

Unsupervised scene adaptation for faster multi-scale pedestrian detection

Speaker

Federico Bartoli¹

Giuseppe Lisanti¹, Svebor Karaman¹, Andrew D. Bagdanov² and Alberto Del Bimbo¹

¹ MICC (Media Integration and Communication Center) - University of Florence, Italy {firstname.lastname}@unifi.it

² CVC (Computer Vision Center) - Autonomous University of Barcelona, Spain bagdanov@cvc.uab.es



Real-time Pedestrian Detection

Application contexts

- Video Surveillance
- Interview Construction Const
- People Re-identification
- Action Recognition

Main critical factors

- Changes of scale and strong view-point dependency
 - Different target locations can produce high scale changes
 - Lost of scene depth information in the image
- Overal and the second secon
 - Different person poses (e.g. front or side view)
 - Changes in illumination intensity
- Scene complexity
 - Indoor or Outdoor
 - Clutter, crowd and partial occlusion





Standard execution pipeline of a multi-scale pedestrian detector





No Maximal Suppression

Federico Bartoli (Unifi::Micc)

Faster Multi-Scale Pedestrian Detection



Standard execution pipeline of a multi-scale pedestrian detector

Four principal phases

- Main bottlenecks:
 - Feature Extraction on Pyramid of Image Channel features [Dollar'14]
 - Oetection Windows Proposal: Sparse or Dense sampling Scene adapted detection windows proposal
 - Classification: Boosting, SVM Soft cascade approximation
 - On Maximal Suppression



No Maximal Suppression

Federico Bartoli (Unifi::Micc)

Faster Multi-Scale Pedestrian Detection





Faster multi-scale pedestrian detection

Question

How to increase the speed of a pre-trained pedestrian detector on a scene?

Framework Proposed

- Speed up the detection process of a Soft-Cascade pedestrian detector
- No a priori information about the scene required
- All learning done by mining statistics about the detector operating on the scene
- Exploit only ROS (Region of Support) information to build the models
- Strategies:
 - Linear Cascade Approximation: acts on classifier domain, for each sample estimate a final score without calculating all stages
 - Generative model for candidate window proposal: acts on pyramid domain, modelling the scene-dependent statistics of detection windows in terms of both location and scale
- The result is a significant reduction in the total number of stages evaluation required in the soft cascade detection process



p_A

Linear Cascade Approximation

Soft Cascade Architecture

Let $x \in \mathbb{R}^D$ be a sample to evaluate and $Y \in \{-1,1\}$ its class label:

- Classifier: $H(x) = \sum_{k=1}^{T} f_k(x)$, where $f_k : \mathbb{R}^D \longrightarrow \mathbb{R}$ is a stage computation
- Partial Score: $H_t(x) = \sum_{k=1}^t f_k(x)$ the sum of the first t stage scores
- x is classified positive $(Y = 1) \iff \Psi(H_t(x), \theta_t) \ge 0 \quad \forall t \in [1, T]$ where Ψ is a *stopping criterion* and $\{\theta_t\}$ are each stage rejection thresholds.

Linear Cascade Approximation

- Objective: For a given test sample x, we want to consider only a reduced number t < T of stages of H(x) in order to assign a score to a *detection window*
- Find $\tilde{H}_{t \to T} \in \mathbb{R}$ that estimates H using only the first t stages of the soft cascade, such that:

 $H(x) \simeq \tilde{H}_{t \to T}(x) \quad \forall x \in \mathcal{P}(I)$



Communication Center

Linear Cascade Approximation

• ex. Average positive traces extracted from a soft cascade of 1024 stages on the Oxford dataset. Traces are colored based on their level membership in the pyramid





<u>r</u>

Linear Cascade Approximation

Strategy

- Grouping all traces respect to their level
- Linear regression to estimate the parameters (slope and intercept) for the interpolation
- Compute the average trace for each group

Final score approximation takes the following form:

$$H_{t \to T}(x) = \bar{\mathbf{w}}_l \cdot \begin{bmatrix} 0 & T-t \end{bmatrix} + H_t(x) + \bar{\epsilon}_l$$

where

- l: level of x
- $\bar{\mathbf{w}}_l \equiv \mathbf{E}[\{\mathbf{w}_l^i\}]$ are the average trace parameters for the level l:

$$\begin{array}{l} \mathbf{w}_{l}^{i} = \arg\min_{\mathbf{w}} ||\mathbf{S}^{\mathsf{T}}\mathbf{w} - \mathbf{h}_{t \rightarrow T}(x^{(i)})|| \\ \mathbf{w} \in \mathbb{R}^{2}, \, \mathbf{w} = \begin{bmatrix} w_{0} & w_{1} \end{bmatrix} \text{ with } w_{0} \text{ the intercept and } w_{1} \text{ the slope} \\ \mathbf{S} = \begin{bmatrix} 1 & \cdots & 1 & \cdots & 1 \\ t & t + \Delta & t + 2\Delta & \cdots & T \end{bmatrix} \\ \mathbf{h} \, \mathbf{h}_{t \rightarrow T}^{\mathsf{T}}(x^{(i)}) = \begin{bmatrix} H_{t}(x^{(i)}) & H_{t + \Delta}(x^{(i)}) & \cdots & T \end{bmatrix} \\ \mathbf{\Delta} \text{: sampling step for the stages used in regression} \end{array}$$

• $\bar{\epsilon}_l$ = average interpolation error on the stage T



Observations:

- The presence and scale of targets is highly dependent on the geometry of the scene.
- Only detection windows in a limited scale range can be detected in a sub-region of frame
- The complete evaluation of all possible scales in all sub-regions of the image is wasteful

Idea

Only evaluate detection windows with a high likelihood to be a local maxima considering the geometric and scale statistics on the scene





- Leveraging Region of Support (ROS) information:
- The ROS is indicative of both the detector precision and the scene geometry:
 - The cardinality of each ROS is a good estimate of true positive: objects with a low rank are often false positive.
 - The location and scale of strongs can be considered to learn a model able to describe the geometry and perspective of the scene
- ROS information are discriminative and can be extracted at no additional cost during the non maximum suppression process.
- ex. Some strongs (and their ROS) from a soft cascade classifier on a frame from Oxford:





2 Scene Model (M_n)



$$M_n = (\mathcal{G}_n, \{\tilde{\mathcal{H}}_b^l\}, \{\mu_b^l, \Sigma_b^l\}, \{E_b\})$$

where:

- n: grid of n^2 blocks
- $1 \leq l \leq L$ pyramid levels
- $\tilde{\mathcal{H}}_{b}^{l}$: \mathcal{H}_{b}^{l} normalized over all levels l in block b

•
$$E_b = \frac{\sum_{l=1}^{L} \mathcal{H}_b^l}{\sum_{\tilde{b} \in \mathcal{G}_n} \sum_{l=1}^{L} \mathcal{H}_{\tilde{b}}^l}$$

Observations:

- Training of Model weakly-supervised
- Search differentiated according to the sub-region of frame (block)
- Generation of detection windows based on: spatial position $(\{\mu_{b,l}\}, \{\Sigma_{b,l}\})$, scale $(\{\mathcal{H}_b\})$ and energy $(\{E_b\})$
- No need of calibration

Federico Bartoli (Unifi::Micc)



Sandidate windows proposal at detection time

Algorithm

For each block b and scale l of Images Pyramid $\mathcal{P}(I)$:

- Compute the total number of detection windows to genereate: $N = \gamma |\mathcal{P}(\mathcal{I})|E_b \mathcal{H}_b^l$
- If not enough information $(\mathcal{H}_{b}^{l} < \tau) \Longrightarrow$ uniform extraction in the block region
- Else randomly sample from normal distribution $\mathcal{N}(\mu_{h}^{l}, \Sigma_{h}^{l})$ with covariance expansion:
 - Strategy round-based
 - For each round the covariance matrix is expanded by a factor (using \mathcal{X}^2_{lpha} distribution)
 - Iteration until the total number of obtained detection windows is approximately N
 - Reduction of duplicate samples

Parameter $\gamma \in [0, 1]$: Proportion of detection windows of a pyramid to be evaluated

- An estimate of the final speedup we want from the resulting detector
- $\bullet\,$ Tradeoff between between speed $(\gamma \rightarrow 0)$ and accuracy $(\gamma \rightarrow 1)$ of the detector



Test and Results

Baseline:

- \bullet 3 Soft cascade with 1024 stages, Images Pyramid from 3 octaves with 8 levels each
- Features used by stages:
 - ► HOG with 6 bin for orientation (0° 360°)
 - Gradient Histogram
 - Color Channels LUV

Dataset:

- $\bullet\,$ seq. Oxford: sampling 1 fps from video Oxford (3 min) and frame reshape at 640×480
- $\bullet\,$ seq. <code>PETS</code>: uniform extraction of 200 images from PETS (795 frames) and reshape at a 640×480

Speed of proposed Framework in terms of stages saving:

$$\delta = \frac{\sum_{\forall x \in \mathcal{P}} \left[H(x) \right]}{\sum_{\forall x \in \mathcal{X}} \mathbf{1}_{\{c=0\}} \left[H(x) \right] + \mathbf{1}_{\{c=1\}} \left[\tilde{H}_{t \to T}(x) \right]}$$



Performance with Linear Cascade Approximation

- Dataset considerati: seq. Oxford e PETS
- 7 values for t, uniform extracted from [64, 961]



- Results:
 - ▶ seq. Oxford: reductions between 1% 24% with maxima accuracy loss lower 5%
 - ▶ seq. *PETS*: same results of saving, but error not more than 4% with t > 257





Limited savings with Linear Cascade Approximation

Considering only stages reduction during detecteion windows evaluation is not enough to obtain high saving values:

- $|\mathcal{X}_P| \ll |\mathcal{X}_N|$ (two order of magnitude)
- $\bullet\,$ Same cost on evaluation of \mathcal{X}_P and \mathcal{X}_N
- $\bullet\,$ Maximum savings lower than 50% respect to full evaluation of Pyramid





Performance with Generative Model

- Dataset: seq. Oxford
- $\mathcal{G}_n \in \{2, 3, 4, 5, 6\}$
- $SpeedUp \in [4 \times, 8 \times, 16 \times, 32 \times, 64 \times]$



- Results:
 - \blacktriangleright With all configurations we obtain a savings greater then 50%
 - ▶ For grid size of 2 × 2, the minimum and maximum saving values is 65% (2.85×) and 95% (19.44×) respectively.

Federico Bartoli (Unifi::Micc)

Faster Multi-Scale Pedestrian Detection

Media Integration and Communication Center

Performance with our Framework





Federico Bartoli (Unifi::Micc)

Faster Multi-Scale Pedestrian Detection

28 August 2014 16/18



Performance Comparison: Our Framework vs Main Person Detectors

| | Oxford | | PETS | |
|------------------------|--------------------------------|-------------------|--------------------------------|-------------------|
| Detectors | $Miss\operatorname{-rate}(\%)$ | $Savings(\delta)$ | $Miss\operatorname{-rate}(\%)$ | $Savings(\delta)$ |
| DPM | 80 | - | 34 | - |
| ACF | 97 | $1 \times$ | 51 | $1 \times$ |
| Baseline | 99 | $1 \times$ | 9 | $1 \times$ |
| Linear Cascade App. | 98.1 | $1.17 \times$ | 10.4 | $1.28 \times$ |
| Candidate Windows Pro. | 98.9 | $11.09 \times$ | 7.4 | 3.37 	imes |
| With both | 98.5 | $12.73 \times$ | 11.4 | $4.19 \times$ |



Conclusions

- In the classifier and Pyramid domains can be applied approximation strategies for complexity reduction
- Maximum saving is obtained means sparse sampling of detection windows to evaluate by classifier
- The ROS information proves to be effective data for modelling both geometry and statistics for a scene

Our Framework

- The proposed strategies are weakly-supervised
- The great reduction of stages to evaluate allows to run the detector in real-time
- The framework implementation results very easy, also no need of dedicated hardware to run (ex. GPU)