

THE UNIVERSITY OF BRITISH COLUMBIA

CPSC 425: Computer Vision



Lecture 23: Object Detection (cont)

Menu for Today (March 28, 2019)

Topics:

- Deformable part models
- Object Proposals

Redings:

- **Today's** Lecture: Forsyth & Ponce (2nd ed.) 15.1, 15.2, 17.2
- Next Lecture: Deep Learning (N/A)

Reminders:

- Assignment 5: Scene Recognition with Bag of Words due April 4th
- Performance guidelines
- Office hours: 3-4pm today



- Grouping - Image Segmentation



Today's "fun" Example:

Audio-Visual Scene Analysis with Self-Supervised Multisensory Features

Andrew Owens Alexei A. Efros UC Berkeley



Today's "fun" Example: DensePose

Dense Pose: Dense Human Pose Estimation In The Wild



Riza Alp Güler * INRIA, CentraleSupélec Natalia NeverovaIasonas KokkinosFacebook Al ResearchFacebook Al Research

Riza Alp Güler was with Facebook Al Research during this work.

4

Today's "fun" Example: Pose Estimation



[Vondrak et al., CVPR 2008]

One common strategy to obtain a better classifier is to combine multiple classifiers.

A simple approach is to train an ensemble of independent classifiers, and average their predictions.

Boosting is another approach.

— Train an ensemble of classifiers sequentially.

 Bias subsequent classifiers to correctly predict training examples that previous classifiers got wrong.

- The final boosted classifier is a weighted combination of the individual classifiers.

1. Evaluate each rectangle filter on each example



Weak classifier $h_j(x) = \begin{cases} 1 & \text{if } f_j(x) > \theta_j \end{cases}$ threshold



Image Credit: Ioannis (Yannis) Gkioulekas (CMU)

2. Select best filter/threshold combination

a. Normalize the weights

b. For each feature, j

c. Choose the classifier, h_t with the lowest error \mathcal{E}

3. Reweight examples

$$W_{t+1,i} = W_{t,i} \beta_t^{1-|h_t(x_i)-y_i|} \qquad \beta_t = \frac{\varepsilon_t}{1-\varepsilon_t}$$

$$w_{t,i} \leftarrow \frac{w_{t,i}}{\sum_{j=1}^{n} w_{t,j}}$$

$$h_j(x) = \begin{cases} 1 & \text{if } f_j(x) > \theta_j \\ 0 & \text{otherwise} \end{cases}$$

$$\mathcal{E}_j = \sum_i w_i \left| h_j(x_i) - y_i \right|$$

Image Credit: Ioannis (Yannis) Gkioulekas (CMU)

Viola & Jones algorithm

4. The final strong classifier is

$$h(x) = \begin{cases} 1 & \sum_{t=1}^{T} \alpha_t h_t(x) \ge \frac{1}{2} \sum_{t=1}^{T} \alpha_t \\ 0 & \text{otherwise} \end{cases} \quad \alpha_t = \log \frac{1}{\beta_t}$$

The final hypothesis is a weighted linear combination of the T hypotheses where the weights are inversely proportional to the training errors

Image Credit: Ioannis (Yannis) Gkioulekas (CMU)



Pedestrian Detection

pedestrians tend to take characteristic poses, (e.g. standing, walking)



Image window; Visualisation of HOG features; HOG features weighted by positive weights; HOG features weighted by negative weights

Fig. 17.7 in Forsyth & Ponce (2nd ed). Original source: Dalal and Triggs, 2005.

The sliding window approach applies naturally to pedestrian detection because



Sliding window detectors tend to fail when the object is not well described by a rigid template



Many complex objects are better represented using a parts model

Felzenszwalb et al., 2010



A deformable part model consists of a root and a set of parts

- **Root**: an approximate model that gives the overall location of the object

at somewhat different locations on the root for different instances

Felzenszwalb et al., 2010

- Parts: object components that have reliable appearance but might appear



Finding a window that looks a lot like the part close to that part's natural location relative to the root yields evidence that the object is present

Felzenszwalb et al., 2010

Each part has an appearance model and a natural location relative to the root







A parts model for a bicycle, containing a root and 6 parts

Figure source: Felzenszwalb et al., 2010









The learned root model is a set of linear weights $\beta^{(r)}$ applied to the feature descriptor of the root window

The i-th learned part model consists of

- a natural location (offset) relative to the root $\mathbf{v}^{(p_i)} = (u^{(p_i)}, v^{(p_i)})$
- a set of distance weights $\mathbf{d}^{(p_i)} = (d_1^{(p_i)}, d_2^{(p_i)}, d_3^{(p_i)}, d_4^{(p_i)})$



Figure source: Felzenszwalb et al., 2010

- a set of linear weights $\beta^{(p_i)}$ applied to the feature descriptor of the part window



Sliding Window with Deformable Part Model

The overall score of the deformable p the sum of several scores

- A root score compares the root to the window
 Each part has its own score, consisting of an appearance score and a
- location score

Model score = Root

The overall score of the deformable parts model at a particular window will be

t score +
$$\sum_{i}$$
 Part i score

Sliding Window with Deformable Part Model

to the root.

offset relative to the root.

The score for part i at offset (x, y) is given by: $S^{(p_i)}(x, y; \beta^{(p_i)}, \mathbf{d}^{(p_i)}, \mathbf{v}^{(p_i)}) = \beta^{(p_i)}\phi(x, y)$

Denote by $\phi(x, y)$ the feature descriptor of a part window at offset (x, y) relative

Denote by $(dx, dy) = (u^{(p_i)}, v^{(p_i)}) - (x, y)$ the difference from the part's natural

 $-\left(d_1^{(p_i)}dx + d_2^{(p_i)}dy + d_3^{(p_i)}(dx)^2 + d_4^{(p_i)}(dy)^2\right)$



Sliding Window with Deformable Part Model

to the root.

offset relative to the root.

The score for part i at offset (x, y) is given by: $S^{(p_i)}(x,y;\beta^{(p_i)},\mathbf{d}^{(p_i)},\mathbf{v}^{(p_i)}) = \beta^{(p_i)}\phi(x,y)$ $-\left(d_1^{(p_i)}dx + d_2^{(p_i)}dy + d_3^{(p_i)}(dx)^2 + d_4^{(p_i)}(dy)^2\right)$

The final part i score is the best score found over all possible offsets (x, y)

Part *i* score = $max_{(x,y)}$

Denote by $\phi(x, y)$ the feature descriptor of a part window at offset (x, y) relative

Denote by $(dx, dy) = (u^{(p_i)}, v^{(p_i)}) - (x, y)$ the difference from the part's natural

$$S^{(p_i)}(x,y;\beta^{(p_i)},\mathbf{d}^{(p_i)},\mathbf{v}^{(p_i)})$$



Learning the model can be tricky. Why?

- Learning the model can be tricky. Why?
- A class model can consist of multiple component models representing different canonical views
- e.g. a front and lateral model of a bicycle
- example

We do not know which component model should respond to which training

- Learning the model can be tricky. Why?
- A class model can consist of multiple component models representing different canonical views
- e.g. a front and lateral model of a bicycle
- We do not know which component model should respond to which training example
- We also do not know the locations of the parts in the training examples

However, notice that if the component and the part locations for each training example are given (fixed), we can simply train a **linear SVM** as usual

- However, notice that if the component and the part locations for each training example are given (fixed), we can simply train a **linear SVM** as usual
- This observation leads to the following iterative strategy: — Assume components and part locations are given (fixed). Compute appearance and offset models.
- Assume appearance and offset models are given (fixed). Re-estimate components and part locations.

Deformable Part Models: Hard Negative Mining

- Even a small false positive rate becomes noticeable

remain computationally feasible

Hard negative mining: As we train the classifier, apply it to the negative examples (e.g. 'not a bicycle') and keep track of ones that get a strong round of training.

- Sliding window detectors must search over an immense number of windows
- As a result, we want to train on as many negative examples as possible, but

response (e.g. are mistakenly detected as bicycles). Include these in the next

Deformable Part Model: Examples

A learned car model



Figure source: Felzenszwalb et al., 2010



٦

)

Deformable Part Model: Examples

A learned cat model







Figure source: Felzenszwalb et al., 2010



Deformable Part Models are Convolutional Neural Networks

Ross Girshick¹ Forrest Iandola² Trevor Darrell² Jitendra Malik² ¹Microsoft Research ²UC Berkeley

rbg@microsoft.com {forresti,trevor,malik}@eecs.berkeley.edu

Abstract

Deformable part models (DPMs) and convolutional neural networks (CNNs) are two widely used tools for visual recognition. They are typically viewed as distinct approaches: DPMs are graphical models (Markov random fields), while CNNs are "black-box" non-linear classifiers. In this paper, we show that a DPM can be formulated as a CNN, thus providing a synthesis of the two ideas. Our construction involves unrolling the DPM inference algorithm and mapping each step to an equivalent CNN layer. From this perspective, it is natural to replace the standard image features used in DPMs with a learned feature extractor. We call the resulting model a DeepPyramid DPM and experimentally validate it on PASCAL VOC object detection. We find that DeepPyramid DPMs significantly outperform DPMs based on histograms of oriented gradients features (HOG) and slightly outperforms a comparable version of the recently introduced R-CNN detection system, while running significantly faster.

CNN. In other words, deformable part models *are* convolutional neural networks. Our construction relies on a new network layer, *distance transform pooling*, which generalizes max pooling.

DPMs typically operate on a scale-space pyramid of gradient orientation feature maps (HOG [5]). But we now know that for object detection this feature representation is suboptimal compared to features computed by deep convolutional networks [17]. As a second innovation, we replace HOG with features learned by a fully-convolutional network. This "front-end" network generates a pyramid of deep features, analogous to a HOG feature pyramid. We call the full model a *DeepPyramid DPM*.

We experimentally validate DeepPyramid DPMs by measuring object detection performance on PASCAL VOC [9]. Since traditional DPMs have been tuned for HOG features over many years, we first analyze the differences between HOG feature pyramids and deep feature pyramids. We then select a good model structure and train a Deep-Pyramid DPM that significantly outperforms the best HOGbased DPMs. While we don't expect our approach to outperform a fine-tuned R-CNN detector [17] we do find that it

Recall: Sliding Window

Train an image classifier as described previously. 'Slide' a fixed-sized detection window across the image and evaluate the classifier on each window.



Image credit: KITTI Vision Benchmark

Recall: Sliding Window

Train an image classifier as described previously. 'Slide' a fixed-sized window.



looking for.

detection window across the image and evaluate the classifier on each

Image credit: KITTI Vision Benchmark

This is a lot of possible windows! And most will not contain the object we are

- object-like properties
- background texture

The object detector then considers these candidate regions only, instead of exhaustive sliding window search

Object proposal algorithms generate a short list of regions that have generic

- These regions are likely to contain some kind of foreground object instead of

First introduced by Alexe et al., who asked 'what is an object?' and defined an 'objectness' score based on several visual cues





First introduced by Alexe et al., who asked 'what is an object?' and defined an 'objectness' score based on several visual cues



This work argued that objects typically - are unique within the image and stand out as salient have a contrasting appearance from surroundings and/or - have a well-defined closed boundary in space



Multiscale Saliency

- Favors regions with a unique appearance within the image







High scale

Low scale

Successful Case

Failure Case





Colour Contrast

- Favors regions with a contrasting colour appearance from immediate surroundings



Successful Cases

Failure Case





Superpixels Straddling

- Favors regions with a well-defined closed boundary
- contain pixels both inside and outside of the window



(b)



— Measures the extent to which superpixels (obtained by image segmentation)

(c)





Superpixels Straddling

- Favors regions with a well-defined closed boundary
- contain pixels both inside and outside of the window







Successful Cases Failure Case

— Measures the extent to which superpixels (obtained by image segmentation)

(c)

(b)





TABLE 2: For each detector [11, 18, 33] we report its performance (left column) and that of our algorithm 1 using the same window scoring function (right column). We show the average number of windows evaluated per image #win and the detection performance as the mean average precision (mAP) over all 20 classes.

	[11] O	BJ- [11]	[18] C	BJ- [18]	ESS-BOW[33]	OBJ-BOW
mAP	0.186	0.162	0.268	0.225	0.127	0.125
#win	79945	1349	18562 -	1358	183501	

Speeding up [11] HOG pedestrian detector [18] Deformable part model detector [33] Bag of words detector

 Table credit: Alexe et al., 2012

.

Summary

Detection scores in the deformable part model are based on both appearance and location

The deformable part model is trained iteratively by alternating the steps 1. Assume components and part locations given; compute appearance and

- 1. Assume components and part lo offset models
- 2. Assume appearance and offset part locations

An object **proposal** algorithm generates a short list of regions with generic object-like properties that can be evaluated by an object detector in place of an exhaustive sliding window search

2. Assume appearance and offset models given; compute components and