Journal of the American*  
473 *Statistical Association* 56 (293) (1961) 52–64.
- 474 [7] W. Fan, S. J. Stolfo, J. Zhang, P. K. Chan, AdaCost: Misclassification  
475 Cost-Sensitive Boosting, in: *Proceedings of the Sixteenth International*  
476 *Conference on Machine Learning, ICML '99*, Morgan Kaufmann Publish-  
477 ers Inc., San Francisco, CA, USA, 1999, pp. 97–105.
- 478 [8] Y. Freund, R. Iyer, R. E. Schapire, Y. Singer, An efficient boosting algo-  
479 rithm for combining preferences, *Journal of Machine Learning Research* 4  
480 (2003) 933–969.
- 481 [9] Y. Freund, R. E. Schapire, A Decision-Theoretic Generalization of On-Line  
482 Learning and an Application to Boosting, *Journal of Computer and System*  
483 *Sciences* 55 (1) (1997) 119–139.
- 484 [10] R. Lienhart, E. Kuranov, V. Pisarevsky, Empirical Analysis of Detection  
485 Cascades of Boosted Classifiers for Rapid Object Detection, in: *DAGM*  
486 *25th Pattern Recognition Symposium, 2003*, pp. 297–304.
- 487 [11] H. Masnadi-Shirazi, N. Vasconcelos, Cost-Sensitive Boosting, *Pattern Anal-*  
488 *ysis and Machine Intelligence, IEEE Transactions on* 33 (2) (2011) 294–309.

- 489 [12] M. Molinara, M. Ricamato, F. Tortorella, Facing Imbalanced Classes  
490 through Aggregation of Classifiers, in: *Image Analysis and Processing,*  
491 2007. ICIAP 2007. 14th International Conference on, 2007, pp. 43–48.
- 492 [13] M.-T. Pham, V.-D. Hoang, T.-J. ChamP, Detection with multi-exit asym-  
493 metric boosting, in: *Computer Vision and Pattern Recognition,* 2008.  
494 CVPR 2008. IEEE Conference on, 2008, pp. 1–8.
- 495 [14] P. Phillips, H. Moon, S. Rizvi, P. Rauss, The FERET evaluation method-  
496 ology for face-recognition algorithms, *Pattern Analysis and Machine Intel-*  
497 *ligence, IEEE Transactions on* 22 (10) (2000) 1090–1104.
- 498 [15] P. Prentasic, S. Loncaric, Z. Vatavuk, G. Bencic, M. Subasic, T. Petkovic,  
499 L. Dujmovic, M. Malenica-Ravlic, N. Budimlija, R. Tadic, Diabetic  
500 Retinopathy Image Database (DRiDB): a new database for diabetic  
501 retinopathy screening programs research, in: *Image and Signal Process-*  
502 *ing and Analysis (ISPA), 8th International Symposium on, IEEE,* 2013,  
503 pp. 711–716.
- 504 [16] H. Rowley, S. Baluja, T. Kanade, Neural network-based face detection,  
505 *Pattern Analysis and Machine Intelligence, IEEE Transactions on* 20 (1)  
506 (1998) 23–38.
- 507 [17] F. W. Samuelson, N. Petrick, Comparing image detection algorithms us-  
508 ing resampling, in: *Biomedical Imaging: Nano to Macro, 2006. 3rd IEEE*  
509 *International Symposium on,* 2006, pp. 1312–1315.
- 510 [18] C. Seiffert, T. Khoshgoftaar, J. Van Hulse, A. Napolitano, RUSBoost: A  
511 Hybrid Approach to Alleviating Class Imbalance, *Systems, Man and Cy-*  
512 *bernetics, Part A: Systems and Humans, IEEE Transactions on* 40 (1)  
513 (2010) 185–197.
- 514 [19] Y. Sun, A. Wong, Y. Wang, Parameter Inference of Cost-Sensitive Boosting  
515 Algorithms, in: P. Perner, A. Imiya (eds.), *Machine Learning and Data*



- 516 Mining in Pattern Recognition, vol. 3587 of Lecture Notes in Computer  
517 Science, Springer Berlin Heidelberg, 2005, pp. 21–30.
- 518 [20] K.-K. Sung, T. Poggio, Example-based learning for view-based human face  
519 detection, *Pattern Analysis and Machine Intelligence*, IEEE Transactions  
520 on 20 (1) (1998) 39–51.
- 521 [21] K. M. Ting, A Comparative Study of Cost-Sensitive Boosting Algorithms,  
522 in: *Proceedings of the Seventeenth International Conference on Machine  
523 Learning, ICML '00*, Morgan Kaufmann Publishers Inc., San Francisco,  
524 CA, USA, 2000, pp. 983–990.
- 525 [22] P. Viola, M. Jones, Robust Real-Time Object Detection, *Int. J. of Comp.*  
526 *Vis.* 57 (2) (2001) 137–154.
- 527 [23] P. Viola, M. Jones, Fast and Robust Classification using Asymmetric Ad-  
528 aBoost and a Detector Cascade, *Adv. in Neur. Inf. Proc. Syst.* 16 (2002)  
529 1311–1318.
- 530 [24] I. Visentini, C. Micheloni, G. L. Foresti, Reactive Learning Strategy for  
531 AsymBoost Based Face Detectors, *14th International Conference on Image  
532 Analysis and Processing* (2007) 357–362.
- 533 [25] P. Wang, C. Shen, H. Zheng, Z. Ren, Training a multi-exit cascade with  
534 linear asymmetric classification for efficient object detection, in: *Image  
535 Processing (ICIP)*, 2010 17th IEEE International Conference on, 2010, pp.  
536 61–64.
- 537 [26] G. Weiss, F. Provost, Learning When Training Data are Costly: The Effect  
538 of Class Distribution on Tree Induction, *Journal of Artificial Intelligence  
539 Research* 19 (2003) 315–354.
- 540 [27] G. M. Weiss, Mining with Rarity: A Unifying Framework, *SIGKDD Explor.*  
541 *News.* 6 (1) (2004) 7–19.

- 542 [28] J. Wu, S. Brubaker, M. Mullin, J. Rehg, Fast Asymmetric Learning for  
543 Cascade Face Detection, *Pattern Analysis and Machine Intelligence, IEEE*  
544 *Transactions on* 30 (3) (2008) 369–382.
- 545 [29] J. Xiao, J. Hays, K. Ehinger, A. Oliva, A. Torralba, SUN database: Large-  
546 scale scene recognition from abbey to zoo, in: *Computer Vision and Pattern*  
547 *Recognition (CVPR), 2010 IEEE Conference on*, 2010, pp. 3485–3492.
- 548 [30] W. A. Yousef, Assessing classifiers in terms of the partial area under the  
549 ROC curve, *Computational Statistics & Data Analysis* 64 (2013) 51–70.