

Training Effective Node Classifiers for Cascade Classification

Chunhua Shen · Peng Wang · Sakrapee Paisitkriangkrai · Anton van den Hengel

December 2012

Abstract Cascade classifiers are widely used in real-time object detection. Different from conventional classifiers that are designed for a low overall classification error rate, a classifier in each node of the cascade is required to achieve an extremely high detection rate and moderate false positive rate. Although there are a few reported methods addressing this requirement in the context of object detection, there is no principled feature selection method that explicitly takes into account this asymmetric node learning objective. We provide such an algorithm here. We show that a special case of the biased minimax probability machine has the same formulation as the linear asymmetric classifier (LAC) of Wu et al. (2005). We then design a new boosting algorithm that directly optimizes the cost function of LAC. The resulting totally-corrective boosting algorithm is implemented by the column generation technique in convex optimization. Experimental results on object detection verify the effectiveness of the proposed boosting algorithm as a node classifier in cascade object detection, and show performance better than that of the current state-of-the-art.

Keywords AdaBoost · Minimax Probability Machine · Cascade Classifier · Object Detection · Human Detection

1 Introduction

Real-time object detection inherently involves searching a large number of candidate image regions for a small number of objects. Processing a single image, for example, can require the interrogation of well over a million scanned windows in order to uncover a single correct detection. This imbalance in the data has an impact on the way that detectors are applied, but also on the training process. This impact is reflected in the need to identify discriminative features from within a large over-complete feature set.

Cascade classifiers have been proposed as a potential solution to the problem of imbalance in the data (Viola and Jones 2004; Bi et al. 2006; Dundar and Bi 2007; Brubaker et al. 2008; Wu et al. 2008), and have received significant attention due to their speed and accuracy. In this work, we propose a principled method by which to train a *boosting*-based cascade of classifiers.

The boosting-based cascade approach to object detection was introduced by Viola and Jones (Viola and Jones 2004; 2002), and has received significant subsequent attention (Li and Zhang 2004; Pham and Cham 2007b; Pham et al. 2008; Paisitkriangkrai et al. 2008; Shen et al. 2008; Paisitkriangkrai et al. 2009). It also underpins the current state-of-the-art (Wu et al. 2005; 2008).

The Viola and Jones approach uses a cascade of increasingly complex classifiers, each of which aims to achieve the best possible classification accuracy while achieving an extremely low false negative rate. These classifiers can be seen as forming the nodes of a degenerate binary tree (see Fig. 1) whereby a negative result from any single such node classifier terminates the interrogation of the current patch. Viola and Jones use AdaBoost to train each node classifier in order to

C. Shen (✉) · P. Wang · S. Paisitkriangkrai · A. van den Hengel
Australian Centre for Visual Technologies, and School of Computer Science, The University of Adelaide, SA 5005, Australia
E-mail: chunhua.shen@adelaide.edu.au
This work was in part supported by Australian Research Council Future Fellowship FT120100969.

achieve the best possible classification accuracy. A low false negative rate is achieved by subsequently adjusting the decision threshold until the desired false negative rate is achieved. This process cannot be guaranteed to produce the best detection performance for a given false negative rate.

Under the assumption that each node of the cascade classifier makes independent classification errors, the detection rate and false positive rate of the entire cascade are: $F_{dr} = \prod_{t=1}^N d_t$ and $F_{fp} = \prod_{t=1}^N f_t$, respectively, where d_t represents the detection rate of classifier t , f_t the corresponding false positive rate and N the number of nodes. As pointed out in (Viola and Jones 2004; Wu et al. 2005), these two equations suggest a *node learning objective*: Each node should have an extremely high detection rate d_t (e.g., 99.7%) and a moderate false positive rate f_t (e.g., 50%). With the above values of d_t and f_t , and a cascade of $N = 20$ nodes, then $F_{dr} \approx 94\%$ and $F_{fp} \approx 10^{-6}$, which is a typical design goal.

One drawback of the standard AdaBoost approach to boosting is that it does not take advantage of the cascade classifier’s special structure. AdaBoost only minimizes the overall classification error and does not particularly minimize the number of false negatives. In this sense, the features selected by AdaBoost are not optimal for the purpose of rejecting as many negative examples as possible. Viola and Jones proposed a solution to this problem in AsymBoost (Viola and Jones 2002) (and its variants (Pham and Cham 2007b; Pham et al. 2008; Wang et al. 2012; Masnadi-Shirazi and Vasconcelos 2007)) by modifying the loss function so as to more greatly penalize false negatives. AsymBoost achieves better detection rates than AdaBoost, but still addresses the node learning goal *indirectly*, and cannot be guaranteed to achieve the optimal solution.

Wu et al. explicitly studied the node learning goal and proposed to use linear asymmetric classifier (LAC) and Fisher linear discriminant analysis (LDA) to adjust the weights on a set of features selected by AdaBoost or AsymBoost (Wu et al. 2005; 2008). Their experiments indicated that with this post-processing technique the node learning objective can be better met, which is translated into improved detection rates. In Viola and Jones’ framework, boosting is used to select features and at the same time to train a strong classifier. Wu et al.’s work separates these two tasks: AdaBoost or AsymBoost is used to select features; and as a second step, LAC or LDA is used to construct a strong classifier by adjusting the weights of the selected features. The node learning objective is only considered at the second step. At the first step—feature selection—the node learning objective is not explicitly considered

at all. We conjecture that *further improvement may be gained if the node learning objective is explicitly taken into account at both steps*. We thus propose new boosting algorithms to implement this idea and verify this conjecture. A preliminary version of this work was published in Shen et al. (2010).

Our major contributions are as follows.

1. Starting from the theory of minimax probability machines (MPMs), we derive a simplified version of the biased minimax probability machine, which has the same formulation as the linear asymmetric classifier of Wu et al. (2005). We thus show the underlying connection between MPM and LAC. Importantly, this new interpretation weakens some of the restrictions on the acceptable input data distribution imposed by LAC.
2. We develop new boosting-like algorithms by directly minimizing the objective function of the linear asymmetric classifier, which results in an algorithm that we label LACBoost. We also propose FisherBoost on the basis of Fisher LDA rather than LAC. Both methods may be used to identify the feature set that optimally achieves the node learning goal when training a cascade classifier. To our knowledge, this is the first attempt to design such a feature selection method.
3. LACBoost and FisherBoost share similarities with LPBoost (Demiriz et al. 2002) in the sense that both use column generation—a technique originally proposed for large-scale linear programming (LP). Typically, the Lagrange dual problem is solved at each iteration in column generation. We instead solve the primal quadratic programming (QP) problem, which has a special structure and entropic gradient (EG) can be used to solve the problem very efficiently. Compared with general interior-point based QP solvers, EG is much faster.
4. We apply LACBoost and FisherBoost to object detection and better performance is observed over other methods (Wu et al. 2005; 2008; Maji et al. 2008). In particular on pedestrian detection, FisherBoost achieves the state-of-the-art, comparing with methods listed in (Dollár et al. 2012) on three benchmark datasets. The results confirm our conjecture and show the effectiveness of LACBoost and FisherBoost. These methods can be immediately applied to other asymmetric classification problems.

Moreover, we analyze the condition that makes the validity of LAC, and show that the multi-exit cascade might be more suitable for applying LAC learning of Wu et al. (2005) and Wu et al. (2008) (and our LACBoost) rather than Viola-Jones’ conventional cascade.

As observed in Wu et al. (2008), in many cases, LDA even performs better than LAC. In our experiments, we have also observed similar phenomena. Paisitkriangkrai et al. (2009) empirically showed that LDA’s criterion can be used to achieve better detection results. An explanation of why LDA works so well for object detection is missing in the literature. Here we demonstrate that in the context of object detection, LDA can be seen as a regularized version of LAC in approximation.

The proposed LACBoost/FisherBoost algorithm differs from traditional boosting algorithms in that it does not minimize a loss function. This opens new possibilities for designing boosting-like algorithms for special purposes. We have also extended column generation for optimizing nonlinear optimization problems. Next we review related work in the context of real-time object detection using cascade classifiers.

1.1 Related Work

The field of object detection has made a significant progress over the last decade, especially after the seminal work of Viola and Jones. Three key components that contribute to their first robust *real-time* object detection framework are:

1. The cascade classifier, which efficiently filters out negative patches in early nodes while maintaining a very high detection rate;
2. AdaBoost that selects informative features and at the same time trains a strong classifier;
3. The use of integral images, which makes the computation of Haar features extremely fast.

This approach has received significant subsequent attention. A number of alternative cascades have been developed including the soft cascade (Bourdev and Brandt 2005), WaldBoost (Sochman and Matas 2005), the dynamic cascade (Xiao et al. 2007), the AND-OR cascade (Dundar and Bi 2007), the multi-exit cascade (Pham et al. 2008), the joint cascade (Lefakis and Fleuret 2010) and recently proposed, the rate constraint embedded cascade (RCECBoost) (Saberian and Vasconcelos 2012). In this work we have adopted the multi-exit cascade of Pham et al. due to its effectiveness and efficiency as demonstrated in Pham et al. (2008). The multi-exit cascade improves classification performance by using the results of all of the weak classifiers applied to a patch so far in reaching a decision at each node of the tree (see Fig. 1). Thus the n -th node classifier uses the results of the weak classifiers associated with node n , but also those associated with the previous $n - 1$ node classifiers in the cascade. We show below that LAC post-processing can enhance the multi-exit cascade, and that

the multi-exit cascade more accurately fulfills the LAC requirement that the margin be drawn from a Gaussian distribution.

In addition to improving the cascade structure, a number of improvements have been made on the learning algorithm for building node classifiers in a cascade. Wu et al., for example, use fast forward feature selection to accelerate the training procedure (Wu et al. 2003). Wu et al. (2005) also showed that LAC may be used to deliver better classification performance. Pham and Cham recently proposed online asymmetric boosting that considerably reduces the training time required (Pham and Cham 2007b). By exploiting the feature statistics, Pham and Cham (2007a) have also designed a fast method to train weak classifiers. Li and Zhang (2004) proposed FloatBoost, which discards redundant weak classifiers during AdaBoost’s greedy selection procedure. Masnadi-Shirazi and Vasconcelos (2011) proposed cost-sensitive boosting algorithms which can be applied to different cost-sensitive losses by means of gradient descent. Liu and Shum (2003) also proposed KLBoost, aiming to select features that maximize the projected Kullback-Leibler divergence and select feature weights by minimizing the classification error. Promising results have also been reported by LogitBoost (Tuzel et al. 2008) that employs the logistic regression loss, and GentleBoost (Torralba et al. 2007) that uses adaptive Newton steps to fit the additive model. Multi-instance boosting has been introduced to object detection (Viola et al. 2005; Dollár et al. 2008; Lin et al. 2009), which does not require precisely labeled locations of the targets in training data.

New features have also been designed for improving the detection performance. Viola and Jones’ Haar features are not sufficiently discriminative for detecting more complex objects like pedestrians, or multi-view faces. Covariance features (Tuzel et al. 2008) and histogram of oriented gradients (HOG) (Dalal and Triggs 2005) have been proposed in this context, and efficient implementation approaches (along the lines of integral images) are developed for each. Shape context, which can also exploit integral images (Aldavert et al. 2010), was applied to human detection in thermal images (Wang et al. 2010). The local binary pattern (LBP) descriptor and its variants have been shown promising performance on human detection (Mu et al. 2008; Zheng et al. 2010). Recently, effort has been spent on combining complementary features, including: simple concatenation of HOG and LBP (Wang et al. 2007), combination of heterogeneous local features in a boosted cascade classifier (Wu and Nevatia 2008), and Bayesian integration of intensity, depth and motion features in a mixture-of-experts model (Enzweiler et al. 2010).

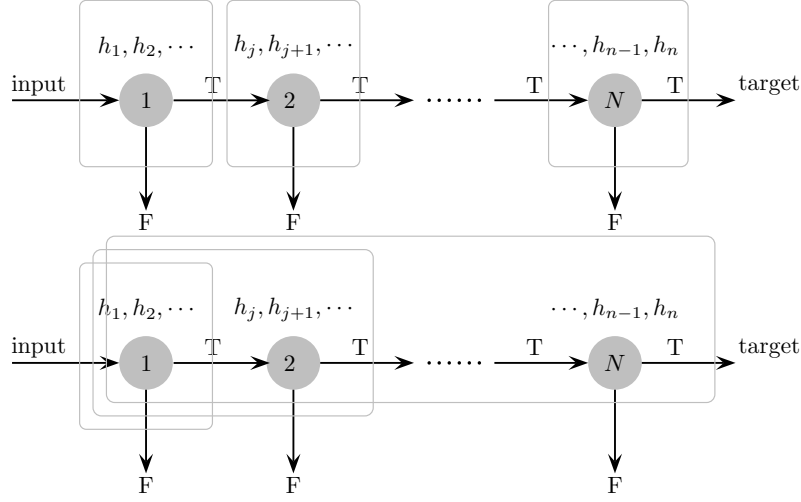


Fig. 1: Cascade classifiers. The first one is the standard cascade of Viola and Jones (2004). The second one is the multi-exit cascade proposed in Pham et al. (2008). Only those classified as true detection by all nodes will be true targets.

The rest of the paper is organized as follows. We briefly review the concept of minimax probability machine and derive the new simplified version of biased minimax probability machine in Section 2. Linear asymmetric classification and its connection to the minimax probability machine is discussed in Section 3. In Section 4, we show how to design new boosting algorithms (LACBoost and FisherBoost) by rewriting the optimization formulations of LAC and Fisher LDA. The new boosting algorithms are applied to object detection in Section 5 and we conclude the paper in Section 6.

1.2 Notation

The following notation is used. A matrix is denoted by a bold upper-case letter (\mathbf{X}); a column vector is denoted by a bold lower-case letter (\mathbf{x}). The i th row of \mathbf{X} is denoted by $\mathbf{X}_{i:}$ and the i th column $\mathbf{X}_{:,i}$. The identity matrix is \mathbf{I} and its size should be clear from the context. $\mathbf{1}$ and $\mathbf{0}$ are column vectors of 1's and 0's, respectively. We use \succ, \preceq to denote component-wise inequalities.

Let $\mathcal{T} = \{(\mathbf{x}_i, y_i)\}_{i=1, \dots, m}$ be the set of training data, where $\mathbf{x}_i \in \mathcal{X}$ and $y_i \in \{-1, +1\}$, $\forall i$. The training set consists of m_1 positive training points and m_2 negative ones; $m_1 + m_2 = m$. Let $h(\cdot) \in \mathcal{H}$ be a weak classifier that projects an input vector \mathbf{x} into $\{-1, +1\}$. Note that here we consider only classifiers with discrete outputs although the developed methods can use real-valued weak classifiers too. We assume that \mathcal{H} , the set from which $h(\cdot)$ is selected, is finite and has n elements.

Define the matrix $\mathbf{H}^{\mathcal{Z}} \in \mathbb{R}^{m \times n}$ such that the (i, j) entry $\mathbf{H}_{ij}^{\mathcal{Z}} = h_j(\mathbf{x}_i)$ is the label predicted by weak classifier $h_j(\cdot)$ for the datum \mathbf{x}_i , where \mathbf{x}_i the i th element of the set \mathcal{Z} . In order to simplify the notation we eliminate the superscript when \mathcal{Z} is the training set, so $\mathbf{H}^{\mathcal{Z}} = \mathbf{H}$. Therefore, each column $\mathbf{H}_{:,j}$ of the matrix \mathbf{H} consists of the output of weak classifier $h_j(\cdot)$ on all the training data; while each row $\mathbf{H}_{i:}$ contains the outputs of all weak classifiers on the training datum \mathbf{x}_i . Define similarly the matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ such that $\mathbf{A}_{ij} = y_i h_j(\mathbf{x}_i)$. Note that boosting algorithms entirely depends on the matrix \mathbf{A} and do not directly interact with the training examples. Our following discussion will thus largely focus on the matrix \mathbf{A} . We write the vector obtained by multiplying a matrix \mathbf{A} with a vector \mathbf{w} as \mathbf{Aw} and its i th entry as $(\mathbf{Aw})_i$. If we let \mathbf{w} represent the coefficients of a selected weak classifier then the margin of the training datum \mathbf{x}_i is $\rho_i = \mathbf{A}_{i:} \mathbf{w} = (\mathbf{Aw})_i$ and the vector of such margins for all of the training data is $\boldsymbol{\rho} = \mathbf{Aw}$.

2 Minimax Probability Machines

Before we introduce our boosting algorithm, let us briefly review the concept of minimax probability machines (MPM) (Lanckriet et al. 2002) first.

2.1 Minimax Probability Classifiers

Let $\mathbf{x}_1 \in \mathbb{R}^n$ and $\mathbf{x}_2 \in \mathbb{R}^n$ denote two random vectors drawn from two distributions with means and covari-

ances $(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)$ and $(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)$, respectively. Here $\boldsymbol{\mu}_1, \boldsymbol{\mu}_2 \in \mathbb{R}^n$ and $\boldsymbol{\Sigma}_1, \boldsymbol{\Sigma}_2 \in \mathbb{R}^{n \times n}$. We define the class labels of \mathbf{x}_1 and \mathbf{x}_2 as +1 and -1, w.l.o.g. The minimax probability machine (MPM) seeks a robust separation hyperplane that can separate the two classes of data with the maximal probability. The hyperplane can be expressed as $\mathbf{w}^\top \mathbf{x} = b$ with $\mathbf{w} \in \mathbb{R}^n \setminus \{\mathbf{0}\}$ and $b \in \mathbb{R}$. The problem of identifying the optimal hyperplane may then be formulated as

$$\max_{\mathbf{w}, b, \gamma} \gamma \quad \text{s.t.} \quad \begin{cases} \inf_{\mathbf{x}_1 \sim (\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)} \Pr\{\mathbf{w}^\top \mathbf{x}_1 \geq b\} \geq \gamma, \\ \inf_{\mathbf{x}_2 \sim (\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)} \Pr\{\mathbf{w}^\top \mathbf{x}_2 \leq b\} \geq \gamma. \end{cases} \quad (1)$$

Here γ is the lower bound of the classification accuracy (or the worst-case accuracy) on test data. This problem can be transformed into a convex problem, more specifically a second-order cone program (SOCP) (Boyd and Vandenberghe 2004) and thus can be solved efficiently (Lanckriet et al. 2002).

2.2 Biased Minimax Probability Machines

The formulation (1) assumes that the classification problem is balanced. It attempts to achieve a high recognition accuracy, which assumes that the losses associated with all mis-classifications are identical. However, in many applications this is not the case.

Huang et al. (2004) proposed a biased version of MPM through a slight modification of (1), which may be formulated as

$$\max_{\mathbf{w}, b, \gamma} \gamma \quad \text{s.t.} \quad \begin{cases} \inf_{\mathbf{x}_1 \sim (\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)} \Pr\{\mathbf{w}^\top \mathbf{x}_1 \geq b\} \geq \gamma, \\ \inf_{\mathbf{x}_2 \sim (\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)} \Pr\{\mathbf{w}^\top \mathbf{x}_2 \leq b\} \geq \gamma_o. \end{cases} \quad (2)$$

Here $\gamma_o \in (0, 1)$ is a prescribed constant, which is the acceptable classification accuracy for the less important class. The resulting decision hyperplane prioritizes the classification of the important class \mathbf{x}_1 over that of the less important class \mathbf{x}_2 . Biased MPM is thus expected to perform better in biased classification applications.

Huang et al. showed that (2) can be iteratively solved via solving a sequence of SOCPs using the fractional programming (FP) technique. Clearly it is significantly more computationally demanding to solve (2) than (1).

Next we show how to re-formulate (2) into a simpler quadratic program (QP) based on the recent theoretical results in (Yu et al. 2009).

2.3 Simplified Biased Minimax Probability Machines

In this section, we are interested in simplifying the problem of (2) for a special case of $\gamma_o = 0.5$, due to its important application in object detection (Viola and Jones 2004; Wu et al. 2005). In the following discussion, for simplicity, we only consider $\gamma_o = 0.5$ although some algorithms developed may also apply to $\gamma_o < 0.5$.

Theoretical results in (Yu et al. 2009) show that, the worst-case constraint in (2) can be written in different forms when \mathbf{x} follows arbitrary, symmetric, symmetric unimodal or Gaussian distributions (see Appendix A). Both the MPM (Lanckriet et al. 2002) and the biased MPM (Huang et al. 2004) are based the most general form of the four cases shown in Appendix A, i.e., Equation (27) for arbitrary distributions, as they do not impose constraints upon the distributions of \mathbf{x}_1 and \mathbf{x}_2 .

However, one may take advantage of structural information whenever available. For example, it is shown in (Wu et al. 2005) that, for the face detection problem, weak classifier outputs can be well approximated by the Gaussian distribution. In other words, the constraint for arbitrary distributions does not utilize any type of *a priori* information, and hence, for many problems, considering arbitrary distributions for simplifying (1) and (2) is too *conservative*. Since both the MPM (Lanckriet et al. 2002) and the biased MPM (Huang et al. 2004) do not assume any constraints on the distribution family, they fail to exploit this structural information.

Let us consider the special case of $\gamma_o = 0.5$. It is easy to see that the worst-case constraint in (2) becomes a simple linear constraint for symmetric, symmetric unimodal, as well as Gaussian distributions (see Appendix A). As pointed out in (Yu et al. 2009), such a result is the immediate consequence of symmetry because the worst-case distributions are forced to put probability mass arbitrarily far away on both sides of the mean. In such case, any information about the covariance is neglected.

We now apply this result to the biased MPM as represented by (2). Our main result is the following theorem.

Theorem 1 *With $\gamma_o = 0.5$, the biased minimax problem (2) can be formulated as an unconstrained problem:*

$$\max_{\mathbf{w}} \frac{\mathbf{w}^\top (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)}{\sqrt{\mathbf{w}^\top \boldsymbol{\Sigma}_1 \mathbf{w}}}, \quad (3)$$

under the assumption that \mathbf{x}_2 follows a symmetric distribution. The optimal b can be obtained through:

$$b = \mathbf{w}^\top \boldsymbol{\mu}_2. \quad (4)$$

The worst-case classification accuracy for the first class, γ^* , is obtained by solving

$$\varphi(\gamma^*) = \frac{-b^* + a^{*\top} \mu_1}{\sqrt{\mathbf{w}^{*\top} \Sigma_1 \mathbf{w}^*}}, \quad (5)$$

where

$$\varphi(\gamma) = \begin{cases} \sqrt{\frac{\gamma}{1-\gamma}} & \text{if } \mathbf{x}_1 \sim (\mu_1, \Sigma_1), \\ \sqrt{\frac{1}{2(1-\gamma)}} & \text{if } \mathbf{x}_1 \sim (\mu_1, \Sigma_1)_S, \\ \frac{2}{3} \sqrt{\frac{1}{2(1-\gamma)}} & \text{if } \mathbf{x}_1 \sim (\mu_1, \Sigma_1)_{SU}, \\ \phi^{-1}(\gamma) & \text{if } \mathbf{x}_1 \sim \mathcal{G}(\mu_1, \Sigma_1). \end{cases} \quad (6)$$

and $\{\mathbf{w}^*, b^*\}$ is the optimal solution of (3) and (4).

Please refer to Appendix A for the proof of Theorem 1.

We have derived the biased MPM algorithm from a different perspective. We reveal that only the assumption of symmetric distributions is needed to arrive at a simple unconstrained formulation. Compared with the approach in (Huang et al. 2004), we have used more information to simplify the optimization problem. More importantly, as will be shown in the next section, this unconstrained formulation enables us to design a new boosting algorithm.

There is a close connection between our algorithm and the linear asymmetric classifier (LAC) in (Wu et al. 2005). The resulting problem (3) is exactly the same as LAC in (Wu et al. 2005). Removing the inequality in this constraint leads to a problem solvable by eigen-decomposition. We have thus shown that the results of Wu et al. may be generalized from the Gaussian distributions assumed in (Wu et al. 2005) to symmetric distributions.

3 Linear Asymmetric Classification

We have shown that starting from the biased minimax probability machine, we are able to obtain the same optimization formulation as shown in Wu et al. (2005), while much weakening the underlying assumption (symmetric distributions versus Gaussian distributions). Before we propose our LACBoost and FisherBoost, however, we provide a brief overview of LAC.

Wu et al. (2008) proposed linear asymmetric classification (LAC) as a post-processing step for training nodes in the cascade framework. In (Wu et al. 2008), it is stated that LAC is guaranteed to reach an optimal solution under the assumption of Gaussian data distributions. We now know that this Gaussianity condition may be relaxed.

Suppose that we have a linear classifier

$$f(\mathbf{x}) = \text{sign}(\mathbf{w}^\top \mathbf{x} - b).$$

We seek a $\{\mathbf{w}, b\}$ pair with a very high accuracy on the positive data \mathbf{x}_1 and a moderate accuracy on the negative \mathbf{x}_2 . This can be expressed as the following problem:

$$\begin{aligned} \max_{\mathbf{w} \neq \mathbf{0}, b} \quad & \Pr_{\mathbf{x}_1 \sim (\mu_1, \Sigma_1)} \{\mathbf{w}^\top \mathbf{x}_1 \geq b\}, \\ \text{s.t.} \quad & \Pr_{\mathbf{x}_2 \sim (\mu_2, \Sigma_2)} \{\mathbf{w}^\top \mathbf{x}_2 \leq b\} = \lambda. \end{aligned} \quad (7)$$

In (Wu et al. 2005), λ is set to 0.5 and it is assumed that for any \mathbf{w} , $\mathbf{w}^\top \mathbf{x}_1$ is Gaussian and $\mathbf{w}^\top \mathbf{x}_2$ is symmetric, (7) can be approximated by (3). Again, these assumptions may be relaxed as we have shown in the last section. Problem (3) is similar to LDA's optimization problem

$$\max_{\mathbf{w} \neq \mathbf{0}} \frac{\mathbf{w}^\top (\mu_1 - \mu_2)}{\sqrt{\mathbf{w}^\top (\Sigma_1 + \Sigma_2) \mathbf{w}}}. \quad (8)$$

Problem (3) can be solved by eigen-decomposition and a closed-form solution can be derived:

$$\mathbf{w}^* = \Sigma_1^{-1}(\mu_1 - \mu_2), \quad b^* = \mathbf{w}^{*\top} \mu_2. \quad (9)$$

On the other hand, each node in cascaded boosting classifiers has the following form:

$$f(\mathbf{x}) = \text{sign}(\mathbf{w}^\top \mathbf{H}(\mathbf{x}) - b). \quad (10)$$

We override the symbol $\mathbf{H}(\mathbf{x})$ here, which denotes the output vector of all weak classifiers over the datum \mathbf{x} . We can cast each node as a linear classifier over the feature space constructed by the binary outputs of all weak classifiers. For each node in a cascade classifier, we wish to maximize the detection rate while maintaining the false positive rate at a moderate level (for example, around 50.0%). That is to say, the problem (3) represents the node learning goal. Boosting algorithms such as AdaBoost can be used as feature selection methods, and LAC is used to learn a linear classifier over those binary features chosen by boosting as in Wu et al. (2005). The advantage of this approach is that LAC considers the asymmetric node learning explicitly.

However, there is a precondition on the validity of LAC that for any \mathbf{w} , $\mathbf{w}^\top \mathbf{x}_1$ is a Gaussian and $\mathbf{w}^\top \mathbf{x}_2$ is symmetric. In the case of boosting classifiers, $\mathbf{w}^\top \mathbf{x}_1$ and $\mathbf{w}^\top \mathbf{x}_2$ can be expressed as the margin of positive data and negative data, respectively. Empirically Wu et al. (2008) verified that $\mathbf{w}^\top \mathbf{x}$ is approximately Gaussian for a cascade face detector. We discuss this issue in more detail in Section 5. Shen and Li (2010b) theoretically proved that under the assumption that weak classifiers are independent, the margin of AdaBoost follows the Gaussian distribution, as long as the number of weak classifiers is *sufficiently large*. In Section 5 we verify this theoretical result by performing the normality test on nodes with different number of weak classifiers.

4 Constructing Boosting Algorithms from LDA and LAC

In kernel methods, the original data are nonlinearly mapped to a feature space by a mapping function $\Psi(\cdot)$. The function need not be known, however, as rather than being applied to the data directly, it acts instead through the inner product $\Psi(\mathbf{x}_i)^\top \Psi(\mathbf{x}_j)$. In boosting (Rätsch et al. 2002), however, the mapping function can be seen as being explicitly known, as $\Psi(\mathbf{x}) : \mathbf{x} \mapsto [h_1(\mathbf{x}), \dots, h_n(\mathbf{x})]$. Let us consider the Fisher LDA case first because the solution to LDA will generalize to LAC straightforwardly, by looking at the similarity between (3) and (8).

Fisher LDA maximizes the between-class variance and minimizes the within-class variance. In the binary-class case, the more general formulation in (8) can be expressed as

$$\max_{\mathbf{w}} \frac{(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^2}{\sigma_1 + \sigma_2} = \frac{\mathbf{w}^\top \mathbf{C}_b \mathbf{w}}{\mathbf{w}^\top \mathbf{C}_w \mathbf{w}}, \quad (11)$$

where \mathbf{C}_b and \mathbf{C}_w are the between-class and within-class scatter matrices; $\boldsymbol{\mu}_1$ and $\boldsymbol{\mu}_2$ are the projected centers of the two classes. The above problem can be equivalently reformulated as

$$\min_{\mathbf{w}} \mathbf{w}^\top \mathbf{C}_w \mathbf{w} - \theta(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2), \quad (12)$$

for some certain constant θ and under the assumption that $\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2 \geq 0$.¹ Now in the feature space, our data are $\Psi(\mathbf{x}_i)$, $i = 1 \dots m$. Define the vectors $\mathbf{e}, \mathbf{e}_1, \mathbf{e}_2 \in \mathbb{R}^m$ such that $\mathbf{e} = \mathbf{e}_1 + \mathbf{e}_2$, the i -th entry of \mathbf{e}_1 is $1/m_1$ if $y_i = +1$ and 0 otherwise, and the i -th entry of \mathbf{e}_2 is $1/m_2$ if $y_i = -1$ and 0 otherwise. We then see that

$$\begin{aligned} \boldsymbol{\mu}_1 &= \frac{1}{m_1} \mathbf{w}^\top \sum_{y_i=1} \Psi(\mathbf{x}_i) = \frac{1}{m_1} \sum_{y_i=1} \mathbf{A}_{i:} \mathbf{w} \\ &= \frac{1}{m_1} \sum_{y_i=1} (\mathbf{A} \mathbf{w})_i = \mathbf{e}_1^\top \mathbf{A} \mathbf{w}, \end{aligned} \quad (13)$$

and

$$\boldsymbol{\mu}_2 = \frac{1}{m_2} \mathbf{w}^\top \sum_{y_i=-1} \Psi(\mathbf{x}_i) = \frac{1}{m_2} \sum_{y_i=-1} \mathbf{H}_{i:} \mathbf{w} = -\mathbf{e}_2^\top \mathbf{A} \mathbf{w}, \quad (14)$$

For ease of exposition we order the training data according to their labels so the vector $\mathbf{e} \in \mathbb{R}^m$:

$$\mathbf{e} = [1/m_1, \dots, 1/m_2, \dots]^\top, \quad (15)$$

and the first m_1 components of $\boldsymbol{\rho}$ correspond to the positive training data and the remaining ones correspond

¹ In our object detection experiment, we found that this assumption can always be satisfied.

to the m_2 negative data. We now see that $\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2 = \mathbf{e}^\top \boldsymbol{\rho}$, $\mathbf{C}_w = m_1/m \cdot \boldsymbol{\Sigma}_1 + m_2/m \cdot \boldsymbol{\Sigma}_2$ with $\boldsymbol{\Sigma}_{1,2}$ the covariance matrices. Noting that

$$\mathbf{w}^\top \boldsymbol{\Sigma}_{1,2} \mathbf{w} = \frac{1}{m_{1,2}(m_{1,2} - 1)} \sum_{i>k, y_i=y_k=\pm 1} (\rho_i - \rho_k)^2,$$

we can easily rewrite the original problem (11) (and (12)) into:

$$\begin{aligned} \min_{\mathbf{w}, \boldsymbol{\rho}} \quad & \frac{1}{2} \boldsymbol{\rho}^\top \mathbf{Q} \boldsymbol{\rho} - \theta \mathbf{e}^\top \boldsymbol{\rho}, \\ \text{s.t.} \quad & \mathbf{w} \succcurlyeq \mathbf{0}, \mathbf{1}^\top \mathbf{w} = 1, \\ & \rho_i = (\mathbf{A} \mathbf{w})_i, i = 1, \dots, m. \end{aligned} \quad (16)$$

Here $\mathbf{Q} = \begin{bmatrix} \mathbf{Q}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{Q}_2 \end{bmatrix}$ is a block matrix with

$$\mathbf{Q}_1 = \begin{bmatrix} \frac{1}{m} & -\frac{1}{m(m_1-1)} & \cdots & -\frac{1}{m(m_1-1)} \\ -\frac{1}{m(m_1-1)} & \frac{1}{m} & \cdots & -\frac{1}{m(m_1-1)} \\ \vdots & \vdots & \ddots & \vdots \\ -\frac{1}{m(m_1-1)} & -\frac{1}{m(m_1-1)} & \cdots & \frac{1}{m} \end{bmatrix},$$

and \mathbf{Q}_2 is similarly defined by replacing m_1 with m_2 in \mathbf{Q}_1 :

$$\mathbf{Q}_2 = \begin{bmatrix} \frac{1}{m} & -\frac{1}{m(m_2-1)} & \cdots & -\frac{1}{m(m_2-1)} \\ -\frac{1}{m(m_2-1)} & \frac{1}{m} & \cdots & -\frac{1}{m(m_2-1)} \\ \vdots & \vdots & \ddots & \vdots \\ -\frac{1}{m(m_2-1)} & -\frac{1}{m(m_2-1)} & \cdots & \frac{1}{m} \end{bmatrix}.$$

Also note that we have introduced a constant $\frac{1}{2}$ before the quadratic term for convenience. The normalization constraint $\mathbf{1}^\top \mathbf{w} = 1$ removes the scale ambiguity of \mathbf{w} . Without it the problem is ill-posed.

We see from the form of (3) that the covariance of the negative data is not involved in LAC and thus that if we set $\mathbf{Q} = \begin{bmatrix} \mathbf{Q}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}$ then (16) becomes the optimization problem of LAC.

At this stage, it remains unclear about how to solve the problem (16) because we do not know all the weak classifiers. There may be extremely (or even infinitely) many weak classifiers in \mathcal{H} , the set from which $h(\cdot)$ is selected, meaning that the dimension of the optimization variable \mathbf{w} may also be extremely large. So (16) is a semi-infinite quadratic program (SIQP). We show how column generation can be used to solve this problem. To make column generation applicable, we need to derive a specific Lagrange dual of the primal problem.

4.1 The Lagrange Dual Problem

We now derive the Lagrange dual of the quadratic problem (16). Although we are only interested in the variable \mathbf{w} , we need to keep the auxiliary variable $\boldsymbol{\rho}$ in order to obtain a meaningful dual problem. The Lagrangian of (16) is

$$L(\underbrace{\mathbf{w}}_{\text{primal}}, \underbrace{\boldsymbol{\rho}, \mathbf{u}, r}_{\text{dual}}) = \frac{1}{2}\boldsymbol{\rho}^\top \mathbf{Q}\boldsymbol{\rho} - \theta \mathbf{e}^\top \boldsymbol{\rho} + \mathbf{u}^\top (\boldsymbol{\rho} - \mathbf{A}\mathbf{w}) - \mathbf{q}^\top \mathbf{w} + r(\mathbf{1}^\top \mathbf{w} - 1), \quad (17)$$

with $\mathbf{q} \succcurlyeq \mathbf{0}$. $\sup_{\mathbf{u}, r} \inf_{\mathbf{w}, \boldsymbol{\rho}} L(\mathbf{w}, \boldsymbol{\rho}, \mathbf{u}, r)$ gives the following Lagrange dual:

$$\begin{aligned} \max_{\mathbf{u}, r} \quad & -r - \overbrace{\frac{1}{2}(\mathbf{u} - \theta \mathbf{e})^\top \mathbf{Q}^{-1}(\mathbf{u} - \theta \mathbf{e})}^{\text{regularization}}, \\ \text{s.t.} \quad & \sum_{i=1}^m u_i \mathbf{A}_i \preceq r \mathbf{1}^\top. \end{aligned} \quad (18)$$

In our case, \mathbf{Q} is rank-deficient and its inverse does not exist (for both LDA and LAC). Actually for both \mathbf{Q}_1 and \mathbf{Q}_2 , they have a zero eigenvalue with the corresponding eigenvector being all ones. This is easy to see because for \mathbf{Q}_1 and \mathbf{Q}_2 , the sum of each row (or each column) is zero. We can simply regularize \mathbf{Q} with $\mathbf{Q} + \tilde{\delta} \mathbf{I}$ with $\tilde{\delta}$ a small positive constant. Actually, \mathbf{Q} is a diagonally dominant matrix but not strict diagonal dominance. So $\mathbf{Q} + \tilde{\delta} \mathbf{I}$ with any $\tilde{\delta} > 0$ is strict diagonal dominance and by the Gershgorin circle theorem, a strictly diagonally dominant matrix must be invertible.

One of the KKT optimality conditions between the dual and primal is

$$\boldsymbol{\rho}^* = -\mathbf{Q}^{-1}(\mathbf{u}^* - \theta \mathbf{e}), \quad (19)$$

which can be used to establish the connection between the dual optimum and the primal optimum. This is obtained by the fact that the gradient of L w.r.t. $\boldsymbol{\rho}$ must vanish at the optimum, $\partial L / \partial \rho_i = 0$, $\forall i = 1 \dots n$.

Problem (18) can be viewed as a regularized LP-Boost problem. Compared with the hard-margin LP-Boost (Demiriz et al. 2002), the only difference is the regularization term in the cost function. The duality gap between the primal (16) and the dual (18) is zero. In other words, the solutions of (16) and (18) coincide. Instead of solving (16) directly, one calculates the most violated constraint in (18) iteratively for the current solution and adds this constraint to the optimization problem. In theory, any column that violates dual feasibility can be added. To speed up the convergence, we

add the most violated constraint by solving the following problem:

$$h'(\cdot) = \operatorname{argmax}_{h(\cdot)} \sum_{i=1}^m u_i y_i h(\mathbf{x}_i). \quad (20)$$

This is exactly the same as the one that standard Adaboost and LPBoost use for producing the best weak classifier at each iteration. That is to say, to find the weak classifier that has the minimum weighted training error. We summarize the LACBoost/FisherBoost algorithm in Algorithm 1. By simply changing \mathbf{Q}_2 , Algorithm 1 can be used to train either LACBoost or FisherBoost. Note that to obtain an actual strong classifier, one may need to include an offset b , i.e. the final classifier is $\sum_{j=1}^n h_j(\mathbf{x}) - b$ because from the cost function of our algorithm (12), we can see that the cost function itself does not minimize any classification error. It only finds a projection direction in which the data can be maximally separated. A simple line search can find an optimal b . Moreover, when training a cascade, we need to tune this offset anyway as shown in (10).

The convergence of Algorithm 1 is guaranteed by general column generation or cutting-plane algorithms, which is easy to establish:

Theorem 2 *The column generation procedure decreases the objective value of problem (16) at each iteration and hence in the limit it solves the problem (16) globally to a desired accuracy.*

The proof is deferred to Appendix B. In short, when a new $h'(\cdot)$ that violates dual feasibility is added, the new optimal value of the dual problem (maximization) would decrease. Accordingly, the optimal value of its primal problem decreases too because they have the same optimal value due to zero duality gap. Moreover the primal cost function is convex, therefore in the end it converges to the global minimum.

At each iteration of column generation, in theory, we can solve either the dual (18) or the primal problem (16). Here we choose to solve an equivalent variant of the primal problem (16):

$$\min_{\mathbf{w}} \frac{1}{2} \mathbf{w}^\top (\mathbf{A}^\top \mathbf{Q} \mathbf{A}) \mathbf{w} - (\theta \mathbf{e}^\top \mathbf{A}) \mathbf{w}, \quad \text{s.t. } \mathbf{w} \in \Delta_n, \quad (21)$$

where Δ_n is the unit simplex, which is defined as $\{\mathbf{w} \in \mathbb{R}^n : \mathbf{1}^\top \mathbf{w} = 1, \mathbf{w} \succcurlyeq \mathbf{0}\}$.

In practice, it could be much faster to solve (21) since

1. Generally, the primal problem has a smaller size, hence faster to solve. The number of variables of (18) is m at each iteration, while the number of variables is the number of iterations for the primal

Algorithm 1 Column generation for SIQP.

Input: Labeled training data $(\mathbf{x}_i, y_i), i = 1 \dots m$; termination threshold $\varepsilon > 0$; regularization parameter θ ; maximum number of iterations n_{\max} .

- 1 **Initialization:** $m = 0$; $\mathbf{w} = \mathbf{0}$; and $u_i = \frac{1}{m}, i = 1 \dots m$.
- 2 **for** iteration = 1 : n_{\max} **do**
- 3 – Check for the optimality:
 if iteration > 1 and $\sum_{i=1}^m u_i y_i h'(\mathbf{x}_i) < r + \varepsilon$,
 then
 break; and the problem is solved;
- 4 – Add $h'(\cdot)$ to the restricted master problem, which corresponds to a new constraint in the dual;
- 5 – Solve the dual problem (18) (or the primal problem (16)) and update r and u_i ($i = 1 \dots m$).
- 6 – Increment the number of weak classifiers $n = n + 1$.

Output: The selected features are h_1, h_2, \dots, h_n . The final strong classifier is: $F(\mathbf{x}) = \sum_{j=1}^n w_j h_j(\mathbf{x}) - b$. Here the offset b can be learned by a simple line search.

problem. For example, in Viola-Jones' face detection framework, the number of training data $m = 10,000$ and $n_{\max} = 200$. In other words, the primal problem has at most 200 variables in this case;

2. The dual problem (18) is a standard QP problem. It has no special structure to exploit. As we will show, the primal problem (21) belongs to a special class of problems and can be efficiently solved using entropic/exponentiated gradient descent (EG) (Beck and Teboulle 2003; Collins et al. 2008). See Appendix C for details of the EG algorithm. A fast QP solver is extremely important for training our object detector since we need to solve a few thousand QP problems. Compared with standard QP solvers like Mosek (MOSEK 2010), EG is much faster. EG makes it possible to train a detector using almost the same amount of time as using standard AdaBoost because the majority of time is spent on weak classifier training and bootstrapping.

We can recover both of the dual variables \mathbf{u}^*, r^* easily from the primal variable \mathbf{w}^*, ρ^* :

$$\mathbf{u}^* = -\mathbf{Q}\rho^* + \theta\mathbf{e}; \quad (22)$$

$$r^* = \max_{j=1 \dots n} \left\{ \sum_{i=1}^m u_i^* \mathbf{A}_{ij} \right\}. \quad (23)$$

The second equation is obtained by the fact that in the dual problem's constraints, at optimum, there must exist at least one u_i^* such that the equality holds. That is to say, r^* is the largest *edge* over all weak classifiers.

In summary, when using EG to solve the primal problem, Line 5 of Algorithm 1 is:

– Solve the primal problem (21) using EG, and update the dual variables \mathbf{u} with (22), and r with (23).

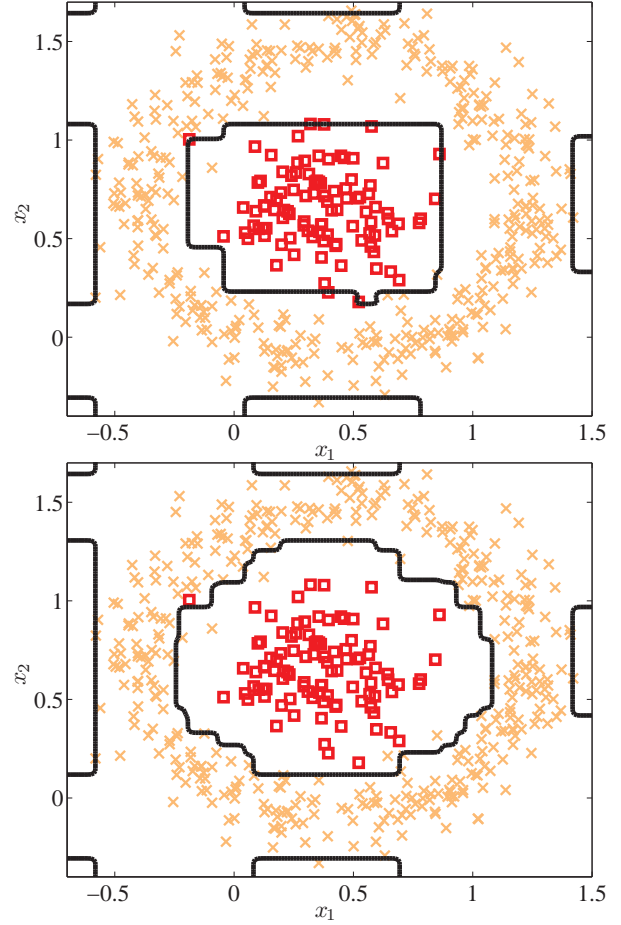


Fig. 2: Decision boundaries of AdaBoost (top) and FisherBoost (bottom) on 2D artificial data generated from the Gaussian distribution (positive data represented by \square 's and negative data by \times 's). Weak classifiers are vertical and horizontal decision stumps. FisherBoost emphasizes more on positive samples than negative samples. As a result, the decision boundary of FisherBoost is more similar to the Gaussian distribution than the decision boundary of AdaBoost.

5 Experiments

In this section, we perform our experiments on both synthetic and challenging real-world data sets, *e.g.*, face and pedestrian detection.

5.1 Synthetic Testing

We first illustrate the performance of FisherBoost on an asymmetrical synthetic data set where there are a large number of negative samples compared to the positive ones. Fig. 2 demonstrates the subtle difference in classification boundaries between AdaBoost and FisherBoost. It can be observed that FisherBoost places more

emphasis on positive samples than negative samples to ensure these positive samples would be classified correctly. AdaBoost, on the other hand, treat both positive and negative samples equally. This might be due to the fact that AdaBoost only optimizes the overall classification accuracy. This finding is consistent with our results reported earlier in (Paisitkriangkrai et al. 2009; Shen et al. 2011).

5.2 Comparison With Other Asymmetric Boosting

In this experiment, FisherBoost and LACBoost are compared against several asymmetric boosting algorithms, namely, AdaBoost with LAC or Fisher LDA post-processing (Wu et al. 2008), AsymBoost (Viola and Jones 2002), cost-sensitive AdaBoost (CS-ADA) (Masnadi-Shirazi and Vasconcelos 2011) and rate constrained boosting (RCBoost) (Saberian and Vasconcelos 2012). The results of AdaBoost are also presented as the baseline. For each algorithm, we train a strong classifier consisting of 100 weak classifiers along with their coefficients. The threshold was determined such that the false positive rate of test set is 50%. For every method, the experiment is repeated 5 times and the average detection rate on positive class is reported. For FisherBoost and LACBoost, the parameter θ is chosen from $\{1/10, 1/12, 1/15, 1/20\}$ by cross-validation. For AsymBoost, we choose k (asymmetric factor) from $\{2^{0.1}, 2^{0.2}, \dots, 2^{0.5}\}$ by cross-validation. For CS-ADA, we set the cost for misclassifying positive and negative data as follows. We assign the asymmetric factor $k = C_1/C_2$ and restrict $0.5(C_1 + C_2) = 1$. We choose k from $\{1.2, 1.65, 2.1, 2.55, 3\}$ by cross-validation. For RCBoost, we conduct two experiments. In the first experiment, we use the same training set to enforce the target detection rate, while in the second experiment; we use 75% of the training data to train the model and the other 25% to enforce the target detection rate. We set the target detection rate, D_T , to 99.5%, the barrier coefficient, γ , to 2 and the number of iterations before halving γ , N_d , to 10.

We tested the performance of all algorithms on five real-world data sets, including both machine learning (USPS) and vision data sets (cars, faces, pedestrians, scenes). We categorized USPS data sets into two classes: even digits and odd digits. For faces, we use face data sets from (Viola and Jones 2004) and randomly extract 5000 negative patches from background images. We apply principle component analysis (PCA) to preserve 95% total variation. The new data set has a dimension of 93. For UIUC car (Agarwal et al. 2004), we downsize the original image from 40×100 pixels to 20×50 pixels and apply PCA. The projected data capture

95% total variation and has a final dimension of 228. For Daimler-Chrysler pedestrian data sets (Munder and Gavrilu 2006), we apply PCA to the original 18×36 pixels. The projected data capture 95% variation and has a final dimension of 139. For indoor/outdoor scene, we divide the 15-scene data set used in (Lazebnik et al. 2006) into 2 groups: indoor and outdoor scenes. We use CENTRIST as our feature descriptors and build 50 visual code words using the histogram intersection kernel (Wu and Rehg 2011). Each image is represented in a spatial hierarchy manner. Each image consists of 31 sub-windows. In total, there are 1550 feature dimensions per image. All 5 classifiers are trained to remove 50% of the negative data, while retaining almost all positive data. We compare their detection rates in Table 1. From our experiments, FisherBoost demonstrates the best performance on most data sets. However, LACBoost does not perform as well as expected. We suspect that the poor performance might partially due to numerical issues, which can cause overfitting. We will discuss this in more detail in Section 5.6.

5.3 Face Detection Using a Cascade Classifier

In this experiments, eight asymmetric boosting methods are evaluated with the multi-exit cascade (Pham et al. 2008), which are FisherBoost/LACBoost, AdaBoost alone or with LDA/LAC post-processing (Wu et al. 2008), AsymBoost alone or with LDA/LAC post-processing. We have also implemented Viola-Jones' face detector (AdaBoost with the conventional cascade) as the baseline (Viola and Jones 2004). Furthermore, our face detector is also compared with state-of-the-art including some cascade design methods, *i.e.*, WaldBoost (Sochman and Matas 2005), FloatBoost (Li and Zhang 2004), Boosting Chain (Xiao et al. 2003) and the extension of (Saberian and Vasconcelos 2010), RCECBoost (Saberian and Vasconcelos 2012). The algorithm for training a multi-exit cascade is summarized in Algorithm 2.

We first illustrate the validity of adopting LAC and Fisher LDA post-processing to improve the node learning objective in the cascade classifier. As described above, LAC and LDA assume that the margin of the training data associated with the node classifier in such a cascade exhibits a Gaussian distribution. We demonstrate this assumption on the face detection task in Fig. 3. Fig. 3 shows the normal probability plot of the margins of the positive training data for the first three node classifiers in the multi-exit LAC classifier. The figure reveals that the larger the number of weak classifiers used the more closely the margins follow the Gaussian

	AdaBoost	LAC	FLDA	AsymBoost	CS-ADA	RCBoost ¹	RCBoost ²	LACBoost	FisherBoost
Digits	99.30 (0.10)	99.30 (0.21)	99.37 (0.08)	99.40 (0.11)	99.37 (0.09)	99.36 (0.17)	99.27 (0.15)	99.12 (0.07)	99.40 (0.13)
Faces	98.70 (0.14)	98.78 (0.42)	98.86 (0.22)	98.73 (0.14)	98.71 (0.20)	98.75 (0.18)	98.66 (0.23)	98.63 (0.29)	98.89 (0.15)
Cars	97.02 (1.55)	97.07 (1.34)	97.02 (1.50)	97.11 (1.36)	97.47 (1.31)	96.84 (0.87)	96.62 (1.08)	96.80 (1.47)	97.78 (1.27)
Pedestrians	98.54 (0.34)	98.59 (0.71)	98.69 (0.28)	98.55 (0.45)	98.51 (0.36)	98.67 (0.29)	98.65 (0.39)	99.12 (0.35)	98.73 (0.33)
Scenes	99.59 (0.10)	99.54 (0.21)	99.57 (0.12)	99.66 (0.12)	99.68 (0.10)	99.61 (0.19)	99.62 (0.16)	97.50 (1.07)	99.66 (0.10)
Average	98.63	98.66	98.70	98.69	98.75	98.64	98.56	98.23	98.89

Table 1: Test errors (%) on five real-world data sets. All experiments are run 5 times with 100 boosting iterations. The average detection rate and standard deviation (in percentage) at 50% false positives are reported. Best average detection rate is shown in boldface.

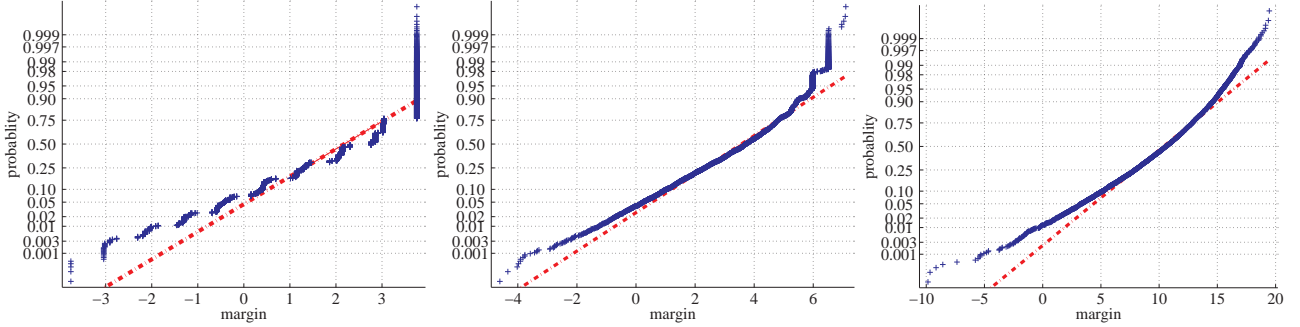


Fig. 3: Normality test (normal probability plot) for the face data’s margin distribution of nodes 1, 2, 3. The 3 nodes contains 7, 22, 52 weak classifiers respectively. The data are plotted against a theoretical normal distribution such that the data which follows the normal distribution model should form a straight line. Curves deviated from the straight line (the red line) indicate departures from normality. The larger the number of weak classifiers, the more closely the margin follow the Gaussian distribution.

distribution. From this, we infer that LAC/LDA post-processing and thus LACBoost and FisherBoost, can be expected to achieve a better performance when a larger number of weak classifiers are used. We therefore apply LAC/LDA only within the later nodes (for example, 9 onwards) of a multi-exit cascade as these nodes contain more weak classifiers. We choose multi-exit due to its property² and effectiveness as reported in (Pham et al. 2008). We have compared the multi-exit cascade with LDA/LAC post-processing against the conventional cascade with LDA/LAC post-processing in (Wu et al. 2008) and performance improvement has been observed.

As in (Wu et al. 2008), five basic types of Haar-like features are calculated, resulting in a 162,336 dimensional over-complete feature set on an image of 24×24 pixels. To speed up the weak classifier training, as in (Wu et al. 2008), we uniformly sample 10% of features for training weak classifiers (decision stumps). The face data set consists of 9,832 mirrored 24×24 images (Viola and Jones 2004) (5,000 images used for training and 4,832 imaged used for validation) and 7,323 larger resolution background images, as used in (Wu et al. 2008).

² Since the multi-exit cascade makes use of all previous weak classifiers in earlier nodes, it would meet the Gaussianity requirement better than the conventional cascade classifier.

Several multi-exit cascades are trained with various algorithms described above. In order to ensure a fair comparison, we have used the same number of multi-exit stages and the same number of weak classifiers. Each multi-exit cascade consists of 22 exits and 2,923 weak classifiers. The indices of exit nodes are pre-determined to simplify the training procedure.

For our FisherBoost and LACBoost, we have an important parameter θ , which is chosen from $\{\frac{1}{10}, \frac{1}{12}, \frac{1}{15}, \frac{1}{20}, \frac{1}{25}, \frac{1}{30}, \frac{1}{40}, \frac{1}{50}\}$. We have not carefully tuned this parameter using cross-validation. Instead, we train a 10-node cascade for each candidate θ , and choose the one with the best *training* accuracy.³ At each exit, negative examples misclassified by current cascade are discarded, and new negative examples are bootstrapped from the background images pool. In total, billions of negative examples are extracted from the pool. The positive training data and validation data keep unchanged during the training process.

Our experiments are performed on a workstation with 8 Intel Xeon E5520 CPUs and 32GB RAM. It takes about 3 hours to train the multi-exit cascade with AdaBoost or AsymBoost. For FisherBoost and LACBoost, it takes less than 4 hours to train a complete

³ To train a complete 22-node cascade and choose the best θ on cross-validation data may give better detection rates.

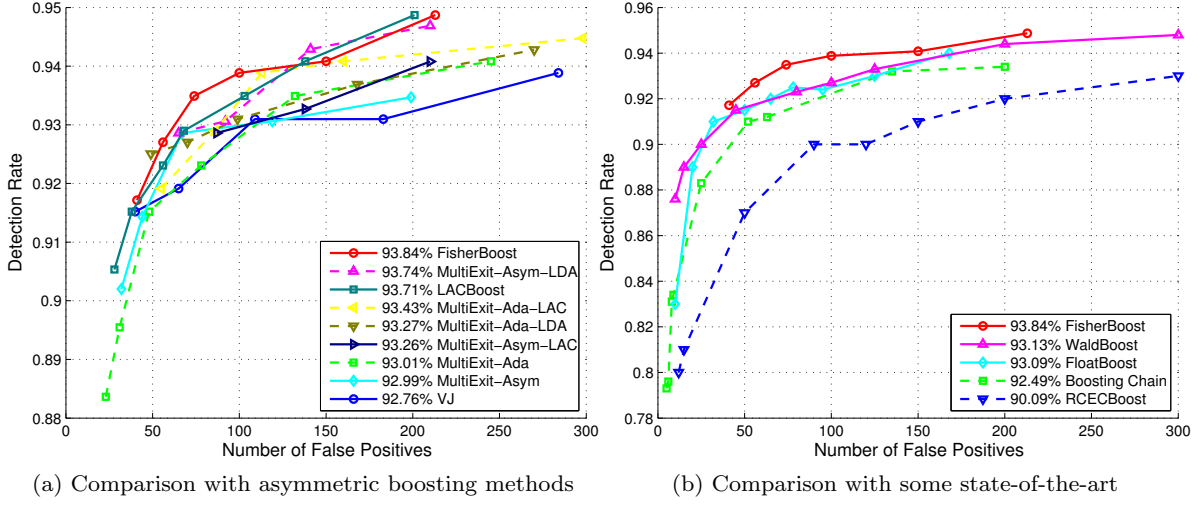


Fig. 4: Our face detectors are compared with other asymmetric boosting methods (a) and some state-of-the-art including cascade design methods (b) on MIT+CMU frontal face test data using ROC curves (number of false positives versus detection rate). “Ada” and “Asym” mean that features are selected using AdaBoost and AsymBoost, respectively. “VJ” implements Viola and Jones’ cascade using AdaBoost (Viola and Jones 2004). “MultiExit” means the multi-exit cascade (Pham et al. 2008). The ROC curves of compared methods in (b) are quoted from their original papers (Sochman and Matas 2005; Li and Zhang 2004; Xiao et al. 2003; Saberian and Vasconcelos 2012). Compared methods are ranked in the legend, based on the average of detection rates.

multi-exit cascade.⁴ In other words, our EG algorithm takes less than 1 hour to solve the primal QP problem (we need to solve a QP at each iteration). As an estimation of the computational complexity, suppose that the number of training examples is m , number of weak classifiers is n . At each iteration of the cascade training, the complexity of solving the primal QP using EG is $O(mn + kn^2)$ with k the iterations needed for EG’s convergence. The complexity for training the weak classifier is $O(md)$ with d the number of all Haar-feature patterns. In our experiment, $m = 10,000$, $n \approx 2900$, $d = 160,000$, $k < 500$. So the majority of the computational cost of the training process is bound up in the weak classifier training.

We have also experimentally observed the speedup of EG against standard QP solvers. We solve the primal QP defined by (21) using EG and Mosek (MOSEK 2010). The QP’s size is 1,000 variables. With the same accuracy tolerance (Mosek’s primal-dual gap is set to 10^{-7} and EG’s convergence tolerance is also set to 10^{-7}), Mosek takes 1.22 seconds and EG is 0.0541 seconds on a standard desktop. So EG is about 20 times faster. Moreover, at iteration $n+1$ of training the cascade, EG can take advantage of the last iteration’s solution by starting EG from a small perturbation of the previous

solution. Such a warm-start gains a 5 to 10 \times speedup in our experiment, while the current QP solver in Mosek does not support warm-start (MOSEK 2010, Chapter 7).

We evaluate the detection performance on the MIT+CMU frontal face test set. This dataset is made up of 507 frontal faces in 130 images with different background.

If one positive output has less than 50% variation of shift and scale from the ground-truth, we treat it as a true positive, otherwise a false positive.

In the test phase, the scale factor of the scanning window is set to 1.2 and the stride step is set to 1 pixel.

The Receiver operating characteristic (ROC) curves in Fig. 4 show the entire cascade’s performance. The average detection rate (similar with the one used in (Dollár et al. 2012)) are used to rank the compared methods, which is the mean of detection rates sampled evenly from 50 to 200 false positives. Note that multiple factors impact on the cascade’s performance, however, including: the classifier set, the cascade structure, bootstrapping *etc.* Fig. 4 (a) demonstrate the superior performance of FisherBoost to other asymmetric boosting methods in the face detection task. We can also find that LACBoost perform worse than FisherBoost. Wu et al. have observed that LAC post-processing does not outperform LDA post-processing in some cases either.

⁴ Our implementation is in C++ and only the weak classifier training part is parallelized using OpenMP.

Algorithm 2 The procedure for training a multi-exit cascade with LACBoost or FisherBoost.

Input:

- A training set with m examples, which are ordered by their labels (m_1 positive examples followed by m_2 negative examples);
- d_{\min} : minimum acceptable detection rate per node;
- f_{\max} : maximum acceptable false positive rate per node;
- F_{fp} : target overall false positive rate.

```

1 Initialize:
   $t = 0$ ; (node index)
   $n = 0$ ; (total selected weak classifiers up to the current node)
   $D_t = 1$ ;  $F_t = 1$ . (overall detection rate and false positive rate up to the current node)
2 while  $F_{\text{fp}} < F_t$  do
3    $t = t + 1$ ; (increment node index)
4   while  $d_t < d_{\min}$  do
5     (current detection rate  $d_t$  is not acceptable yet)
6     –  $n = n + 1$ , and generate a weak classifier and update all the weak classifiers' linear coefficient using LACBoost or FisherBoost.
7     – Adjust threshold  $b$  of the current boosted strong classifier
8
9     
$$F^t(\mathbf{x}) = \sum_{j=1}^n w_j^t h_j(\mathbf{x}) - b$$

10    such that  $f_t \approx f_{\max}$ .
11    – Update the detection rate of the current node  $d_t$  with the learned boosted classifier.
12  Update  $D_{t+1} = D_t \times d_t$ ;  $F_{t+1} = F_t \times f_t$ 
13  Remove correctly classified negative samples from negative training set.
14  if  $F_{\text{fp}} < F_t$  then
15    Evaluate the current cascaded classifier on the negative images and add misclassified samples into the negative training set; (bootstrap)
16 Output: A multi-exit cascade classifier with  $n$  weak classifiers and  $t$  nodes.

```

We have also compared our methods with the boosted greedy sparse LDA (BGS LDA) in (Paisitkriangkrai et al. 2009; Shen et al. 2011), which is considered one of the state-of-the-art. FisherBoost and LACBoost outperform BGS LDA with AdaBoost/AsymBoost in the detection rate. Note that BGS LDA uses the standard cascade.

From Fig. 4 (b), we can see the performance of FisherBoost is better than the other considered cascade design methods. However, since the parameters of cascade structure (*e.g.*, node thresholds, number of nodes, number of weak classifiers per node) are not carefully tuned, our method can not guarantee an optimal trade-off between accuracy and speed. We believe that the boosting method and the cascade design strategy compensate each other. Actually in (Saberian and Vasconcelos 2010), the authors also incorporate some cost-sensitive boosting algorithms, *e.g.*, cost-sensitive AdaBoost (Masnadi-Shirazi and Vasconcelos 2011), Asym-

Boost (Viola and Jones 2002), with their cascade design method.

5.4 Pedestrian Detection Using a Cascade Classifier

We run our experiments on a pedestrian detection with a minor modification to visual features being used. We evaluate our approach on INRIA data set (Dalal and Triggs 2005). The training set consists of 2,416 cropped mirrored pedestrian images and 1,200 large resolution background images. The test set consists of 288 images containing 588 annotated pedestrians and 453 non-pedestrian images. Each training sample is scaled to 64×128 pixels with an additional of 16 pixels added to each border to preserve human contour information. During testing, the detection scanning window is resized to 32×96 pixels to fit the human body. We use histogram of oriented gradient (HOG) features in our experiments. Instead of using fixed-size blocks (105 blocks of size 16×16 pixels) as in Dalal and Triggs (Dalal and Triggs 2005), we define blocks with various scales (from 12×12 pixels to 64×128 pixels) and width-length ratios ($1 : 1$, $1 : 2$, $2 : 1$, $1 : 3$, and $3 : 1$). Each block is divided into 2×2 cells, and HOG features in each cell are summarized into 9 bins. Hence 36-dimensional HOG feature is generated from each block. In total, there are 7,735 blocks from a 64×128 -pixels patch. ℓ_1 -norm normalization is then applied to the feature vector. Furthermore, we use integral histograms to speed up the computation as in (Zhu et al. 2006). At each iteration, we randomly sample 10% of all the possible blocks for training a weak classifier. We have used weighted linear discriminant analysis (WLDA) as weak classifiers, same as in (Paisitkriangkrai et al. 2008). Zhu et al. used linear support vector machines as weak classifiers (Zhu et al. 2006), which can also be used as weak classifiers here.

In this experiment, all cascade classifiers have the same number of nodes and weak classifiers. For the same reason described in the face detection section, the FisherBoost/LACBoost and Wu et al.'s LDA/LAC post-processing are applied to the cascade from the 3-rd node onwards, instead of the first node. The positive examples remain the same for all nodes while the negative examples in later nodes are obtained by a bootstrap approach. The parameter θ of our FisherBoost and LACBoost is selected from $\{\frac{1}{10}, \frac{1}{12}, \frac{1}{14}, \frac{1}{16}, \frac{1}{18}, \frac{1}{20}\}$. We have not carefully selected θ in this experiment. Ideally, cross-validation should be used to pick the best value of θ by using an independent cross-validation data set. Since there are not many labeled positive training data in the INRIA data set, we use the same 2,416 positive examples for validation. We collect 500 additional

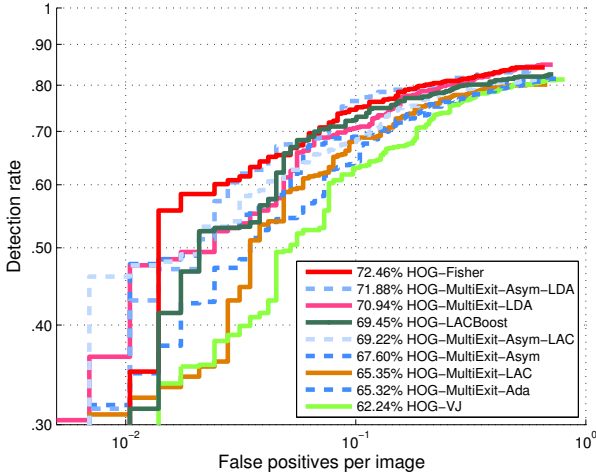


Fig. 5: FisherBoost (HOG-Fisher) and LACBoost (HOG-LACBoost) are compared with other cascade pedestrian detectors on the INRIA data set. All cascades are trained with the same number of weak classifiers and nodes, using HOG features. In the legend, detectors are sorted based on their log-average detection rates. FisherBoost performs best compared to other cascades.

negative examples by bootstrapping for validation. Further improvement is expected if the positive data used during validation is different from those used during training. During evaluation, we use a step stride of 4×4 pixels with 10 scales per octave (a scale ratio of 1.0718). The performance of different cascade detectors is evaluated using a protocol described in (Dollár et al. 2012). A technique known as pairwise maximum suppression (Dollár 2012) is applied to suppress less confident detection windows. A confidence score is needed for each detection window as the input of pairwise maximum suppression. In this work, this confidence is simply calculated as the mean of decision scores of the last five nodes in the cascade.

The ROC curves are plotted in Fig. 5. Same as (Dollár et al. 2012), the log-average detection rate is used to summarize overall detection performance, which is the mean of detection rates sampled evenly at 9 positions from 0.01 to 1. In general, FisherBoost (HOG-Fisher) outperforms all other cascade detectors. Similar to our previous experiments, LAC and LDA post-processing further improve the performance of AdaBoost. However, we observe that both FisherBoost and LDA post-processing have a better generalization performance than LACBoost and LAC post-processing. We will discuss this issue at the end of the experiments.

5.5 Comparison with State-of-the-art Pedestrian Detectors

In this experiment, we compare FisherBoost with state-of-the-art pedestrian detectors on several public data sets. In (Dollár et al. 2012), the authors compare various pedestrian detectors and conclude that combining multiple discriminative features can often significantly boost the performance of pedestrian detection. This is not surprising since a similar conclusion was drawn in (Gehler and Nowozin 2009) on an object recognition task. Clearly, the pedestrian detector, which relies solely on the HOG feature, is unlikely to outperform those using a combination of features.

To this end, we train our pedestrian detector by combining both HOG features (Dalal and Triggs 2005) and covariance features (Tuzel et al. 2008)⁵. For HOG, we use the same experimental settings as our previous experiment. For covariance features, we use the following image statistics $[x, y, I, |I_x|, |I_y|, \sqrt{I_x^2 + I_y^2}, |I_{xx}|, |I_{yy}|, \arctan(|I_x|/|I_y|)]$, where x and y are the pixel location, I is the pixel intensity, I_x and I_y are first order intensity derivatives, I_{xx} and I_{yy} are second order intensity derivatives and the edge orientation. Each pixel is mapped to a 9-dimensional feature image. We then calculate 36 correlation coefficients in each block and concatenate these features to previously computed HOG features. The new feature not only encodes the gradient histogram (edges) but also information of the correlation of defined statistics inside each spatial layout (texture). Similar to the previous experiment, we project these new features to a line using weighted linear discriminant analysis. Except for new features, other training and test implementations are the same with those in the previous pedestrian detection experiments.

We first compare FisherBoost (HOGCOV-Fisher) with two baseline detectors trained with AdaBoost. The first baseline detector is trained with the conventional cascade (HOGCOV-VJ) while the second baseline detector is trained with the multi-exit cascade (HOGCOV-MultiExit-Ada). All detectors are trained with both HOG and covariance features on INRIA training set. The results on INRIA test sets using the protocol in (Dollár et al. 2012) are reported in Fig. 6 (a). Similar to previous results, FisherBoost outperforms both baseline detectors.

⁵ Covariance features capture the relationship between different image statistics and have been shown to perform well in our previous experiments. However, other discriminative features can also be used here instead, *e.g.*, Haar-like features, Local Binary Pattern (LBP) (Mu et al. 2008) and self-similarity of low-level features (CSS) (Walk et al. 2010).

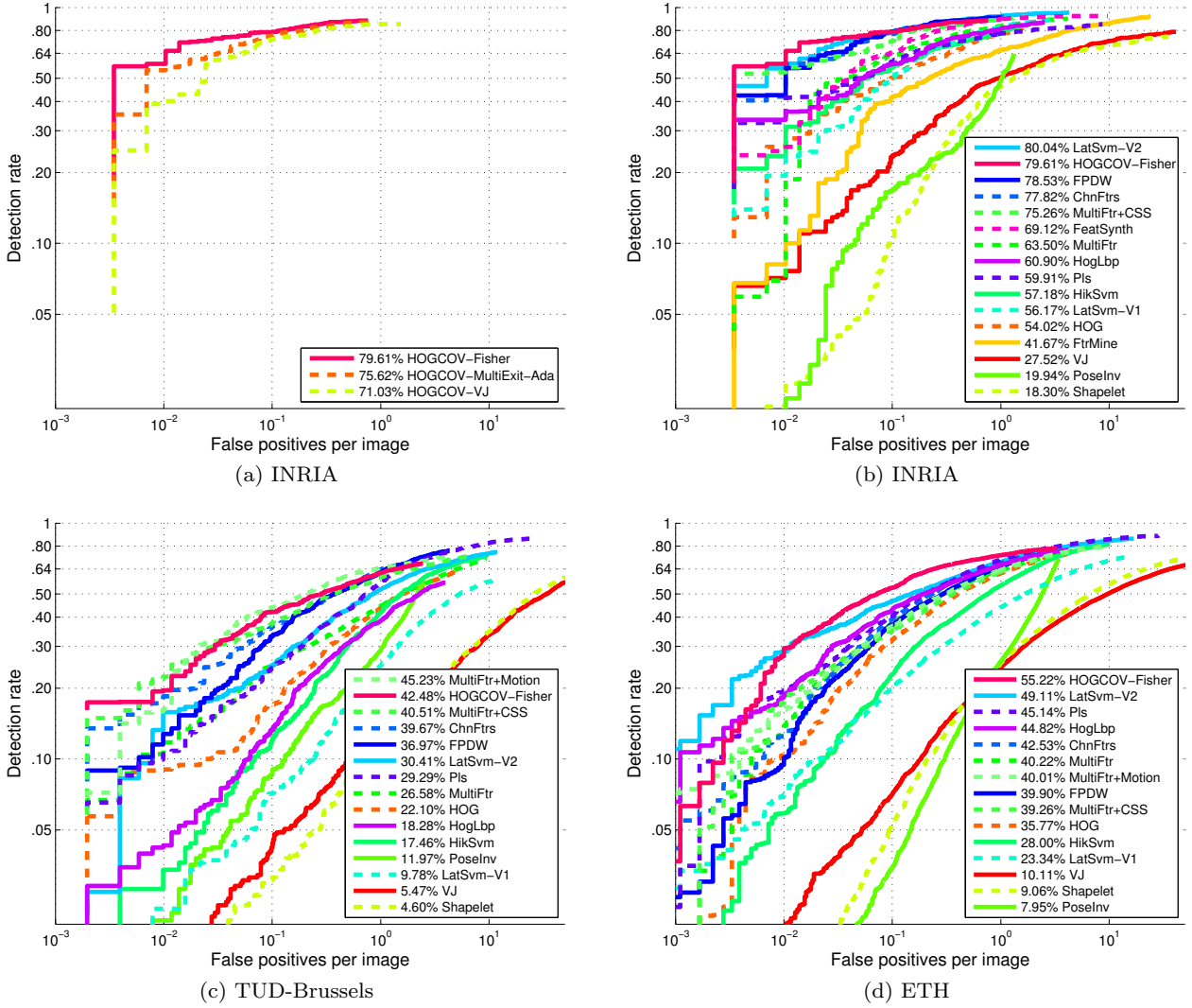


Fig. 6: The performance of our pedestrian detector (HOGCOV-Fisher) compared with (a) baseline detectors and (b, c, d) state-of-the-art detectors on publicly available pedestrian data sets. Our detector uses HOG and covariance features. The performances are ranked using log-average detection rates in the legend. Our detector performs best on the INRIA (Dalal and Triggs 2005) data set, second best on the TUD-Brussels (Wojek et al. 2009) and ETH (Ess et al. 2007) data sets. Note that the best one on the latter two data sets has either used many more features or used a more sophisticated part-based model.

Our detector is then compared with existing pedestrian detectors listed in (Dollár et al. 2012), on the INRIA, TUD-Brussels and ETH data sets. For the TUD-Brussels and ETH data sets, since sizes of ground-truths are smaller than that in INRIA training set, we up-sample the original image to 1280×960 pixels before applying our pedestrian detector. ROC curves and log-average detection rates are reported in Fig. 6 (b), (c) and (d). On the ETH data set, FisherBoost outperforms all the other 14 compared detectors. On the TUD-Brussels data set, our detector is the second best, only inferior to MultiFtr+Motion (Walk et al. 2010) that uses

more discriminative features (gradient, self-similarity and motion) than ours. On the INRIA data set, FisherBoost’s performance is also ranked the second, and only worse than the part-based detector (Felzenszwalb et al. 2010) which uses a much more complex model (deformable part models) and training process (latent SVM). We believe that by further combining with more discriminative features, *e.g.*, CSS features as used in (Walk et al. 2010), the overall detection performance of our method can be further improved. In summary, despite the use of simple HOG plus covariance features, our FisherBoost pedestrian detector still achieves the

	avg. features	frames/sec.
FisherBoost + multi-exit	10.89	0.186
AdaBoost + multi-exit	11.35	0.166
AdaBoost + VJ cascade	21.00	0.109

Table 2: Average features required per detection window and average frames processed per second for different pedestrian detectors on CalTech images of 640×480 pixels (based on our own implementation).

state-of-the-art performance on public benchmark data sets.

Finally, we report an average number of features evaluated per scanning window in Table 2. We compare FisherBoost with our implementation of AdaBoost with the traditional cascade and AdaBoost with the multi-exit cascade. Each image is scanned with 4×4 pixels step stride and 10 scales per octave. There are 90,650 patches to be classified per image. On a single-core Intel i7 CPU 2.8 GHz processor, our detector achieves an average speed of 0.186 frames per second (on 640×480 pixels CalTech images), which is ranked eighth compared with 15 detectors evaluated in (Dollár et al. 2012). Currently, 90% of the total evaluation time is spent on extracting both HOG and covariance features (60% of the evaluation time is spent on extracting raw HOG and covariance features while another 30% of the evaluation time is spent on computing integral images for fast feature calculation during scanning phase).

The major bottleneck of our pedestrian detector lies in the feature extraction part. In our implementation, we make use of multi-threading to speed up the runtime of our pedestrian detector. Using all 8 cores of Intel i7 CPU, we are able to speed up an average processing time to less than 1 second per frame. We believe that by using a special purpose hardware, such as Graphic Processing Unit (GPU), the speed of our detector can be significantly improved.

5.5.1 Discussion

Impact of varying the number of weak classifiers In the next experiment, we vary the number of weak classifiers in each cascade node to evaluate their impact on the final detection performance. We train three different pedestrian detectors (Fisher4/5/6, see Table 3 for details) on the INRIA data set. We limit the maximum number of weak classifiers in each multi-exit node to be 80. The first two nodes is trained using AdaBoost and subsequent nodes are trained using FisherBoost. Fig. 7 shows ROC curves of different detectors. Although we observe a performance improvement as the number of weak classifiers increases, this improvement is minor

compared to a significant increase in the average number of features required per detection window. This experiment indicates the robustness of FisherBoost to the number of weak classifiers in the multi-exit cascade. Note that Fisher5 is used in our previous experiments on pedestrian detection.

Impact of training FisherBoost from an early node In the previous section, we conjecture that FisherBoost performs well when the margin follows the Gaussian distribution. As a result, we apply FisherBoost in the later node of a multi-exit cascade (as these nodes often contain a large number of weak classifiers). In this experiment, we show that it is possible to start training FisherBoost from the first node of the cascade. To achieve this, one can train an additional 50 weak classifiers in the first node (to guarantee the margin approximately follow the Gaussian distribution). We conduct an experiment by training two FisherBoost detectors. In the first detector (Fisher50), FisherBoost is applied from the first node onwards. The number of weak classifiers in each node is 55, 60 (with 55 weak classifiers from the first node), 70 (60 weak classifiers from previous nodes), 80 (70 weak classifiers from previous nodes), etc. In the second detector (Fisher5), we apply AdaBoost in the first two nodes and apply FisherBoost from the third node onwards. The number of weak classifiers in each node is 5, 10 (with 5 weak classifiers from the first node), 20 (10 from previous nodes), 30 (20 from previous nodes), etc. Both detectors use the same node criterion, i.e., each node should discard at least 50% background samples. All other configurations are kept to be the same.

We report the performance of both detectors in Fig. 7. From the results, Fisher50 performs slightly better than Fisher5 (log-average detection rate of 80.38% vs. 79.61%). Based on these results, classifiers in early nodes of the cascade may be heuristically chosen such that a large number of easy negative patches can be quickly discarded. In other words, the first few nodes can significantly affect the efficiency of the visual detector but do not play a significant role in the final detection performance. Actually, one can always apply simple classifiers to remove a large percentage of negative windows to speed up the detection.

5.6 Why LDA Works Better Than LAC

Wu et al. observed that in many cases, LDA post-processing gives better detection rates on MIT+CMU face data than LAC (Wu et al. 2008). When using the LDA criterion to select Haar features, Shen et al. (2011)

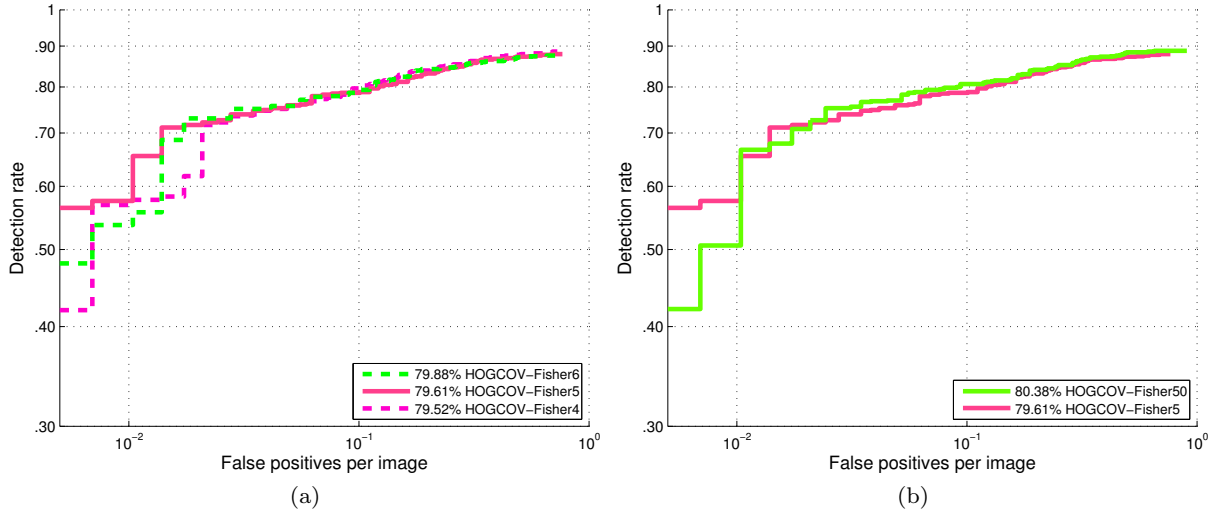


Fig. 7: Performance comparison. (a) We vary the number of weak classifiers in each multi-exit node. When more weak classifiers are used in each node, the accuracy can be slightly improved. (b) We start training FisherBoost from the first node (HOGCOV-Fisher50). HOGCOV-Fisher50 can achieve a slightly better detection rate than HOGCOV-Fisher5.

	Node 1	2	3	4	5	6	7	8	9	10	11	12 onwards	avg. features	log-average det. rate
Fisher4	4	4	8	8	16	16	32	32	64	64	80	80	26.4	79.52%
Fisher5	5	5	10	10	20	20	40	40	80	80	80	80	26.2	79.61%
Fisher6	6	6	12	12	24	24	48	48	80	80	80	80	30.6	79.88%

Table 3: We compare the performance of FisherBoost by varying the number of weak classifiers in each multi-exit node. Average features required per detection window and log-average detection rates on the INRIA pedestrian dataset are reported. When more weak classifiers in each multi-exit node are used, slightly improved accuracy can be achieved at the price of more features being evaluated.

tried different combinations of the two classes' covariance matrices for calculating the within-class matrix: $\mathbf{C}_w = \Sigma_1 + \delta \Sigma_2$ with δ being a nonnegative constant. It is easy to see that $\delta = 1$ and $\delta = 0$ correspond to LDA and LAC, respectively. They found that setting $\delta \in [0.5, 1]$ gives best results on the MIT+CMU face detection task (Paisitkriangkrai et al. 2009; Shen et al. 2011).

According to the analysis in this work, LAC is optimal if the distribution of $[h_1(\mathbf{x}), h_2(\mathbf{x}), \dots, h_n(\mathbf{x})]$ on the negative data is symmetric. In practice, this requirement may not be perfectly satisfied, especially for the first several node classifiers. This may explain why in some cases the improvement of LAC is not significant. However, this does not explain why LDA (FisherBoost) works; and sometimes it performs even better than LAC (LACBoost). At the first glance, LDA (or FisherBoost) by no means explicitly considers the imbalanced node learning objective. Wu et al. did not have a plausible explanation either (Wu et al. 2008; 2005).

Proposition 1 *For object detection problems, the Fisher linear discriminant analysis can be viewed as a regularized version of linear asymmetric classifier. In other words, linear discriminant analysis has already considered the asymmetric learning objective. In FisherBoost, this regularization is equivalent to having a ℓ_2 -norm penalty on the primal variable \mathbf{w} in the objective function of the QP problem in Section 4. Having the ℓ_2 -norm regularization, $\|\mathbf{w}\|_2^2$, avoids over-fitting and increases the robustness of FisherBoost. This similar penalty is also used in machine learning algorithms such as Ridge regression (also known as Tikhonov regularization).*

For object detection such as face and pedestrian detection considered here, the covariance matrix of the negative class is close to a scaled identity matrix. In theory, the negative data can be anything other than the target. Let us look at one of the off-diagonal ele-

	$\delta = 0$ (LACBoost)	$\delta = 0.1$	$\delta = 0.2$	$\delta = 0.5$	$\delta = 1$ (FisherBoost)
Digits	99.12 (0.1)	99.57 (0.2)	99.57 (0.1)	99.55 (0.1)	99.40 (0.1)
Faces	98.63 (0.3)	98.82 (0.3)	98.84 (0.2)	98.48 (0.4)	98.89 (0.2)
Cars	96.80 (1.5)	97.47 (1.1)	97.69 (1.2)	97.96 (1.2)	97.78 (1.2)
Pedestrians	99.12 (0.4)	99.31 (0.1)	99.22 (0.1)	99.13 (0.3)	98.73 (0.3)
Scenes	97.50 (1.1)	98.30 (0.6)	98.62 (0.7)	99.16 (0.4)	99.66 (0.1)
Average (%)	98.23	98.69	98.79	98.86	98.89

Table 4: The average detection rate and its standard deviation (in %) at 50% false positives. We vary the value of δ , which balances the ratio between positive and negative class’s covariance matrices.

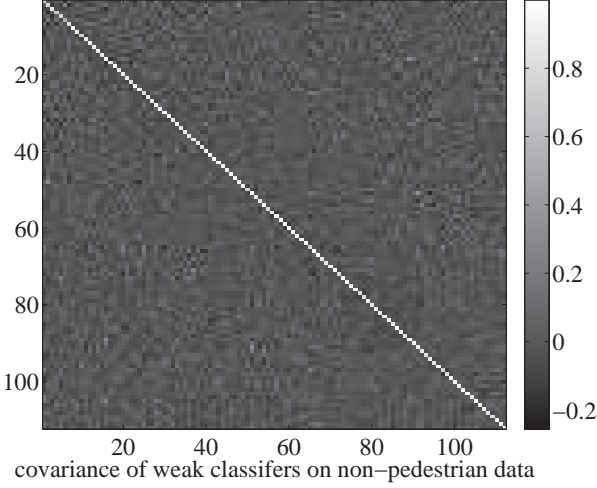


Fig. 8: The covariance matrix of the first 112 weak classifiers selected by FisherBoost on non-pedestrian data. It may be approximated by a scaled identity matrix. On average, the magnitude of diagonal elements is 20 times larger than those off-diagonal elements.

ments

$$\begin{aligned}\Sigma_{ij, i \neq j} &= \mathbb{E}[(h_i(\mathbf{x}) - \mathbb{E}[h_i(\mathbf{x})])(h_j(\mathbf{x}) - \mathbb{E}[h_j(\mathbf{x})])] \\ &= \mathbb{E}[h_i(\mathbf{x})h_j(\mathbf{x})] \approx 0.\end{aligned}\quad (24)$$

Here \mathbf{x} is the image feature of the negative class. We can assume that \mathbf{x} is i.i.d. and approximately, \mathbf{x} follows a symmetric distribution. So $\mathbb{E}[h_{i,j}(\mathbf{x})] = 0$. That is to say, on the negative class, the chance of $h_{i,j}(\mathbf{x}) = +1$ or $h_{i,j}(\mathbf{x}) = -1$ is the same, which is 50%. Note that this does not apply to the positive class because \mathbf{x} of the positive class is not symmetrically distributed, in general. The last equality of (24) uses the fact that weak classifiers $h_i(\cdot)$ and $h_j(\cdot)$ are approximately statistically independent. Although this assumption may not hold in practice as pointed out in (Shen and Li 2010b), it could be a plausible approximation.

Therefore, the off-diagonal elements of Σ are almost all zeros; and Σ is a diagonal matrix. Moreover in object detection, it is a reasonable assumption that the diagonal elements $\mathbb{E}[h_j(\mathbf{x})h_j(\mathbf{x})]$ ($j = 1, 2, \dots$) have

similar values. Hence, $\Sigma_2 \approx v\mathbf{I}$ holds, with v being a small positive constant.

So for object detection, the only difference between LAC and LDA is that, for LAC, $\mathbf{C}_w = \frac{m_1}{m}\Sigma_1$ and for LDA, $\mathbf{C}_w = \frac{m_1}{m}\Sigma_1 + v \cdot \frac{m_2}{m}\mathbf{I}$.

In summary, LDA-like approaches (*e.g.*, LDA post-processing and FisherBoost) perform better than LAC-like approaches (*e.g.*, LAC and LACBoost) in object detection due to two main reasons. The first reason is that LDA is a regularized version of LAC. The second reason is that the negative data are not necessarily symmetrically distributed. Particularly, in latter nodes, bootstrapping forces the negative data to be visually similar the positive data. In this case, ignoring the negative data’s covariance information is likely to deteriorate the detection performance.

Fig. 8 shows some empirical evidence that Σ_2 is close to a scaled identity matrix. As we can see, the diagonal elements are much larger than those off-diagonal elements (off-diagonal ones are close to zeros).

In this experiment, we evaluate the impact of the regularization parameter by varying the value of δ , which balances the ratio between positive and negative class’s covariance matrices, *i.e.*, $\mathbf{C}_w = \Sigma_1 + \delta\Sigma_2$; and also $\mathbf{Q} = \begin{bmatrix} \mathbf{Q}_1 & \mathbf{0} \\ \mathbf{0} & \delta\mathbf{Q}_2 \end{bmatrix}$. Setting $\delta = 0$ corresponds to LACBoost, $\begin{bmatrix} \mathbf{Q}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}$, while setting $\delta = 1$ corresponds to FisherBoost, $\mathbf{Q} = \begin{bmatrix} \mathbf{Q}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{Q}_2 \end{bmatrix}$.

We conduct our experiments on 5 visual data sets by setting the value of δ to be $\{0, 0.1, 0.2, 0.5, 1\}$. All 5 classifiers are trained to remove 50% of the negative data, while retaining almost all positive data. We compare their detection rate in Table 4. First, in general, we observe performance improvement when we set δ to be a small positive value. Since setting δ to be 1 happens to coincide with the LDA objective criterion, the LDA classifier also inherits the node learning goal of LAC in the context of object detection. Second, on different datasets, in theory this parameter should be cross validated and setting it to be 1 (FisherBoost) does not

always give the best performance, which is not surprising.

At this point, a hypothesis naturally arises: *If regularization is really the reason why LACBoost underperforms FisherBoost, then applying other forms of regularization to LACBoost would also be likely to improve LACBoost.* Our last experiment tries to verify this hypothesis.

Here we regularize the matrix \mathbf{Q} by adding an appropriately scaled identity matrix $\mathbf{Q} + \tilde{\delta}\mathbf{I}$. As discussed in Section 4.1, from a numerical stability point of view, difficulties arise when \mathbf{Q} is rank-deficient, causing the dual solution to (18) to be non-uniquely defined. This issue is much worse for LACBoost because the lower-block of \mathbf{Q} , *i.e.*, \mathbf{Q}_2 is a zero matrix. In that case, a well-defined problem can be obtained by replacing \mathbf{Q} with $\mathbf{Q} + \tilde{\delta}\mathbf{I}$. This can be interpreted as corresponding to the primal-regularized QP (refer to (16)):

$$\begin{aligned} \min_{\mathbf{w}, \boldsymbol{\rho}} \quad & \frac{1}{2} \boldsymbol{\rho}^\top \mathbf{Q} \boldsymbol{\rho} - \theta \mathbf{e}^\top \boldsymbol{\rho} + \tilde{\delta} \|\boldsymbol{\rho}\|_2^2, \\ \text{s.t. } \quad & \mathbf{w} \succcurlyeq \mathbf{0}, \mathbf{1}^\top \mathbf{w} = 1, \\ & \rho_i = (\mathbf{A}\mathbf{w})_i, i = 1, \dots, m. \end{aligned} \quad (25)$$

Clearly here in the primal, we are applying the Tikhonov ℓ_2 norm regularization to the variable $\boldsymbol{\rho}$. Also we expect accuracy improvement with this regularization because the margin variance is minimized by minimizing the ℓ_2 norm of the margin while maximizing the weighted mean of the margin, *i.e.*, $\mathbf{e}^\top \boldsymbol{\rho}$. Thus a better margin distribution may be achieved (Shen and Li 2010a;b).

Now we evaluate the impact of the regularization parameter $\tilde{\delta}$ by running experiments on the same datasets as in the last experiment. We vary the values of $\tilde{\delta}$ and the results of detection accuracy are reported in Table 5. Again, the 5 classifiers are trained to remove 50% of the negative data, while correctly classifying as most positive data as possible. As can be seen, indeed, regularization often improves the results. Note that in the experiments, we have solved the primal optimization problem so that even when \mathbf{Q} is not invertible, we can still obtain a solution. Having the primal solutions, the dual solutions are obtained using (23). This experiment demonstrates that other formats of regularization indeed improves LACBoost too.

6 Conclusion

By explicitly taking into account the node learning goal in cascade classifiers, we have designed new boosting algorithms for more effective object detection.

Experiments validate the superiority of the methods developed, which we have labeled FisherBoost and

LACBoost. We have also proposed the use of entropic gradient descent to efficiently implement FisherBoost and LACBoost. The proposed algorithms are easy to implement and can be applied to other asymmetric classification tasks in computer vision. We aim in future to design new asymmetric boosting algorithms by exploiting asymmetric kernel classification methods such as (Tu and Lin 2010). Compared with stage-wise Adaboost, which is parameter-free, our boosting algorithms need to tune a parameter.

We are also interested in developing parameter-free stage-wise boosting that considers the node learning objective. Moreover, the developed boosting algorithms only work for the case $\gamma_o \leq 0.5$ in (2). How can we make it work for $\gamma_o \geq 0.5$? Last, to relax the symmetric distribution requirement for the feature responses of the negative class is also a topic of interest.

References

- S. Agarwal, A. Awan, and D. Roth. Learning to detect objects in images via a sparse, part-based representation. *IEEE Trans. Pattern Anal. Mach. Intell.*, 26(11):1475–1490, 2004.
- D. Aldavert, A. Ramisa, R. Toledo, and R. Lopez de Mantaras. Fast and robust object segmentation with the integral linear classifier. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, San Francisco, US, 2010.
- A. Beck and M. Teboulle. Mirror descent and nonlinear projected subgradient methods for convex optimization. *Oper. Res. Lett.*, 31(3):167–175, 2003.
- J. Bi, S. Periaswamy, K. Okada, T. Kubota, G. Fung, M. Salganicoff, and R. B. Rao. Computer aided detection via asymmetric cascade of sparse hyperplane classifiers. In *Proc. ACM Int. Conf. Knowledge Discovery & Data Mining*, pages 837–844, Philadelphia, PA, USA, 2006.
- L. Bourdev and J. Brandt. Robust object detection via soft cascade. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pages 236–243, San Diego, CA, US, 2005.
- S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, 2004.
- S. C. Brubaker, J. Wu, J. Sun, M. D. Mullin, and J. M. Rehg. On the design of cascades of boosted ensembles for face detection. *Int. J. Comp. Vis.*, 77(1–3):65–86, 2008.
- M. Collins, A. Globerson, T. Koo, X. Carreras, and P. L. Bartlett. Exponentiated gradient algorithms for conditional random fields and max-margin Markov networks. *J. Mach. Learn. Res.*, pages 1775–1822, 2008.
- N. Dalal and B. Triggs. Histograms of oriented gradients for human detection. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, volume 1, pages 886–893, San Diego, CA, 2005.
- A. Demiriz, K.P. Bennett, and J. Shawe-Taylor. Linear programming boosting via column generation. *Mach. Learn.*, 46(1-3):225–254, 2002.
- P. Dollár. Piotr’s image and video Matlab toolbox. <http://vision.ucsd.edu/~pdollar/toolbox/doc/>, 2012.
- P. Dollár, B. Babenko, S. Belongie, P. Perona, and Z. Tu. Multiple component learning for object detection. In *Proc. Eur. Conf. Comp. Vis.*, pages 211–224, Marseille, France, 2008.

	$\tilde{\delta} = 0$ (LACBoost)	$\tilde{\delta} = 5 \times 10^{-4}$	$\tilde{\delta} = 2 \times 10^{-4}$	$\tilde{\delta} = 10^{-4}$	$\tilde{\delta} = 5 \times 10^{-5}$	$\tilde{\delta} = 2 \times 10^{-5}$	$\tilde{\delta} = 10^{-5}$
Digits	99.12 (0.1)	99.50 (0.1)	99.41 (0.2)	99.59 (0.2)	99.60 (0.2)	99.50 (0.3)	99.11 (0.5)
Faces	98.63 (0.3)	98.73 (0.0)	98.87 (0.0)	99.02 (0.0)	98.38 (0.0)	98.84 (0.0)	99.04 (0.0)
Cars	96.80 (1.5)	96.62 (1.5)	96.80 (1.5)	96.80 (1.5)	96.67 (1.4)	96.58 (1.5)	96.71 (1.5)
Pedestrians	99.12 (0.4)	96.62 (1.5)	99.32 (0.1)	99.22 (0.2)	98.97 (0.4)	98.97 (0.3)	98.81 (0.4)
Scenes	97.50 (1.1)	98.96 (0.4)	98.36 (0.5)	97.88 (0.6)	98.38 (0.6)	98.25 (0.8)	97.1 (0.8)
Average (%)	98.23	98.09	98.55	98.50	98.40	98.43	98.20

Table 5: The average detection rate and its standard deviation (in %) at 50% false positives of various regularized LACBoosts. We vary the value of $\tilde{\delta}$, i.e., $\mathbf{Q} + \tilde{\delta}\mathbf{I}$. Regularization often improves the overall detection accuracy.

- P. Dollár, C. Wojek, B. Schiele, and P. Perona. Pedestrian detection: An evaluation of the state of the art. *IEEE Trans. Pattern Anal. Mach. Intell.*, 34(4):743–761, 2012.
- M. Dundar and J. Bi. Joint optimization of cascaded classifiers for computer aided detection. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, Minneapolis, MN, USA, 2007.
- M. Enzweiler, A. Eigenstetter, B. Schiele, and D. M. Gavrila. Multi-cue pedestrian classification with partial occlusion handling. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, San Francisco, US, 2010.
- A. Ess, B. Leibe, and L. Van Gool. Depth and appearance for mobile scene analysis. In *Proc. IEEE Int. Conf. Comp. Vis.*, 2007.
- P. Felzenszwalb, R. Girshick, D. McAllester, and D. Ramanan. Object detection with discriminatively trained part based models. *IEEE Trans. Pattern Anal. Mach. Intell.*, 32(9):1627–1645, 2010.
- P. Gehler and S. Nowozin. On feature combination for multi-class object classification. In *Proc. IEEE Int. Conf. Comp. Vis.*, 2009.
- K. Huang, H. Yang, I. King, M. Lyu, and L. Chan. The minimum error minimax probability machine. *J. Mach. Learn. Res.*, 5:1253–1286, Dec. 2004.
- G. R. G. Lanckriet, L. El Ghaoui, C. Bhattacharyya, and M. I. Jordan. A robust minimax approach to classification. *J. Mach. Learn. Res.*, 3:555–582, Dec. 2002.
- S. Lazebnik, C. Schmid, and J. Ponce. Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, New York City, USA, 2006.
- L. Lefakis and F. Fleuret. Joint cascade optimization using a product of boosted classifiers. In *Proc. Adv. Neural Inf. Process. Syst.*, 2010.
- S. Z. Li and Z. Zhang. FloatBoost learning and statistical face detection. *IEEE Trans. Pattern Anal. Mach. Intell.*, 26(9):1112–1123, 2004.
- Z. Lin, G. Hua, and L. S. Davis. Multiple instance feature for robust part-based object detection. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pages 405–412, Miami, FL, US, 2009.
- C. Liu and H.-Y. Shum. Kullback-Leibler boosting. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, volume 1, pages 587–594, Madison, Wisconsin, June 2003.
- S. Maji, A. C. Berg, and J. Malik. Classification using intersection kernel support vector machines is efficient. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, Anchorage, AK, US, 2008.
- H. Masnadi-Shirazi and N. Vasconcelos. Asymmetric boosting. In *Proc. Int. Conf. Mach. Learn.*, pages 609–619, Corvallis, Oregon, US, 2007.
- H. Masnadi-Shirazi and N. Vasconcelos. Cost-sensitive boosting. *IEEE Trans. Pattern Anal. Mach. Intell.*, 33(2):294–309, 2011.
- MOSEK. The MOSEK optimization toolbox for matlab manual, version 6.0, revision 93, 2010. <http://www.mosek.com/>.
- Y. Mu, S. Yan, Y. Liu, T. Huang, and B. Zhou. Discriminative local binary patterns for human detection in personal album. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, Anchorage, AK, US, 2008.
- S. Munder and D. M. Gavrila. An experimental study on pedestrian classification. *IEEE Trans. Pattern Anal. Mach. Intell.*, 28(11):1863–1868, 2006.
- S. Paisitkriangkrai, C. Shen, and J. Zhang. Fast pedestrian detection using a cascade of boosted covariance features. *IEEE Trans. Circuits Syst. Video Technol.*, 18(8):1140–1151, 2008.
- S. Paisitkriangkrai, C. Shen, and J. Zhang. Efficiently training a better visual detector with sparse Eigenvectors. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, Miami, Florida, US, June 2009.
- M.-T. Pham and T.-J. Cham. Fast training and selection of Haar features using statistics in boosting-based face detection. In *Proc. IEEE Int. Conf. Comp. Vis.*, Rio de Janeiro, Brazil, 2007a.
- M.-T. Pham and T.-J. Cham. Online learning asymmetric boosted classifiers for object detection. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, Minneapolis, MN, 2007b.
- M.-T. Pham, V.-D. D. Hoang, and T.-J. Cham. Detection with multi-exit asymmetric boosting. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, Anchorage, Alaska, 2008.
- G. Rätsch, S. Mika, B. Schölkopf, and K.-R. Müller. Constructing boosting algorithms from SVMs: An application to one-class classification. *IEEE Trans. Pattern Anal. Mach. Intell.*, 24(9):1184–1199, 2002.
- M. Saberian and N. Vasconcelos. Learning optimal embedded cascades. *IEEE Trans. Pattern Anal. Mach. Intell.*, 2012.
- M.J. Saberian and N. Vasconcelos. Boosting classifier cascades. In *Proc. Adv. Neural Inf. Process. Syst.*, 2010.
- C. Shen and H. Li. Boosting through optimization of margin distributions. *IEEE Trans. Neural Networks*, 21(4):659–666, 2010a.
- C. Shen and H. Li. On the dual formulation of boosting algorithms. *IEEE Trans. Pattern Anal. Mach. Intell.*, 32(12):2216–2231, 2010b. IEEE computer Society Digital Library. <http://dx.doi.org/10.1109/TPAMI.2010.47>.
- C. Shen, S. Paisitkriangkrai, and J. Zhang. Face detection from few training examples. In *Proc. Int. Conf. Image Process.*, pages 2764–2767, San Diego, California, USA, 2008.
- C. Shen, P. Wang, and H. Li. LACBoost and FisherBoost: Optimally building cascade classifiers. In *Proc. Eur. Conf. Comp. Vis.*, volume 2, LNCS 6312, pages 608–621, Crete Island, Greece, 2010.
- C. Shen, S. Paisitkriangkrai, and J. Zhang. Efficiently learning a detection cascade with sparse Eigenvectors. *IEEE Trans. Image Process.*, 20(1):2235, 2011. <http://dx.doi.org/10.1109/TIP.2010.2055880>.

J. Sochman and J. Matas. Waldboost - learning for time constrained sequential detection. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2005.

A. Torralba, K. P. Murphy, and W. T. Freeman. Sharing visual features for multiclass and multiview object detection. *IEEE Trans. Pattern Anal. Mach. Intell.*, 29(5):854–869, 2007.

H.-H. Tu and H.-T. Lin. One-sided support vector regression for multiclass cost-sensitive classification. In *Proc. Int. Conf. Mach. Learn.*, Haifa, Israel, 2010.

O. Tuzel, F. Porikli, and P. Meer. Pedestrian detection via classification on Riemannian manifolds. *IEEE Trans. Pattern Anal. Mach. Intell.*, 30(10):1713–1727, 2008.

P. Viola and M. Jones. Fast and robust classification using asymmetric AdaBoost and a detector cascade. In *Proc. Adv. Neural Inf. Process. Syst.*, pages 1311–1318. MIT Press, 2002.

P. Viola and M. J. Jones. Robust real-time face detection. *Int. J. Comp. Vis.*, 57(2):137–154, 2004.

P. Viola, J. C. Platt, and C. Zhang. Multiple instance boosting for object detection. In *Proc. Adv. Neural Inf. Process. Syst.*, pages 1417–1424, Vancouver, Canada, 2005.

S. Walk, N. Majer, K. Schindler, and B. Schiele. New features and insights for pedestrian detection. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, San Francisco, US, 2010.

P. Wang, C. Shen, N. Barnes, and H. Zheng. Fast and robust object detection using asymmetric totally-corrective boosting. *IEEE Trans. Neural Networks & Learn. Syst.*, 23(1):33–46, 2012.

W. Wang, J. Zhang, and C. Shen. Improved human detection and classification in thermal images. In *Proc. Int. Conf. Image Process.*, Hong Kong, 2010.

X. Wang, T. X. Han, and S. Yan. An HOG-LBP human detector with partial occlusion handling. In *Proc. IEEE Int. Conf. Comp. Vis.*, Rio de Janeiro, Brazil, 2007.

C. Wojek, S. Walk, and B. Schiele. Multi-cue onboard pedestrian detection. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2009.

B. Wu and R. Nevatia. Optimizing discrimination-efficiency tradeoff in integrating heterogeneous local features for object detection. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, Anchorage, AK, US, 2008.

J. Wu and J. M. Rehg. CENTRIST: A visual descriptor for scene categorization. *IEEE Trans. Pattern Anal. Mach. Intell.*, 33(8):1489–1501, 2011.

J. Wu, J. M. Rehg, and M. D. Mullin. Learning a rare event detection cascade by direct feature selection. In S. Thrun, L. Saul, and B. Schölkopf, editors, *Proc. Adv. Neural Inf. Process. Syst.*, 2003.

J. Wu, M. D. Mullin, and J. M. Rehg. Linear asymmetric classifier for cascade detectors. In *Proc. Int. Conf. Mach. Learn.*, pages 988–995, Bonn, Germany, 2005.

J. Wu, S. C. Brubaker, M. D. Mullin, and J. M. Rehg. Fast asymmetric learning for cascade face detection. *IEEE Trans. Pattern Anal. Mach. Intell.*, 30(3):369–382, 2008.

R. Xiao, L. Zhu, and H.-J. Zhang. Boosting chain learning for object detection. In *Proc. IEEE Int. Conf. Comp. Vis.*, pages 709–715, Nice, France, 2003.

R. Xiao, H. Zhu, H. Sun, and X. Tang. Dynamic cascades for face detection. In *Proc. IEEE Int. Conf. Comp. Vis.*, Rio de Janeiro, Brazil, 2007.

Y.-L. Yu, Y. Li, D. Schuurmans, and C. Szepesvári. A general projection property for distribution families. In Y. Bengio, D. Schuurmans, J. Lafferty, C. K. I. Williams, and A. Culotta, editors, *Proc. Adv. Neural Inf. Process. Syst.*, pages 2232–2240, Vancouver, Canada, 2009.

Y. Zheng, C. Shen, R. Hartley, and X. Huang. Pyramid center-symmetric local binary, trinary patterns for effective pedestrian detection. In *Proc. Asian Conf. Comp. Vis.*, New Zealand, 2010.

Q. Zhu, S. Avidan, M.-C. Yeh, and K.-T. Cheng. Fast human detection using a cascade of histograms of oriented gradients. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pages 1491–1498, New York City, USA, 2006.

A Proof of Theorem 1

Before we present our results, we introduce an important proposition from (Yu et al. 2009). Note that we have used different notation.

Proposition 2 *For a few different distribution families, the worst-case constraint*

$$\left[\inf_{\mathbf{x} \sim (\boldsymbol{\mu}, \boldsymbol{\Sigma})} \Pr\{\mathbf{w}^\top \mathbf{x} \leq b\} \right] \geq \gamma, \quad (26)$$

can be written as:

1. if $\mathbf{x} \sim (\boldsymbol{\mu}, \boldsymbol{\Sigma})$, i.e., \mathbf{x} follows an arbitrary distribution with mean $\boldsymbol{\mu}$ and covariance $\boldsymbol{\Sigma}$, then

$$b \geq \mathbf{w}^\top \boldsymbol{\mu} + \sqrt{\frac{\gamma}{1-\gamma}} \cdot \sqrt{\mathbf{w}^\top \boldsymbol{\Sigma} \mathbf{w}}; \quad (27)$$

2. if $\mathbf{x} \sim (\boldsymbol{\mu}, \boldsymbol{\Sigma})_S$,⁶ then we have

$$\begin{cases} b \geq \mathbf{w}^\top \boldsymbol{\mu} + \sqrt{\frac{1}{2(1-\gamma)}} \cdot \sqrt{\mathbf{w}^\top \boldsymbol{\Sigma} \mathbf{w}}, & \text{if } \gamma \in (0.5, 1); \\ b \geq \mathbf{w}^\top \boldsymbol{\mu}, & \text{if } \gamma \in (0, 0.5]; \end{cases} \quad (28)$$

3. if $\mathbf{x} \sim (\boldsymbol{\mu}, \boldsymbol{\Sigma})_{SU}$, then

$$\begin{cases} b \geq \mathbf{w}^\top \boldsymbol{\mu} + \frac{2}{3} \sqrt{\frac{1}{2(1-\gamma)}} \cdot \sqrt{\mathbf{w}^\top \boldsymbol{\Sigma} \mathbf{w}}, & \text{if } \gamma \in (0.5, 1); \\ b \geq \mathbf{w}^\top \boldsymbol{\mu}, & \text{if } \gamma \in (0, 0.5]; \end{cases} \quad (29)$$

4. if \mathbf{x} follows a Gaussian distribution with mean $\boldsymbol{\mu}$ and covariance $\boldsymbol{\Sigma}$, i.e., $\mathbf{x} \sim \mathcal{G}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, then

$$b \geq \mathbf{w}^\top \boldsymbol{\mu} + \phi^{-1}(\gamma) \cdot \sqrt{\mathbf{w}^\top \boldsymbol{\Sigma} \mathbf{w}}, \quad (30)$$

where $\phi(\cdot)$ is the cumulative distribution function (c.d.f.) of the standard normal distribution $\mathcal{G}(0, 1)$, and $\phi^{-1}(\cdot)$ is the inverse function of $\phi(\cdot)$.

Two useful observations about $\phi^{-1}(\cdot)$ are: $\phi^{-1}(0.5) = 0$; and $\phi^{-1}(\cdot)$ is a monotonically increasing function in its domain.

⁶ Here $(\boldsymbol{\mu}, \boldsymbol{\Sigma})_S$ denotes the family of distributions in $(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ that are also symmetric about the mean $\boldsymbol{\mu}$. $(\boldsymbol{\mu}, \boldsymbol{\Sigma})_{SU}$ denotes the family of distributions in $(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ that are additionally symmetric and linear unimodal about $\boldsymbol{\mu}$.

We omit the proof of Proposition 2 here and refer the reader to (Yu et al. 2009) for details. Next we begin to prove Theorem 1:

Proof The second constraint of (2) is simply

$$b \geq \mathbf{w}^\top \boldsymbol{\mu}_2. \quad (31)$$

The first constraint of (2) can be handled by writing $\mathbf{w}^\top \mathbf{x}_1 \geq b$ as $-\mathbf{w}^\top \mathbf{x}_1 \leq -b$ and applying the results in Proposition 2. It can be written as

$$-b + \mathbf{w}^\top \boldsymbol{\mu}_1 \geq \varphi(\gamma) \sqrt{\mathbf{w}^\top \Sigma_1 \mathbf{w}}, \quad (32)$$

with (6).

Let us assume that Σ_1 is strictly positive definite (if it is only positive semidefinite, we can always add a small regularization to its diagonal components). From (32) we have

$$\varphi(\gamma) \leq \frac{-b + \mathbf{w}^\top \boldsymbol{\mu}_1}{\sqrt{\mathbf{w}^\top \Sigma_1 \mathbf{w}}}. \quad (33)$$

So the optimization problem becomes

$$\max_{\mathbf{w}, b, \gamma} \gamma, \text{ s.t. (31) and (33)}. \quad (34)$$

The maximum value of γ (which we label γ^*) is achieved when (33) is strictly an equality. To illustrate this point, let us assume that the maximum is achieved when

$$\varphi(\gamma^*) < \frac{-b + \mathbf{w}^\top \boldsymbol{\mu}_1}{\sqrt{\mathbf{w}^\top \Sigma_1 \mathbf{w}}}.$$

Then a new solution can be obtained by increasing γ^* with a positive value such that (33) becomes an equality. Notice that the constraint (31) will not be affected, and the new solution will be better than the previous one. Hence, at the optimum, (5) must be fulfilled.

Because $\varphi(\gamma)$ is monotonically increasing for all the four cases in its domain $(0, 1)$ (see Fig. 9), maximizing γ is equivalent to maximizing $\varphi(\gamma)$ and this results in

$$\max_{\mathbf{w}, b} \frac{-b + \mathbf{w}^\top \boldsymbol{\mu}_1}{\sqrt{\mathbf{w}^\top \Sigma_1 \mathbf{w}}}, \text{ s.t. } b \geq \mathbf{w}^\top \boldsymbol{\mu}_2. \quad (35)$$

As in (Lanckriet et al. 2002; Huang et al. 2004), we also have a scale ambiguity: if (\mathbf{w}^*, b^*) is a solution, $(t\mathbf{w}^*, tb^*)$ with $t > 0$ is also a solution.

An important observation is that the problem (35) must attain the optimum at (4). Otherwise if $b > \mathbf{w}^\top \boldsymbol{\mu}_2$, the optimal value of (35) must be smaller. So we can rewrite (35) as an unconstrained problem (3).

We have thus shown that, if \mathbf{x}_1 is distributed according to a symmetric, symmetric unimodal, or Gaussian distribution, the resulting optimization problem is

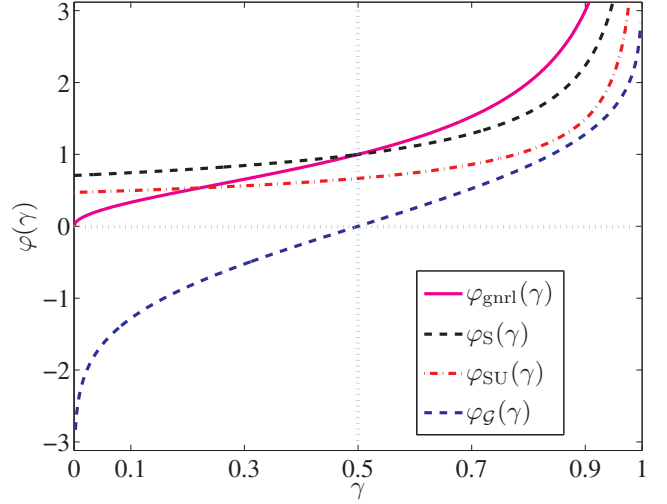


Fig. 9: The function $\varphi(\cdot)$ in (6). The four curves correspond to the four cases. They are all monotonically increasing in $(0, 1)$.

identical. This is not surprising considering the latter two cases are merely special cases of the symmetric distribution family.

At optimality, the inequality (33) becomes an equality, and hence γ^* can be obtained as in (5). For ease of exposition, let us denote the four cases in the right side of (6) as $\varphi_{\text{gnrl}}(\cdot)$, $\varphi_S(\cdot)$, $\varphi_{\text{SU}}(\cdot)$, and $\varphi_G(\cdot)$. For $\gamma \in [0.5, 1)$, as shown in Fig. 9, we have $\varphi_{\text{gnrl}}(\gamma) > \varphi_S(\gamma) > \varphi_{\text{SU}}(\gamma) > \varphi_G(\gamma)$. Therefore, when solving (5) for γ^* , we have $\gamma_{\text{gnrl}}^* < \gamma_S^* < \gamma_{\text{SU}}^* < \gamma_G^*$. That is to say, one can get better accuracy when additional information about the data distribution is available, although the actual optimization problem to be solved is identical.

B Proof of Theorem 2

Let us assume that in the current solution we have selected n weak classifiers and their corresponding linear weights are $\mathbf{w} = [w_1, \dots, w_n]$. If we add a weak classifier $h'(\cdot)$ that is not in the current subset, the corresponding w is zero, then we can conclude that the current weak classifiers and \mathbf{w} are the optimal solution already. In this case, the best weak classifier that is found by solving the subproblem (20) does not contribute to solving the master problem.

Let us consider the case that the optimality condition is violated. We need to show that we are able to find such a weak learner $h'(\cdot)$, which is not in the set of current selected weak classifiers, that its corresponding coefficient $w > 0$ holds. Again assume $h'(\cdot)$ is the most

violated weak learner found by solving (20) and the convergence condition is not satisfied. In other words, we have

$$\sum_{i=1}^m u_i y_i h'(\mathbf{x}_i) \geq r. \quad (36)$$

Now, after this weak learner is added into the master problem, the corresponding primal solution w must be non-zero (positive because we have the nonnegativity constraint on \mathbf{w}).

If this is not the case, then the corresponding $w = 0$. This is not possible because of the following reason. From the Lagrangian (17), at optimality we have $\partial L / \partial \mathbf{w} = \mathbf{0}$, which leads to

$$r - \sum_{i=1}^m u_i y_i h'(\mathbf{x}_i) = q > 0. \quad (37)$$

Clearly (36) and (37) contradict.

Thus, after the weak classifier $h'(\cdot)$ is added to the primal problem, its corresponding w must have a positive solution. This is to say, one more free variable is added into the problem and re-solving the primal problem (16) must reduce the objective value. Therefore a strict decrease in the objective is obtained. In other words, Algorithm 1 must make progress at each iteration. Furthermore, the primal optimization problem is convex, there are no local optimal points. The column generation procedure is guaranteed to converge to the global optimum up to some prescribed accuracy.

C Exponentiated Gradient Descent

Exponentiated Gradient Descent (EG) is a very useful tool for solving large-scale convex minimization problems over the unit simplex. Let us first define the unit simplex $\Delta_n = \{\mathbf{w} \in \mathbb{R}^n : \mathbf{1}^\top \mathbf{w} = 1, \mathbf{w} \succeq \mathbf{0}\}$. EG efficiently solves the convex optimization problem

$$\min_{\mathbf{w}} f(\mathbf{w}), \text{ s.t. } \mathbf{w} \in \Delta_n, \quad (38)$$

under the assumption that the objective function $f(\cdot)$ is a convex Lipschitz continuous function with Lipschitz constant L_f w.r.t. a fixed given norm $\|\cdot\|$. The mathematical definition of L_f is that $|f(\mathbf{w}) - f(\mathbf{z})| \leq L_f \|\mathbf{x} - \mathbf{z}\|$ holds for any \mathbf{x}, \mathbf{z} in the domain of $f(\cdot)$. The EG algorithm is very simple:

1. Initialize with $\mathbf{w}^0 \in$ the interior of Δ_n ;
2. Generate the sequence $\{\mathbf{w}^k\}$, $k = 1, 2, \dots$ with:

$$\mathbf{w}_j^k = \frac{\mathbf{w}_j^{k-1} \exp[-\tau_k f'_j(\mathbf{w}^{k-1})]}{\sum_{j=1}^n \mathbf{w}_j^{k-1} \exp[-\tau_k f'_j(\mathbf{w}^{k-1})]}. \quad (39)$$

Here τ_k is the step-size. $f'(\mathbf{w}) = [f'_1(\mathbf{w}), \dots, f'_n(\mathbf{w})]^\top$ is the gradient of $f(\cdot)$;

3. Stop if some stopping criteria are met.

The learning step-size can be determined by

$$\tau_k = \frac{\sqrt{2 \log n}}{L_f} \frac{1}{\sqrt{k}},$$

following (Beck and Teboulle 2003). In (Collins et al. 2008), the authors have used a simpler strategy to set the learning rate.

In EG there is an important parameter L_f , which is used to determine the step-size. L_f can be determined by the ℓ_∞ -norm of $|f'(\mathbf{w})|$. In our case $f'(\mathbf{w})$ is a linear function, which is trivial to compute. The convergence of EG is guaranteed; see (Beck and Teboulle 2003) for details.