Laurens van der Maaten



Introduction

- Object detection aims to find a particular object in an image
- Most popular object detectors are based on a *discriminative model*:
 - Gather annotated image patches (positive and negative examples)
 - Extract your favorite *image features* from these image patches
 - Train a *classifier* on the features to discriminate object from everything else
 - Classifier is applied on *candidate locations* to determine object presence
- The Dalal-Triggs detector is a commonly used object detector

• Extract *histograms of oriented gradients* (HOG) features from image patch:



• HOG features divide an image into small (8x8) *blocks*, and measure the *gradient orientations* in each of the blocks using a histogram (almost like SIFT)

* Dalal & Triggs, 2005

• Different objects have different HOG features:





• Train a linear SVM on annotated images to predict object presence:

Training:
$$\mathbf{w}^* = \operatorname*{argmin}_{\mathbf{w}} \max \left(0, 1 - y \mathbf{w}^T \phi(\mathbf{I}; \mathbf{x}) \right)$$

Detection: $s(\mathbf{I}; \mathbf{x}) = \mathbf{w}^{*T} \phi(\mathbf{I}; \mathbf{x})$



• Train a linear SVM on annotated images to predict object presence:

Training:
$$\mathbf{w}^* = \operatorname*{argmin}_{\mathbf{w}} \max \left(0, 1 - y \mathbf{w}^T \phi(\mathbf{I}; \mathbf{x}) \right)$$

Detection: $s(\mathbf{I}; \mathbf{x}) = \mathbf{w}^{*T} \phi(\mathbf{I}; \mathbf{x})$



• How do we get the *negative examples* to train the SVM?

• Train a linear SVM on annotated images to predict object presence:

Training:
$$\mathbf{w}^* = \operatorname*{argmin}_{\mathbf{w}} \max \left(0, 1 - y \mathbf{w}^T \phi(\mathbf{I}; \mathbf{x}) \right)$$

Detection: $s(\mathbf{I}; \mathbf{x}) = \mathbf{w}^{*T} \phi(\mathbf{I}; \mathbf{x})$



• How do we get the *negative examples* to train the SVM? Random patches!



• HOG visualization of the SVM weights for a pedestrian detector:





• Example of pedestrian detections using Dalal-Triggs detector:



• What can we do when a part of the object to be detected is occluded?

- What can we do when a part of the object to be detected is occluded?
- Exploit the fact that other parts of the object are still visible!

- What can we do when a part of the object to be detected is occluded?
- Exploit the fact that other parts of the object are still visible!
- Pictorial structures does this by modeling objects as a constellation of parts:



* Fischler & Elschlager, 1973

• Defines a score function that involves parts and part deformations:

 $s(\mathbf{I}; x_0, y_0, \dots, x_{|V|}, y_{|V|}) = \mathbf{w}_0^{\mathrm{T}} \phi(\mathbf{I}; x_0, y_0)$



Global object model

• Defines a score function that involves parts and part deformations:



Global object model

Object part models

* Felzenszwalb et al., 2010

• Defines a score function that involves parts and part deformations:



• Defines a score function that involves parts and part deformations:



Global object model

Object part models

Deformation model

• Deformable template models are much more robust against *partial occlusions* and *deformations* of non-rigid objects

• Find the *optimal* configuration of a pictorial structures (detection) as follows:

$$\max_{x_0, y_0, \dots, x_{|V|}, y_{|V|}} s(\mathbf{I}; x_0, y_0, \dots, x_{|V|}, y_{|V|})$$

• Find the *optimal* configuration of a pictorial structures (detection) as follows:

$$\max_{x_0, y_0, \dots, x_{|V|}, y_{|V|}} s(\mathbf{I}; x_0, y_0, \dots, x_{|V|}, y_{|V|})$$

$$g(x_i) = \min_{x_j} (f(x_j) + (x_i - x_j)^2)$$

• Find the *optimal* configuration of a pictorial structures (detection) as follows:

$$\max_{x_0, y_0, \dots, x_{|V|}, y_{|V|}} s(\mathbf{I}; x_0, y_0, \dots, x_{|V|}, y_{|V|})$$

$$g(x_i) = \min_{x_j} (f(x_j) + (x_i - x_j)^2)$$

final score
with deformations

• Find configuration of pict. structures model by maximizing over part locations:

$$\max_{x_0, y_0, \dots, x_{|V|}, y_{|V|}} s(\mathbf{I}; x_0, y_0, \dots, x_{|V|}, y_{|V|})$$

$$g(x_i) = \min_{\substack{x_j \\ x_j}} (f(x_j) + (x_i - x_j)^2)$$
final score negative part
with deformations model score

• Find the *optimal* configuration of a pictorial structures (detection) as follows:

$$\max_{x_0, y_0, \dots, x_{|V|}, y_{|V|}} s(\mathbf{I}; x_0, y_0, \dots, x_{|V|}, y_{|V|})$$



• Find the *optimal* configuration of a pictorial structures (detection) as follows:

$$\max_{x_0, y_0, \dots, x_{|V|}, y_{|V|}} s(\mathbf{I}; x_0, y_0, \dots, x_{|V|}, y_{|V|})$$

• For squared-error deformation models, this can be done very efficiently:



• Hence, we have a parabola for every pixel x_j rooted at $(x_j, f(x_j))$





• It is straightforward to compute the *intersection* between two parabolas:

$$\frac{(f(x_i) + x_i^2) - (f(x_j) + x_j^2)}{2x_i - 2x_j}$$

• If $x_j < x_i$: parabola corresponding to x_j is *below* that of x_i *left* of the intersection, and above it right of the intersection



- Maintain the lower envelope of the parabolas (parabolas and intersections)
- When adding a new parabola, there are two possibilities:



new intersection right of last intersection: maintain last parabola in the envelope new intersection left of last intersection: remove last parabola from the envelope

- This suggests a simple algorithm that is *linear* in the number of pixels:
 - Maintain list with the lower envelope of the parabolas (indices and intersections)
 - Move from *left to right* through all parabolas; and do for each parabola:
 - Find intersection of parabola with the last parabola in lower envelope
 - If intersection is left of last intersection in lower envelope: remove last parabola from lower envelope, and go back one step
 - Add parabola to lower envelope, starting from intersection



Graph structure

• One can define different graph structures, as long as they are trees:



• The tree structure is fixed, but edge lengths and directions are learned

• Examples of object detections by pictorial-structures models:





Example detections













• Use pictorial structures to prevent trackers from "switching" objects:



Questions?