CS 2770: Computer Vision Object Detection

Prof. Adriana Kovashka University of Pittsburgh February 26, 2019

So far: Image Classification



Other Computer Vision Tasks













Plan for this lecture

- Fully supervised detection
 - Pre-CNN: Deformable part models
 - Detection with region proposals: R-CNN, Fast/er R-CNN
 - Detection without region proposals: YOLO
 - Semantic and instance segmentation: FCN, Mask R-CNN
- Weak or out-of-domain supervision
 - Weakly supervised object detection
 - Domain adaptation

Object Detection



Object detection: basic framework

- Build/train object model
- Generate candidate regions in new image
- Score the candidates

Window-template-based models Building an object model

Given the representation, train a binary classifier



Window-template-based models Generating and scoring candidates



Window-template-based object detection: recap

Training:

- 1. Obtain training data
- 2. Define features
- 3. Define classifier

Given new image:

- 1. Slide window
- 2. Score by classifier





Dalal-Triggs pedestrian detector



- 1. Extract fixed-sized (64x128 pixel) window at multiple positions and scales
- 2. Compute HOG (histogram of gradient) features within each window
- 3. Score the window with a linear SVM classifier
- 4. Perform non-maxima suppression to remove overlapping detections with lower scores

Histograms of oriented gradients (HOG)

Divide image into 8x8 regions

Orientation: 9 bins (for unsigned angles)



Histograms in 8x8 pixel cells



Votes weighted by magnitude



Adapted from Pete Barnum

Train SVM for pedestrian detection using HoG



Adapted from Pete Barnum

Navneet Dalal and Bill Triggs, Histograms of Oriented Gradients for Human Detection, CVPR05



Adapted from Derek Hoiem

Are window templates enough?

• Many objects are articulated, or have parts that can vary in configuration

Images from Caltech-256, D. Ramanan



• Many object categories look very different from different viewpoints, or from instance to instance





Adapted from N. Snavely, D. Tran

Parts-based Models

Define object by collection of parts modeled by

- 1. Appearance
- 2. Spatial configuration





How to model spatial relations?

• One extreme: fixed template



Fixed part-based template

 Object model = sum of scores of features at fixed positions



How to model spatial relations?

• Another extreme: bag of words



How to model spatial relations?

• Star-shaped model



Parts-based Models

- Articulated parts model
 - Object is configuration of parts
 - Each part is detectable and can move around





Discriminative part-based models





P. Felzenszwalb, R. Girshick, D. McAllester, D. Ramanan, <u>Object Detection</u> with Discriminatively Trained Part Based Models, PAMI 32(9), 2010

Lana Lazebnik

Discriminative part-based models

Multiple components



P. Felzenszwalb, R. Girshick, D. McAllester, D. Ramanan, <u>Object Detection</u> with Discriminatively Trained Part Based Models, PAMI 32(9), 2010

Lana Lazebnik

Scoring an object hypothesis

 The score of a hypothesis is the sum of appearance scores minus the sum of deformation costs

 $z=(p_0,...,p_n)$

 p_0 : location of root $p_1,..., p_n$: location of parts





Felzenszwalb et al.

Detection



Felzenszwalb et al.

Car model

Component 1







Component 2







Car detections

high scoring true positives









high scoring false positives





Person model





Person detections

high scoring true positives







high scoring false positives (not enough overlap)





Cat model



Cat detections

high scoring true positives



high scoring false positives (not enough overlap)






Plan for this lecture

- Fully supervised detection
 - Pre-CNN: Deformable part models
 - Detection with region proposals: R-CNN, Fast/er R-CNN
 - Detection without region proposals: YOLO
 - Semantic and instance segmentation: FCN, Mask R-CNN
- Weak or out-of-domain supervision
 - Weakly supervised object detection
 - Domain adaptation

Complexity and the plateau

[Source: http://pascallin.ecs.soton.ac.uk/challenges/VOC/voc20{07,08,09,10,11,12}/results/index.html]



Girshick et al., "Rich Feature Hierarchies for Accurate Object Detection and Semantic Segmentation", CVPR 2014

Impact of Deep Learning



Object Detection as Regression?





CAT: (x, y, w, h)



DOG: (x, y, w, h) DOG: (x, y, w, h) CAT: (x, y, w, h)

DUCK: (x, y, w, h) DUCK: (x, y, w, h)

Object Detection as Regression?





CAT: (x, y, w, h) 4 numbers

 3
 122
 128
 2048

 4
 122
 132
 133
 dense

 1
 1
 132
 133
 dense
 dense

 1
 1
 132
 133
 dense
 dense
 dense

 1
 1
 132
 132
 133
 dense
 dense
 dense

 1
 1
 132
 132
 132
 133
 dense
 dense

DOG: (x, y, w, h) CAT: (x, y, w, h)

DOG: (x, y, w, h)

16 numbers

DUCK: (x, y, w, h) Many DUCK: (x, y, w, h) numbers!

Each image needs a different number of outputs!





Apply a CNN to many different crops of the image, CNN classifies each crop as object or background





Dog? NO Cat? NO Background? YES

Apply a CNN to many different crops of the image, CNN classifies each crop as object or background





Dog? YES Cat? NO Background? NO

Apply a CNN to many different crops of the image, CNN classifies each crop as object or background





Dog? YES Cat? NO Background? NO

Apply a CNN to many different crops of the image, CNN classifies each crop as object or background





Dog? NO Cat? YES Background? NO

Apply a CNN to many different crops of the image, CNN classifies each crop as object or background





Dog? NO Cat? YES Background? NO

Problem: Need to apply CNN to huge number of locations and scales, very computationally expensive!

Region Proposals

- Find "blobby" image regions that are likely to contain objects
- Relatively fast to run; e.g. Selective Search gives 1000 region proposals in a few seconds on CPU



Alexe et al, "Measuring the objectness of image windows", TPAMI 2012 Uijlings et al, "Selective Search for Object Recognition", IJCV 2013 Cheng et al, "BING: Binarized normed gradients for objectness estimation at 300fps", CVPR 2014 Zitnick and Dollar, "Edge boxes: Locating object proposals from edges", ECCV 2014

Speeding up detection: Restrict set of windows we pass through SVM to those w/ high "objectness"



Fig. 1: **Desired behavior of an objectness measure.** The desired objectness measure should score the blue windows, partially covering the objects, lower than the ground truth windows (green), and score even lower the red windows containing only stuff or small parts of objects.

Proposals cue: color contrast at boundary



(a) (b) (c) Fig. 3: **CC success and failure. Success:** the windows containing the objects (cyan) have high color contrast with their surrounding ring (yellow) in images (a) and (b). **Failure:** the color contrast for windows in cyan in image (c) is much lower.

Alexe et al., "Measuring the objectness of image windows", PAMI 2012 and CVPR 2010

Proposals cue: no segments "straddling" the object box



Fig. 5: The SS cue. Given the segmentation (b) of image (a), for a window w we compute $SS(w, \theta_{SS})$ (eq. 4). In (c), most of the surface of w_1 is covered by superpixels contained almost entirely inside it. Instead, all superpixels passing by w_2 continue largely outside it. Therefore, w_1 has a higher SS score than w_2 . The window w_3 has an even higher score as it fits the object tightly.

Proposals cue: many edges wholly contained inside box



Zitnick and Dollar, "Edge Boxes: Locating Object Proposals from Edges", ECCV 2014





Regions of Interest (RoI) from a proposal method (~2k)







Linear Regression for bounding box offsets



R-CNN: Regions with CNN features



image proposals (~2k / image) features (linear SVM)



Input Extract region image proposals (~2k / image)

Proposal-method agnostic, many choices

- Selective Search [van de Sande, Uijlings et al.] (Used in this work)
- Objectness [Alexe etal.]
- Category independent object proposals [Endres & Hoiem]
- CPMC [Carreira & Sminchisescu]









Dilate proposal













Step 4: Object proposal refinement



Linear regression

on CNN features



Original proposal

Predicted object bounding box

Bounding-box regression

R-CNN on ImageNet detection

ILSVRC2013 detection test set mAP



Girshick et al., "Rich Feature Hierarchies for Accurate Object Detection and Semantic Segmentation", CVPR 2014

Linear Regression for bounding box offsets



Post hoc component

What's wrong with slow R-CNN?

- Ad hoc training objectives
 - Fine-tune network with softmax classifier (log loss)
 - Train post-hoc linear SVMs (hingeloss)
 - Train post-hoc bounding-box regressions (least squares)
- Training is slow (84h), takes a lot of disk space
- Inference (detection) is slow
 - 47s / image with VGG16 [Simonyan & Zisserman, ICLR15]



Girshick, "Fast R-CNN", ICCV 2015

~2000 ConvNet forward passes per image

Fast R-CNN

- Fast test time
- One network, trained in one stage
- Higher mean average precision

Fast R-CNN



Girshick, "Fast R-CNN", ICCV 2015

Fast R-CNN



Girshick, "Fast R-CNN", ICCV 2015








Fast R-CNN (Training)



Fast R-CNN (Training)



Fast R-CNN vs R-CNN

	Fast R-CNN	R-CNN
Train time (h)	9.5	84
Speedup	8.8x	1x
Test time / image	0.32s	47.0s
Test speedup	146x	1x
mAP	66.9%	66.0%

Timings exclude object proposal time, which is equal for all methods. All methods use VGG16 from Simonyan and Zisserman.

Fast<u>er</u> R-CNN



Ren et al, "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks", NIPS 2015



Accurate object detection is slow!

	Pascal 2007 mAP	Speed	
DPM v5	33.7	.07 FPS	14 s/img
R-CNN	66.0	.05 FPS	20 s/img





Accurate object detection is slow!

	Pascal 2007 mAP	Speed	
DPM v5	33.7	.07 FPS	14 s/img
R-CNN	66.0	.05 FPS	20 s/img
Fast R-CNN	70.0	.5 FPS	2 s/img
Faster R-CNN	73.2	7 FPS	140 ms/img
YOLO	69.0	45 FPS	22 ms/img



Plan for this lecture

- Fully supervised detection
 - Pre-CNN: Deformable part models
 - Detection with region proposals: R-CNN, Fast/er R-CNN
 - Detection without region proposals: YOLO
 - Semantic and instance segmentation: FCN, Mask R-CNN
- Weak or out-of-domain supervision
 - Weakly supervised object detection
 - Domain adaptation

Detection without Proposals: YOLO / SSD



Input image 3 x H x W



Divide image into grid 7 x 7

Image a set of **base boxes** centered at each grid cell Here B = 3 Within each grid cell:

- Regress from each of the B base boxes to a final box with 5 numbers:
 - (dx, dy, dh, dw, confidence)
- Predict scores for each of C classes (including background as a class)

Output: 7 x 7 x (5 * B + C)

Redmon et al, "You Only Look Once: Unified, Real-Time Object Detection", CVPR 2016 Liu et al, "SSD: Single-Shot MultiBox Detector", ECCV 2016

Slide by: Justin Johnson

Split the image into a grid



Each cell predicts boxes and confidences: P(Object)



Each cell also predicts a probability P(Class | Object)



Combine the box and class predictions



Finally do NMS and threshold detections



This parameterization fixes the output

size

Each cell predicts:

- For each bounding box:
 - 4 coordinates (x, y, w, h)
 - 1 confidence value
- Some number of class probabilities

For Pascal VOC:

- 7x7 grid
- 2 bounding boxes / cell
- 20 classes





YOLO works across many natural images







It also generalizes well to new domains







YOLOv2: Fast, Accurate Detection

















Each node is a conditional probability



Each node is a conditional probability


















Plan for this lecture

- Fully supervised detection
 - Pre-CNN: Deformable part models
 - Detection with region proposals: R-CNN, Fast/er R-CNN
 - Detection without region proposals: YOLO
 - Semantic and instance segmentation: FCN, Mask R-CNN
- Weak or out-of-domain supervision
 - Weakly supervised object detection
 - Domain adaptation

Semantic Segmentation



Semantic Segmentation



Label each pixel in the image with a category label

Don't differentiate instances, only care about pixels

Semantic Segmentation Idea: Sliding Window



Farabet et al, "Learning Hierarchical Features for Scene Labeling," TPAMI 2013 Pinheiro and Collobert, "Recurrent Convolutional Neural Networks for Scene Labeling", ICML 2014

Semantic Segmentation Idea: Sliding Window



Problem: Very inefficient! Not reusing shared features between overlapping patches

Farabet et al, "Learning Hierarchical Features for Scene Labeling," TPAMI 2013 Pinheiro and Collobert, "Recurrent Convolutional Neural Networks for Scene Labeling", ICML 2014

Design a network as a bunch of convolutional layers to make predictions for pixels all at once!



Design a network as a bunch of convolutional layers to make predictions for pixels all at once!



Design network as a bunch of convolutional layers, with **downsampling** and **upsampling** inside the network!



Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR 2015 Noh et al, "Learning Deconvolution Network for Semantic Segmentation", ICCV 2015



Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR 2015 Noh et al, "Learning Deconvolution Network for Semantic Segmentation", ICCV 2015

In-Network upsampling: "Unpooling"





1	1	2	2
1	1	2	2
3	3	4	4

Input: 2 x 2

Output: 4 x 4



1	0	2	0
0	0	0	0
3	0	4	0
0	0	0	0

Input: 2 x 2

Output: 4 x 4

In-Network upsampling: "Max Unpooling"



Learnable Upsampling: Transpose Convolution

3 x 3 transpose convolution, stride 2 pad 1





Input: 2 x 2

Output: 4 x 4

Learnable Upsampling: Transpose Convolution

3 x 3 transpose convolution, stride 2 pad 1



Input: 2 x 2

Output: 4 x 4

Learnable Upsampling: Transpose Convolution



Transpose Convolution: 1D Example



Output

Output contains copies of the filter weighted by the input, summing at where at overlaps in the output

Adapted from Justin Johnson



Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR 2015 Noh et al, "Learning Deconvolution Network for Semantic Segmentation", ICCV 2015

Instance Segmentation



Mask R-CNN

He et al, "Mask R-CNN", ICCV 2017

What is Mask R-CNN: Parallel Heads





Slide by: Kaiming He

Mask R-CNN





Adapted from Justin Johnson

Plan for this lecture

- Fully supervised detection
 - Pre-CNN: Deformable part models
 - Detection with region proposals: R-CNN, Fast/er R-CNN
 - Detection without region proposals: YOLO
 - Semantic and instance segmentation: FCN, Mask R-CNN
- Weak or out-of-domain supervision
 - Weakly supervised object detection
 - Domain adaptation

What if no bounding boxes to train?

- Weakly supervised object detection
 - Image-level class labels
 - Image-level captions



- Let f_k(x, y) be the activation in the k-th map at location (x, y)
- Global average pooling: $F^k = \Sigma_{x,y} f_k(x, y)$
- Input to softmax is $S_c = \Sigma_k w_k^c F^k$ where w_k^c is the weight for class c and map k

$$S_c = \sum_k w_k^c \sum_{x,y} f_k(x,y) = \sum_{x,y} \sum_k w_k^c f_k(x,y)$$

k

• Map for class c: $M_c(x,y) = \sum w_k^c f_k(x,y)$



Class activation maps of top 5 predictions



Class activation maps for one object class

Table 3.	Localization	error on	the	ILSVRC	test	set	for	various
weakly-	and fully- sup	ervised n	neth	ods.				

Method	supervision	top-5 test error
GoogLeNet-GAP (heuristics)	weakly	37.1
GoogLeNet-GAP	weakly	42.9
Backprop [23]	weakly	46.4
GoogLeNet [25]	full	26.7
OverFeat [22]	full	29.9
AlexNet [25]	full	34.2



Figure 5. Class activation maps from CNN-GAPs and the class-specific saliency map from the backpropagation methods.



Figure 6. a) Examples of localization from GoogleNet-GAP. b) Comparison of the localization from GooleNet-GAP (upper two) and the backpropagation using AlexNet (lower two). The ground-truth boxes are in green and the predicted bounding boxes from the class activation map are in red.

Localization from captions



Ye et al., "Learning to discover and localize visual objects with open vocabulary", arxiv 2018

Localization from captions

• Learn via triplet loss

$$L(\theta) = \sum \left[Sim^{agr}(\boldsymbol{x}, \boldsymbol{t}') - Sim^{agr}(\boldsymbol{x}, \boldsymbol{t}) + \alpha \right]_{+}$$

- Aggregate similarity, all words and regions $Sim^{agr}(\boldsymbol{x}, \boldsymbol{t}) = \sum [S^{img}(\boldsymbol{x})S^{txt}(\boldsymbol{t})^T \odot Sim^{ind}(\boldsymbol{x}, \boldsymbol{t})]$
- Individual word/region similarity $Sim^{ind}(\boldsymbol{x}_i, t_j) = \frac{\langle G^{img}(\boldsymbol{f}_i), G^{txt}(t_j) \rangle}{\|G^{img}(\boldsymbol{f}_i)\|_2 \|G^{txt}(t_j)\|_2}$ Ye et al., "Learning to discover and localize visual objects with open vocabulary", arxiv 2018

Localization from captions



Ye et al., "Learning to discover and localize visual objects with open vocabulary", arxiv 2018

Localization from sound



Harwath et al., "Jointly Discovering Visual Objects and Spoken Words from Raw Sensory Input", ECCV 2018

Localization from sound



Fig. 4: Speech-prompted localization maps for several word/object pairs. From top to bottom and from left to right, the queries are instances of the spoken words "WOMAN," "BRIDGE,", "SKYLINE", "TRAIN", "CLOTHES" and "VEHICLES" extracted from each image's accompanying speech caption.

Harwath et al., "Jointly Discovering Visual Objects and Spoken Words from Raw Sensory Input", ECCV 2018

Localization from sound



Fig. 7: On the left are shown two images and their speech signals. Each color corresponds to one connected component derived from two matchmaps from a fully random MISA network. The masks on the right display the segments that correspond to each speech segment. We show the caption words obtained from the ASR transcriptions below the masks. Note that those words were never used for learning, only for analysis.

Detection from documentaries



Chen et al., "Discover and Learn New Objects from Documentaries", CVPR 2017
What if test data very diff from train?

• Adapt detection models

Adapting detectors



Figure 1. Illustration of different datasets for autonomous driving: From top to bottom-right, example images are taken from: *KITTI*[17], *Cityscapes*[5], *Foggy Cityscapes*[49], *SIM10K*[30]. Though all datasets cover urban scenes, images in those dataset vary in style, resolution, illumination, object size, *etc.* The visual difference between those datasets presents a challenge for applying an object detection model learned from one domain to another domain.

Adapting detectors



Figure 2. An overview of our Domain Adaptive Faster R-CNN model: we tackle the domain shift on two levels, the image level and the instance level. A domain classifier is built on each level, trained in an adversarial training manner. A consistency regularizer is incorporated within these two classifiers to learn a domain-invariant RPN for the Faster R-CNN model.

Adapting detectors

	img	ins	cons	car AP
Faster R-CNN				30.12
Ours	1			33.03
		~		35.79
	-	~		37.86
	 Image: A set of the set of the	~	✓	38.97

Table 1. The average precision (AP) of *Car* on the *Cityscapes* validation set. The models are trained using the *SIM 10k* dataset as the source domain and the *Cityscapes* training set as the target domain. *img* is short for *image-level alignment*, *ins* for *instance-level alignment* and *cons* is short for our *consistency loss*



Adapting classifiers



Thomas and Kovashka, "Artistic Object Recognition by Unsupervised Style Adaptation", ACCV 2018

Adapting classifiers



Fig. 2. Training with multiple modalities and style-invariance constraint. We train networks on real and synthetic data. We show an example of style transfer transforming photos into labeled synthetic cartoons. The style-invariance loss trains the FC2 layer to predict which modality the image came from. During backpropagation, we reverse its gradient before propagating it to the layers used by both classifiers. This encourages those layers to learn style-invariant features.

What's next?