

基于多阈值 Boosting 方法的人脸检测

钟向阳¹, 凌捷²

(1. 嘉应学院计算机系, 梅州 514015; 2. 广东工业大学计算机学院, 广州 510090)

摘要: Adaboost 算法采用单阈值弱分类器, 难以拟合复杂分布, 其训练过程收敛速度较慢。针对该问题设计一种多阈值弱学习器, 利用平方和减少最大化准则划分节点并生成弱分类器, 在训练数据集上采用 GAB 算法将弱分类器提升为强分类器。实验结果表明, 在弱分类器数目相同的情况下, 该方法的正样本误报率低于 Adaboost 算法。

关键词: 人脸检测; boosting 方法; 实值 Adaboost; 平缓 Adaboost

Face Detection Based on Multiple Thresholds Boosting Method

ZHONG Xiang-yang¹, LING Jie²

(1. Department of Computer, Jiaying University, Meizhou 514015;

2. School of Computer, Guangdong University of Technology, Guangzhou 510090)

【Abstract】 Aiming at the problem that the Adaboost algorithm using simple threshold weak classifiers is too weak to fit complex distributions and its slow convergence rate in training, this paper designs a multiple thresholds weak learner. This learner splits the nodes by using biggest reduction in the sum of squares as the partition criteria and builds a weak classifier. It boosts weak classifiers using GAB algorithm on training dataset. Experimental results show that the false positive rate of this method is lower than Adaboost algorithm at the same number of weak classifiers.

【Key words】 face detection; Boosting method; real Adaboost; gentle Adaboost

1 概述

人脸检测具有重要意义, 可以应用到人脸识别、新一代人机界面、安全访问和视觉监控、基于内容的检索等领域。文献[1]提出一种综合了矩形特征、积分图、Adaboost 算法和 Cascade 分类器的实时人脸检测系统。Adaboost 算法是一种 DAB(Discrete Adaboost)算法, 其目标是从弱分类器空间中自动挑选若干个弱分类器, 并将其整合成一个强分类器, 使生成的强分类器具有比单个学习器更好的性能。采用 DAB 方法需要的训练时间较长。文献[2]在此基础上提出一种实值 RAB(Real Adaboost)算法, 该算法将 DAB 算法从处理二值判定推广到连续置信度输出, 能更精确地逼近实际的指数误差函数且收敛速度更快。文献[3]基于 RAB 算法对 Haar 型特征值域采用等距划分, 并用查找表方法实现多视角人脸检测。文献[4]将人脸检测视为在较高检测率和较低正样本误报率的情况下, 构造一个有效检测器的过程, 它采用前向特征选择方法提高了特征训练速度, 并取得了与 Viola-Jones 检测器相似的性能。KLBoosting 方法以 Kullback-Leibler 离散度最大化作为特征选择准则, 通过计算人脸和非人脸之间对称的 KL 离散度, 以最优 KL 特征构造最优分类器, 但 KL 特征的计算量较大。本文基于 RAB 思想, 采用多阈值和 GAB(Gentle Adaboost)^[5]度量方法设计了一种人脸检测器, 在速度和性能方面取得了较好效果。

2 基于实值 Adaboost 的强分类器

本文根据文献[4], 将强分类器的设计看成是在较高检测率 dr (detection rate) 和较低正样本误报率 fpr (false positive rate) 情况下, 构造一个有效检测器的过程, 即给定 dr 求 Bayes 代价函数 $\lambda \cdot (1 - dr) + fpr$ 最小化的过程, 其中, Bayes 风险系数 λ 设为 1.0。强分类器的训练算法如下:

(1) 给定期望强分类器的最小检测率 dr 和最大误检率 fpr , 输入 m 个训练样本 $(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)$, 其中, 人脸正样本 P 集有 $npos$ 个; 非人脸负样本 N 集有 $nneg$ 个; $y = -1$ 或 $+1$ 分别对应负样本和正样本。

(2) 对每个正样本赋初始权值 $w_p = 1 / (2 \times npos)$, 每个负样本的权值 $w_n = 1 / (2 \times nneg)$, 当前强分类器的误检率 $fpr_temp = 1.0$, 当前弱分类器个数 $t = 1$; $H \leftarrow \phi$ 。

(3) 当 $fpr_temp > fpr$ && $t \leq T_{max}$ 时, 执行如下操作:

1) 调用弱分类器训练算法构造一个弱分类器 $f_t(x)$;

2) 更新训练集上每个样本的权值 $w_{i+1} = w_i \times e^{-y_i f_t(x_i)}$;

3) 权值归一化 $w_{i+1} = w_i / \sum_{i=1}^m w_i$;

4) 令 $F(x) = F(x) + f_t(x)$;

5) 在验证集上对当前强分类器进行检测, 对每个正样本 x^p , 计算数组 $PV = \sum_{i=1}^t f_i(x^p)$, 对于每个负样本 x^n , 计算数组 $NV = \sum_{i=1}^t f_i(x^n)$;

6) 对 PV 数组按从小到大排序, 并令当前强分类器的阈值 $\varphi = PV[npos \cdot (1 - dr)]$, 保证有 $npos \times dr$ 个正样本能通过该阈值;

7) 统计 NV 数组中大于阈值 φ 的负样本数目 $numfalse$, 计算当前强分类器的误检率 $fpr_temp = numfalse / nneg$;

8) 令 $t = t + 1$ 。

基金项目: 广东省科技计划基金资助项目(2007B010200071, 2006B11201014)

作者简介: 钟向阳(1969—), 男, 讲师、硕士, 主研方向: 图像处理, 模式识别; 凌捷, 教授、博士

收稿日期: 2008-12-12 E-mail: zxy@jyu.edu.cn

在上算法中, T_{\max} 为给定的最大迭代次数, 由此生成强分类器 $H(x) = \text{sgn}(\sum_{t=1}^T f_t(x) + \varphi)$ 。

3 弱分类器的设计

弱分类器的训练有多种不同方法。可以根据分类错误最小准则选择最优的 Haar 型特征, 从而生成一个基于阈值的二值型弱分类器。但正负样本特征的实际分布不一定是高斯分布, 对较复杂的函数如混合高斯分布, 用二值型弱分类器逼近会产生较大误差。因此, 本文依据 RAB 方法的思想, 对 Haar 型特征值域用 $N-1$ 个最优阈值作为分界点进行划分, 由划分生成的多个基本特征迭加构造一个弱分类器, 主要包括如下 2 个方面: (1) 区间划分策略; (2) 划分区间内样本特征置信度的估计方法。

3.1 区间划分策略

区间划分策略是对 Haar 型特征值域的划分方法, 文献[3]对 Haar 型特征的值域采用 N 等距划分方法, 若 N 的数量太大, 将降低特征检测速度。本文采用误差测度减少最大化准则对特征值域进行划分, 用尽可能少的划分区间减少整个分类误差。误差测度减少最大化划分准则描述如下:

定义误差测度为 $\varepsilon(t) = \sum_{i=1}^m w_i (y_i - \alpha_i)^2$, 其中, y_i 为已知样本值; α_i 为样本估计值; w_i 为样本权重。设一个节点 t 代表当前树 T 的一个子集, S 为当前节点 t 的一个分割集, 若 S 分割为 t_L, t_R 左右 2 个子节点, 则误差减少程度 $\Delta\varepsilon(s, t) = \varepsilon(s) - \varepsilon(t_L) - \varepsilon(t_R)$ 。最好的分割 s^* 是使 S 中误差测度减少最多的分割, 即划分准则为 $\Delta\varepsilon(s^*, t) = \max_{s \in S} \Delta\varepsilon(s, t)$ 。循环分割 $N-1$ 次, 从而生成一棵由 $N-1$ 个阈值分割构成的 N 叉树。

3.2 区间内的样本特征估计

在对应的划分区间, 可以用不同度量方法估计正负样本特征的置信度。文献[2]以最小化指数误差函数为目标来选择特征, 采用样本所在区间正负样本密度比的对数作弱分类器的度量。KLBoosting 方法以 Kullback-Leibler 离散度作为度量估计特征。由于弱分类器设计中特征和样本数量较多, 因此其计算量较大。而训练算法和检测算法不能耗时过多, 基于速度和性能的考虑, 本文采用 GAB 算法进行弱分类器的度量, 定义弱分类器为

$$f(x) = \sum_{j=1}^n [P_w(y=+1 | x \in \text{bin}_j) - P_w(y=-1 | x \in \text{bin}_j)]$$

其中, $P_w(y=+1 | x \in \text{bin}_j) = \frac{W_{+1}^j}{W_{+1}^j + W_{-1}^j}$; $P_w(y=-1 | x \in \text{bin}_j) = \frac{W_{-1}^j}{W_{+1}^j + W_{-1}^j}$; $W_l^j = P(f_{\text{Haar}}(x) \in \text{bin}_j, y=l)$, $l = \pm 1$, 即 W_l^j 表示第 j 个区间 bin_j 内正样本或负样本权值的总和。

图 1 给出了根据上述划分策略生成的阈值及其置信度。

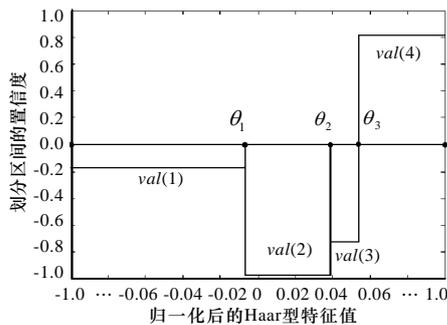


图 1 Haar 型特征值的 4 个区间及其置信度

在图 1 中, 当划分个数 $\text{split_num}=3$ 时, 生成了相应的阈值 $\theta_1, \theta_2, \theta_3$ 以及对应的区间置信度 $\text{val}(1) \sim \text{val}(4)$ 。

3.3 弱分类器的训练算法

弱分类器的训练算法描述如下:

(1) 基于误差测度减少最大化准则创建一棵 N 叉树 nodes , 将 Haar 型特征值域由 θ_j 分割生成 N 个节点, 每个节点对应一个划分区域及其基本分类器 $h_j(x - \theta_j)$, $j = 1, 2, \dots, N-1$ 。

$\text{nodes} = \text{weak_learner}(\text{tree_node}, x, y, w, \text{split_num})$

weak_learner 函数实现步骤如下:

1) 创建树的根节点, 具体描述如下:

$\text{create_root_node}(s)$;

$\text{learn_node_unit}(\text{cur_node}, x, y, w, \alpha_L, \alpha_R, \varepsilon(t_L), \varepsilon(t_R), \theta)$;

$\varepsilon(s) = \varepsilon(t_L) + \varepsilon(t_R)$;

2) 反复将一个节点分割生成新的左、右子节点, 具体描述如下:

for ($i = 1$; $i \leq \text{split_num}$; $i++$)

{ if ($s_L \neq 0$)

{ $s_L \leftarrow t$; $\text{create_left_node}(t)$;

$\text{learn_node_unit}(\text{cur_node}, x, y, w, \alpha_L, \alpha_R, \varepsilon(t_L), \varepsilon(t_R), \theta)$;

$\Delta\varepsilon(s, t) = \varepsilon(s) - \varepsilon(t_L) - \varepsilon(t_R)$;

if ($s_R \neq 0$)

{ $s_R \leftarrow t$; $\text{create_right_node}(t)$;

$\text{learn_node_unit}(\text{cur_node}, x, y, w, \alpha_L, \alpha_R, \varepsilon(t_L), \varepsilon(t_R), \theta)$;

$\Delta\varepsilon(s, t) = \varepsilon(s) - \varepsilon(t_L) - \varepsilon(t_R)$;

搜索使误差测度降低最多的节点 $\Delta\varepsilon(s^*, t) = \max_{s \in S} \Delta\varepsilon(s, t)$, 进入

下一轮分割; }

(2) 累加 N 个划分区间的置信度, 生成一个弱分类器 $f(x)$, 具体描述如下:

for ($j = 1$; $j \leq N$; $j++$)

{ $\text{basic_classifier} = \text{nodes}(j)$;

$\text{h_out} = \text{calculate_output}(\text{basic_classifier}, x)$;

$\text{weak_classifier} = \text{weak_classifier} + \text{h_out}$ }

上述子函数 learn_node_unit 的实现如下过程:

(1) 对所有样本的特征值 $\{x_1, x_2, \dots, x_m\}$ 进行由小到大排序。

(2) 以最小化加权平方和误差为准则, 对已排序的特征值从头到尾扫描一遍, 查找最佳分界点, 求得阈值 θ_j 。

1) 以 $k = 1, 2, \dots, m$ 为界将特征值划分为左子区间 x_1, x_2, \dots, x_k 和右子区间 $x_{k+1}, x_{k+2}, \dots, x_m$, 计算左右子区间的度量 α_L, α_R :

$$\alpha_L = \frac{W_{+1}^L - W_{-1}^L}{W_{+1}^L + W_{-1}^L}$$

$$\alpha_R = \frac{W_{+1}^R - W_{-1}^R}{W_{+1}^R + W_{-1}^R}$$

2) 对所有特征用加权最小平方和计算 α_L, α_R 与实际值 y_i 的误差 $\{\varepsilon_1, \varepsilon_2, \dots, \varepsilon_m\}$:

$$\varepsilon_k(t) = \varepsilon_k(t_L) + \varepsilon_k(t_R), \quad k = 1, 2, \dots, m$$

其中, $\varepsilon_k(t_L) = \sum_{i=1}^k w_i (y_i - \alpha_L)^2$; $\varepsilon_k(t_R) = \sum_{i=k+1}^m w_i (y_i - \alpha_R)^2$ 。

3) 搜索出最小误差对应的特征值 $\theta_j = \arg \min_{\theta} \{\varepsilon_1, \varepsilon_2, \dots, \varepsilon_m\}$ 。

(3) 生成一个基本分类器

$$h_j(x-\theta_j) = \begin{cases} \alpha_L & \text{if } x_i < \theta_j \\ \alpha_R & \text{else} \end{cases}$$

4 实验结果

本文通过如下实验比较单阈值弱分类器(DAB 方法)和多阈值弱分类器(GAB 方法)的性能。实验的训练数据集和验证数据集均采用 2 429 幅人脸图像作为正样本和 2 500 幅非人脸图像作为负样本, 图像大小为 20×20 像素, 选用的 Haar 型特征是文献[1]中的基本矩形特征。

图 2 和图 3 描述了最小 dr 为 99% 且最大 fpr 为 100% 时, 不同弱分类器的训练过程。

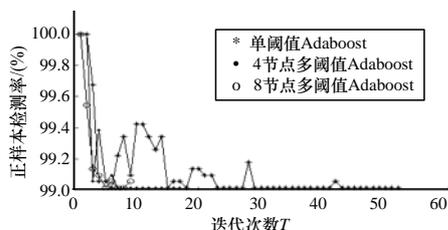


图 2 不同弱分类器的正样本检测率

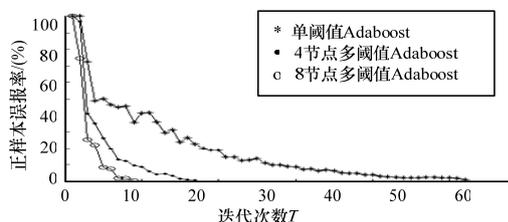


图 3 不同弱分类器的正样本误报率

当 $T=9$ 时, 单阈值方法与 4 节点多阈值方法、8 节点多阈值方法的 fpr 值分别为 35.6%、9.6%、0.9%。当 $T=17$ 时, 单阈值方法与 4 节点多阈值方法的 fpr 值分别为 22.7%、0.7%。

可见, 当 $T \leq 17$ 时, 多阈值方法的 fpr 较单阈值方法的性能可以提高 22% 以上, 多阈值方法收敛速度较快。

5 结束语

本文基于 RAB 思想对人脸检测器的弱分类器设计进行如下改进: (1) 基于误差测度减少最大化准则, 对特征值域用多个阈值进行不等距划分; (2) 采用 GAB 方法估计弱分类器的置信度。上述措施显著提高了单个 Haar 型特征值的检测性能, 提高了检测器的收敛速度, 使整个检测器训练可以在指定的最小 dr 和最大 fpr 下较快地生成, 从而快速生成相应的 Cascade 分类器。

参考文献

- [1] Viola P, Jones M. Rapid Object Detection Using a Boosted Cascade of Simple Features[C]//Proceedings of IEEE CVPR'01. Kauai, Hawaii, USA: IEEE Computer Society Press, 2001: 511-518.
- [2] Schapire R E, Singer Y. Improved Boosting Algorithms Using Confidence-rated Predictions[J]. Machine Learning, 1999, 37(3): 297-336.
- [3] 武 勃, 黄 畅, 艾海舟, 等. 基于连续 Adaboost 算法的多视角人脸检测[J]. 计算机研究与发展, 2005, 42(9): 1612-1621.
- [4] Wu Jianxin, Rehg J M, Mullin M D. Learning a Rare Event Detection Cascade by Direct Feature Selection[C]//Proc. of Annual Conference on Neural Information Processing Systems. Columbia, USA: [s. n.], 2004: 1523-1530.
- [5] Friedman J, Hastie T, Tibshirani R. Additive Logistic Regression: A Statistical View of Boosting[J]. Annals of Statistics, 2000, 28(2): 337-374.

编辑 陈 晖

(上接第 162 页)

为进一步说明算法的效果, 在 $p=20\%$ 和 $N=20$ 情况下将 ASERF 同一类似算法 DPIT^[5]做了比较, 如图 4 所示。可以看到, DPIT 算法无法有效处理模拟中的情况, 正常流量存活率仅维持在 20% 左右, 而笔者的算法却能在攻击较为严重时仍保持 80% 的存活率。

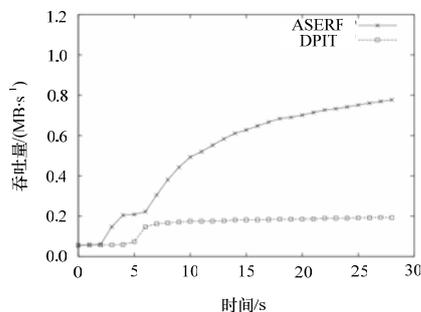


图 4 正常流量存活率对比

4 结束语

随着互联网的高速发展, DDoS 攻击日益成为一个安全隐患, 危害许多网站的安全, 而 DDoS 防御却一直是一个困难的问题。本文基于自治域内管理域的统一性, 提出了基于自治域内边界路由器反馈的分布式 DDoS 防御方法。该方法

具有一定的数学基础。经模拟实验验证, 该方法能够维持正常流量较高的存活率, 具不错的防御效果。

参考文献

- [1] Wan K K K, Chang R K C. Engineering of a Global Defense Infrastructure for DDoS Attacks[C]//Proc. of the 10th International Conference on Network. Los Alamitos, CA, USA: IEEE Computer Society Press, 2002: 419-427.
- [2] Mahajan R, Bellovin S M, Floyd S, et al. Controlling High Bandwidth Aggregates in the Network[J]. Computer Communication Review, 2002, 32(3): 75-85.
- [3] Keromytis A D, Misra V, Rubenstein D. SOS: Secure Overlay Services[J]. Computer Communication Review, 2002, 32(4): 117-129.
- [4] Yang Xiaowei, Wetherall D, Anderson T. A DoS Limiting Network Architecture[J]. Computer Communication Review, 2005, 35(4): 241-252.
- [5] Chen Shigang, Song Qingguo. Perimeter-based Defense Against High Bandwidth DDoS Attacks[J]. IEEE Transactions on Parallel and Distributed Systems, 2005, 16(6): 43-55.

编辑 顾逸斐

