Lecture 11: Object detection

1

Contains slides from S. Lazebnik, R. Girshick, B. Hariharan



Object detection with bounding boxes



What?

Where?

"Object detection"





- For each detection, determine whether it is a true or false positive



Evaluating an object detector

At test time, predict bounding boxes, class labels, and confidence scores

Intersection over union (IoU): Area(GT \cap Det) / Area(GT U Det) > 0.5





Evaluating an object detector

$IoU(A, B) = \frac{|A \cap B|}{|A \cup B|}$

Intersection over union (also known as Jaccard similarity)

Source: B. Hariharan

- (area under the curve)
- Take mean of AP over classes to get mAP

Evaluating an object detector

• For each class, plot Recall-Precision curve and compute Average Precision

Precision:

- true positive detections /
- total detections
- **Recall**:
- true positive detections / total positive test instances

Average precision

Precision

Recall

Source: B. Hariharan

Recall

Average precision

Source: B. Hariharan

Detection as classification

- Run through every possible box and classify
 - Well-localized object of class k or not?
- How many boxes?
 - Every pair of pixels = 1 box

•
$$\begin{pmatrix} N \\ 2 \end{pmatrix} = O(N^2)$$

- For 300×500 image, N = 150K
- 2.25×10^{10} boxes!
- Related challenge: almost all boxes are negative!

Selective search

Stage 1: generate candidate bounding boxes

Input image

Stage 2: apply classifier only to each candidate bounding box

Positive examples

Training Examples

Train

if overlap with positive 20-50%

Source: Torralba, Freeman, Isola

Edge detection

Bounding box proposal

[Zitnick and Dollar, "Edge Boxes...", 2014]

R-CNN: Region proposals + CNN features

R. Girshick, J. Donahue, T. Darrell, and J. Malik, Rich Feature Hierarchies for Accurate Object Detection and Semantic Segmentation, CVPR 2014.

Classify regions with linear classifier

Forward each region through ConvNet

Warped image regions

Region proposals from **selective search** (~2K rectangles that are likely to contain objects)

Input image

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

a. Crop

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

a. Crop

Compute CNN features

227 x 227

b. Scale (anisotropic)

ut Extract region ge proposals (~2k / image)

1. Crop b. Scale (anisotropic)

c. Forward propagate Output: "fc7" features

Extract region

Warped proposal

4096-dimensional fc7 feature vector

linear classifiers (SVM or softmax)

R-CNN at test time: proposal refinement

Original proposal

Linear regression

on CNN features

Predicted object bounding box

Bounding-box regression

Bounding-box regression

h

$(\Delta x \times w + x, \Delta y \times h + h)$

 $\Delta h \times h + h$

Source: R. Girshick

predicted

Non-maximum suppression

If two boxes overlap significantly (e.g. > 50% IoU), drop the one with the lower score. Usually use greedy algorithm.

Source: B. Hariharan

Problems with R-CNN

- 1. Slow! Have to run CNN per window
- 2. Hand-crafted mechanism for region proposal might be suboptimal.

"Fast" R-CNN: reuse features between proposals

Source: R. Girshick

R. Girshick, Fast R-CNN, ICCV 2015

- How do we crop from a feature map?
- Step 1: Resize boxes to account for subsampling

ROI Pooling

- How do we crop from a feature map?
- Step 2: Snap to feature map grid

ROI Pooling

- How do we crop from a feature map?
- Step 3: Overlay a new grid of fixed size

ROI Pooling

- How do we crop from a feature map?
- Step 4: Take max in each cell

See more here: <u>https://deepsense.ai/region-of-interest-pooling-explained/</u>

ROI Pooling

Source: B. Hariharan

S. Ren, K. He, R. Girshick, and J. Sun, Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks, NIPS 2015

RPN: Region Proposal Network

RPN: Region Proposal Network

3x3 "sliding window" -Scans the feature map looking for objects

RPN: Anchor Box

3x3 "sliding window" -Scans the feature map looking for objects

Anchor box: predictions are w.r.t. this box, *not the 3x3* sliding window

RPN: Anchor Box

3x3 "sliding window"-

 \succ Objectness classifier [0, 1]

Box regressor predicting (dx, dy, dh, dw)

Anchor box: predictions are w.r.t. this box, *not the 3x3* sliding window

RPN: Prediction (on object) Anchor box: transformed by box regressor P(object) = 0.94

- 3x3 "sliding window" > Objectness classifier [0, 1]
- Box regressor predicting (dx, dy, dh, dw)

RPN: Prediction (off object)

Objectness score

3x3 "sliding window" > Objectness classifier

Box regressor predicting (dx, dy, dh, dw)

Anchor box: transformed by box regressor

RPN: Multiple Anchors

32

Conv feature map

- 3x3 "sliding window" - \succ K objectness classifiers
- \succ K box regressors

Anchor boxes: *K* anchors per location with different scales and aspect ratios

Source: R. Girshick, K. He, S. Lazebnik

Faster R-CNN results

system	time	07 data	07+12 data
R-CNN	~50s	66.0	-
Fast R-CNN	~2s	66.9	70.0
Faster R-CNN	198ms	69.9	73.2

detection mAP on PASCAL VOC 2007, with VGG-16 pre-trained on ImageNet

Object detection progress

Source: S. Lazebnik

Is it possible do detection in one shot?

Conv feature map of the entire image

Single-stage object detector

• Divide the image into a coarse grid candidate boxes for each grid cell

J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, <u>You Only Look Once: Unified, Real-Time</u> <u>Object Detection</u>, CVPR 2016

Divide the image into a coarse grid and directly predict class label and a few

Source: S. Lazebnik

- **1.** Take conv feature maps at 7x7 resolution
- - For PASCAL, output is $7x7x30(30 = 20 + 2^{*}(4+1))$
- 7x speedup over Faster R-CNN (45-155 FPS vs. 7-18 FPS) but less accurate (e.g. 65% vs. 72 mAP%)

J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, You Only Look Once: Unified, Real-Time Object Detection, CVPR 2016

YOLO detector

2. Predict, at each location, a score for each class and 2 bboxes w/ confidences

Source: S. Lazebnik

Challenges in object detection

Beyond bounding boxes: instance segmentation

Predict segmentation mask for each object From COCO [Lin et al., 2014]

Source: B. Hariharan

Instance segmentation

ROI pooling with tiny change: bilinear interpolation instead of max

[He et al., "Mask R-CNN", 2017]

Image with training proposal

28x28 mask target

Image with training proposal

28x28 mask target

Image with training proposal

28x28 mask target

Image with training proposal

28x28 mask target

person1.00 person1.00 surfboard1.00 fboard1.00

urfboard.98 surfboard1.00

1000

12:22

person1.00

person1.00

dining table.95

wine glass1.00

bottle.97

Sur Itu

wine glass1.00

tv.98 tv.84 person1.00

person.88

wine glass.99

Human Pose

(Not shown: Head architecture is slightly different for keypoints)

> Add keypoint head (28x28x17)

> Predict one "mask" for each keypoint

 \gg Softmax over spatial locations (encodes one keypoint per mask "prior")

Panoptic Segmentation

Instance detection, "things"

Semantic segmentation "stuff"

panoptic segmentation [Kriilov et al. 2018] predict label + instance id per pixel

UPSNet-101-M: Cityscapes

2000

We still need lots of labeled examples

Handle the long tail of the distribution

Person, dog, table, ...

Frequency -

Object cate

Teacup, wreath, birdfeeder, ...

gories	\longrightarrow
$\mathbf{}$	54

Handle the "long tail" of the distribution

From COCO (80 categories) [Lin et al., 2014]

LVIS dataset (1000+ categories) "Few shot" (e.g. < 20 examples) [Gupta et al., 2019]

Image source: R. Girshick

Next time: Action recognition