Object Detection

School of Computer Science The University of Adelaide

S. Hamid Rezatofighi

Introduction

Object detection is a type of classification problem.

Object detection requires that we locate a specific object in an image, if it is there.

- It thus requires deciding whether an object is in an image, and if so determining its bounding box.
- This implies that the object can appear anywhere in the image (whereas for classification images tend to have one item only).

Introduction (cont.)

The challenge in object detection is the number of locations at which an object might appear in a single image.

- ► A megapixel image has 10⁶ locations at which an object can occur. If it can have 10 scales, and 10 rotations, then there are 10⁸ patches that need to be checked.
 - ► A false-positive rate of 1 in 10⁻⁶ will give 100 false positives per class per image.
- Special purpose detectors are constructed to find instances of the target object very rapidly given an input image.

Face detection





Without Face Detect AE



Without Face Detect FE



With Face Detect AE



With Face Detect FE

Pedestrian detection

Volvo demo video:

http://www.youtube.com/watch?v=rO86-ERWsgQ





"Indexian Distribution with 64 and isolate semicha of a used with integrated links the carly gales, a semicer fitted in these of the testor sear-less minor and a control control and. The solar's task is to detect algorithm in the observation for the detects in them. The canona determines what lines of algorithm is to its an energy ency duration, the determines are audited exercise to them. The canona determines what lines of algorithm is to its an energy ency duration, the determines are audited exercises to the data is durating gift in the verdication is had use clusters. It is not an encoded with the basics are pro-clusters, the data was an extension of the average and an audited is investment. Is it handing powers an automatical appendix.





Object Detection: Impact of Deep Learning





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





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



Warped image regions

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





Linear Regression for bounding box offsets



R-CNN: Problems

- Ad hoc training objectives
 - Fine-tune network with softmax classifier (log loss)
 - Train post-hoc linear SVMs (hinge loss)
 - 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]
 - Fixed by SPP-net [He et al. ECCV14]





Input image

Girshick, "Fast R-CNN", ICCV 2015. Figure copyright Ross Girshick, 2015; <u>source</u>. Reproduced with permission.













Girshick, "Fast R-CNN", ICCV 2015. Figure copyright Ross Girshick, 2015; <u>source</u>. Reproduced with permission.



Girshick, "Fast R-CNN", ICCV 2015. Figure copyright Ross Girshick, 2015; <u>source</u>. Reproduced with permission.

Faster R-CNN: Rol Pooling



512 x 18 x 8

(varies per proposal)

Girshick, "Fast R-CNN", ICCV 2015.

R-CNN vs SPP vs Fast R-CNN



Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014. He et al, "Spatial pyramid pooling in deep convolutional networks for visual recognition", ECCV 2014 Girshick, "Fast R-CNN", ICCV 2015

R-CNN vs SPP vs Fast R-CNN



Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014. He et al, "Spatial pyramid pooling in deep convolutional networks for visual recognition", ECCV 2014 Girshick, "Fast R-CNN", ICCV 2015

Faster R-CNN:

Make CNN do proposals!

Insert Region Proposal **Network (RPN)** to predict proposals from features

Jointly train with 4 losses:

- RPN classify object / not object 1.
- 2. **RPN** regress box coordinates
- 3. Final classification score (object classes)
- Final box coordinates 4.

Ren et al, "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks", NIPS 2015 Figure copyright 2015, Ross Girshick; reproduced with permission

loss



Fast<u>er</u> R-CNN: Make CNN do proposals!



Detection without Proposals: YOLO / SSD



Input image 3 x H x W

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



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)

Detection without Proposals: YOLO / SSD

Go from input image to tensor of scores with one big convolutional network!



Input image 3 x H x W

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

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)

Object Detection: Lots of variables ...

Base Network VGG16 ResNet-101 Inception V2 Inception V3 Inception ResNet MobileNet

Object Detection architecture Faster R-CNN R-FCN SSD

Takeaways Faster R-CNN is slower but more accurate

Image Size # Region Proposals SSD is much faster but not as accurate

Huang et al, "Speed/accuracy trade-offs for modern convolutional object detectors", CVPR 2017

R-FCN: Dai et al, "R-FCN: Object Detection via Region-based Fully Convolutional Networks", NIPS 2016 Inception-V2: loffe and Szegedy, "Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift", ICML 2015 Inception V3: Szegedy et al, "Rethinking the Inception Architecture for Computer Vision", arXiv 2016 Inception ResNet: Szegedy et al, "Inception-V4, Inception-ResNet and the Impact of Residual Connections on Learning", arXiv 2016 MobileNet: Howard et al, "Efficient Convolutional Neural Networks for Mobile Vision Applications", arXiv 2017

...

Aside: Object Detection + Captioning = Dense Captioning



Johnson, Karpathy, and Fei-Fei, "DenseCap: Fully Convolutional Localization Networks for Dense Captioning", CVPR 2016 Figure copyright IEEE, 2016. Reproduced for educational purposes.



C x 14 x 14

Mask R-CNN: Very Good Results!



He et al, "Mask R-CNN", arXiv 2017 Figures copyright Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick, 2017. Reproduced with permission.



C x 14 x 14