

Deep Learning 4

Object Detection



Michal Drozdal
mdrozdal@fb.com

About me



Index

Part 1
Introduction

Part 2
Basic blocks & concepts

Part 3
Models

Bonus material
Weakly supervised
localization

19/3/19

Deep learning 4 object detection

2

Computer Vision tasks

1. Classification



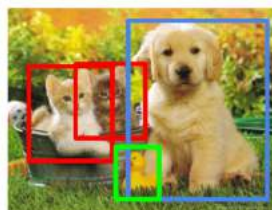
CAT

2. Localization:



CAT

3. Detection:



CAT, DOG, DUCK

4. Segmentation:



CAT, DOG, DUCK

19/3/19

Deep learning 4 object detection

<https://chaosmail.github.io/deeplearning/2016/10/22/intro-to-deep-learning-for-computer-vision/>

3

Datasets



19/3/19

Deep learning 4 object detection

4



2005-2012

Classification **Detection** **Segmentation**



- 20 categories
- 6k training images (17k objects)
- 6k validation + 10k test

26/2/18

Deep learning 4 object detection

5

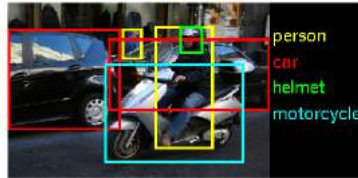


Localization



- 1000 categories
- 1.2M training images
- 150k validation + test images

Detection



- 200 categories
- 456k training images
- 60k validation + test images

Detection from video



- 30 categories
- 6k videos

19/3/19

Deep learning 4 object detection

6



Detection

Segmentation



- 80 categories
- 160k images
- 1M instances (350k people)

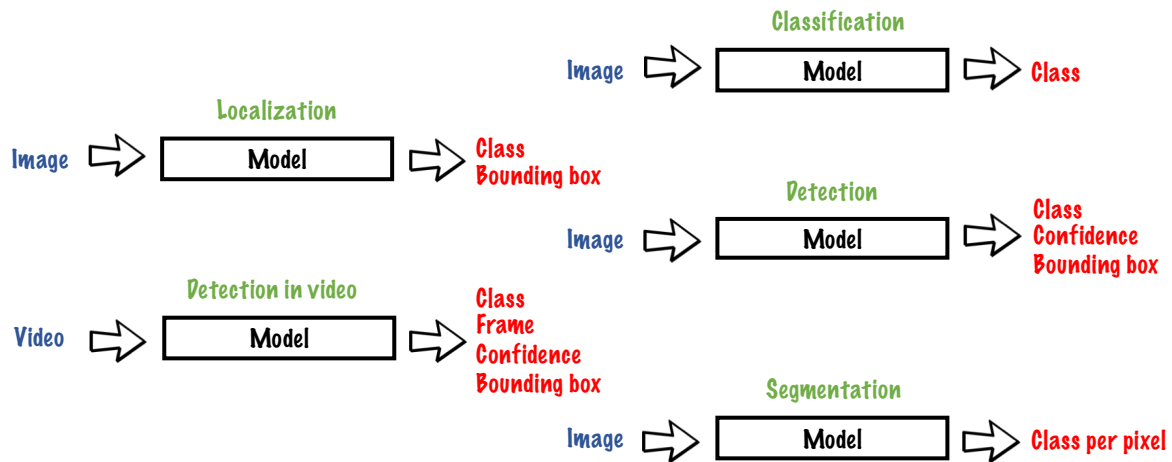
19/3/19

Deep learning 4 object detection

7

<http://presentations.cocodataset.org/COCO17-Detect-Overview.pdf>

Let's define an **input** and an **output** for each **task**

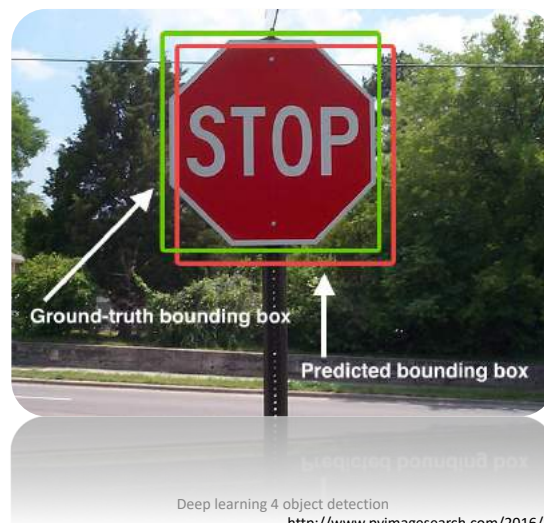


19/3/19

Deep learning 4 object detection

8

Evaluation



19/3/19

Deep learning 4 object detection

<http://www.pyimagesearch.com/2016/11/07/intersection-over-union-iou-for-object-detection/>

9

Intersection over Union (IoU)

$$\text{IoU} = \frac{\text{Intersection}}{\text{Union}} = \frac{\text{Diagram of two overlapping squares with the intersection shaded black}}{\text{Diagram of the union of two overlapping squares, shaded black}}$$

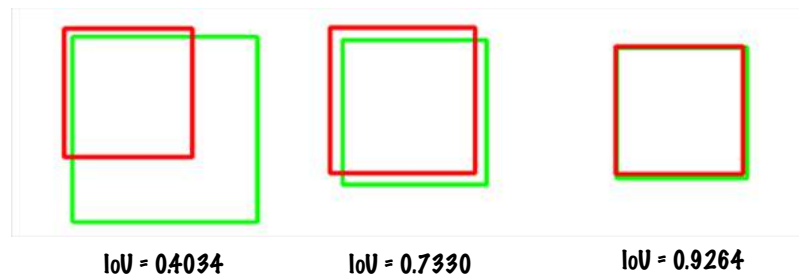
19/3/19

Deep learning 4 object detection

<http://www.pyimagesearch.com/2016/11/07/intersection-over-union-iou-for-object-detection/>

10

Intersection over Union (IoU)



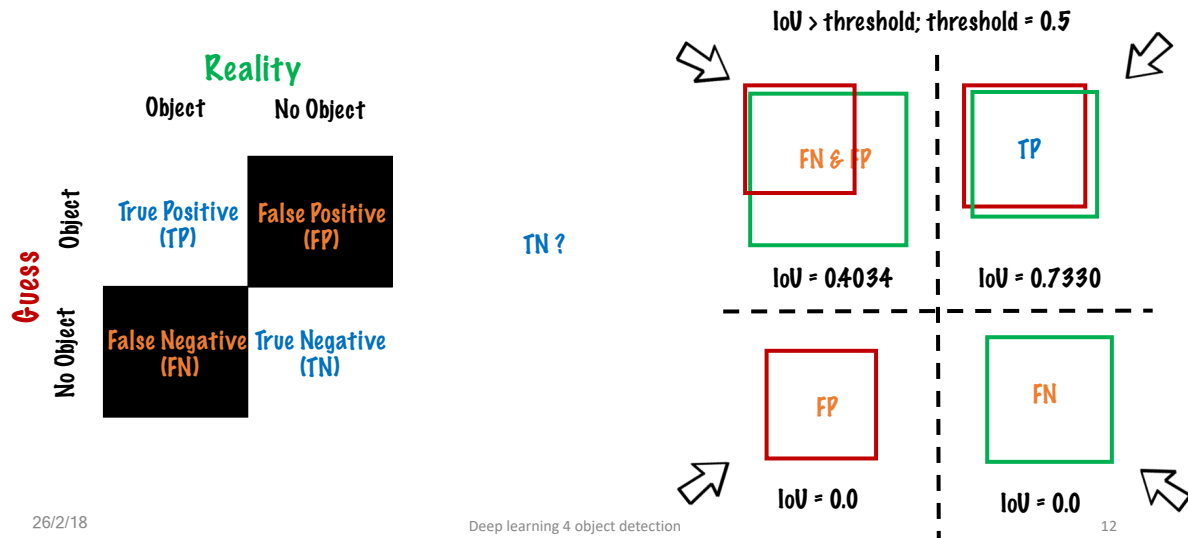
19/3/19

Deep learning 4 object detection

<http://www.pyimagesearch.com/2016/11/07/intersection-over-union-iou-for-object-detection/>

11

Being **correct**, being **wrong**



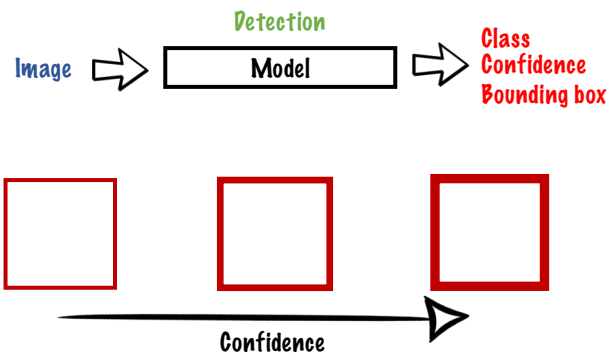
Localization

Localization error:
Wrong class or bad localization (IoU < 0.5)



http://image-net.org/challenges/talks/2016/ILSVRC2016_10_09_clsloc.pdf

Detection



Metrics:

- 1) Average Precision (AP)
- 2) Average Recall (AR)

19/3/19

Deep learning 4 object detection

14

Computing Average Precision (A four step procedure)

1. Order the predictions using confidence
2. Compute Precision and Recall
3. Plot Precision Recall plot (optional)
4. Compute Average Precision (AP) and Average Recall (AR)

19/3/19

Deep learning 4 object detection

15

1. Predictions ordering

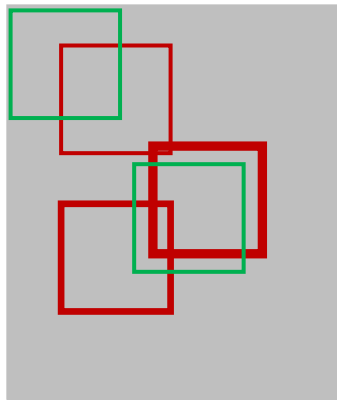
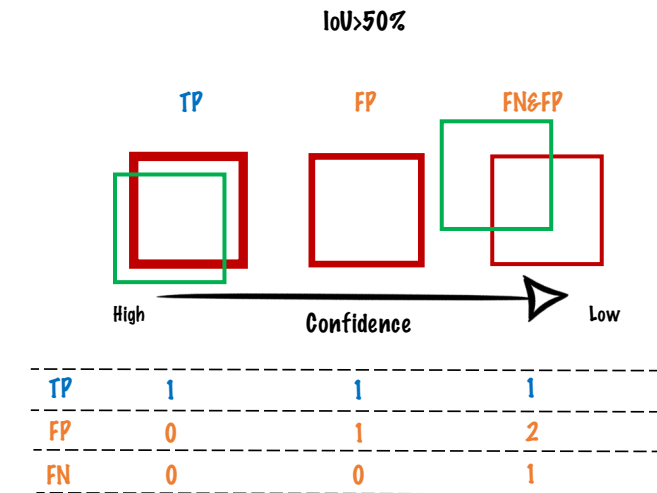


Image
Annotations
Predictions

19/3/19



Deep learning 4 object detection

16

1. Predictions ordering

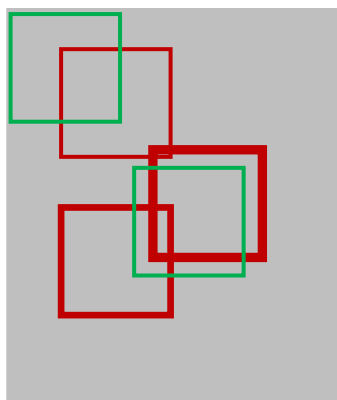
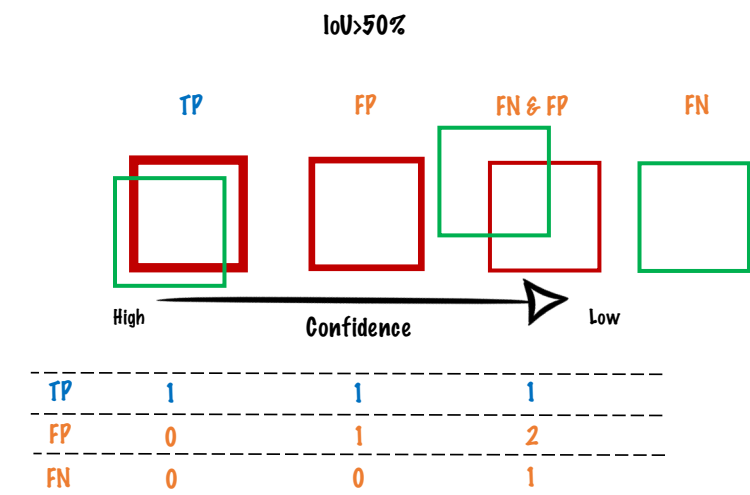


Image
Annotations
Predictions

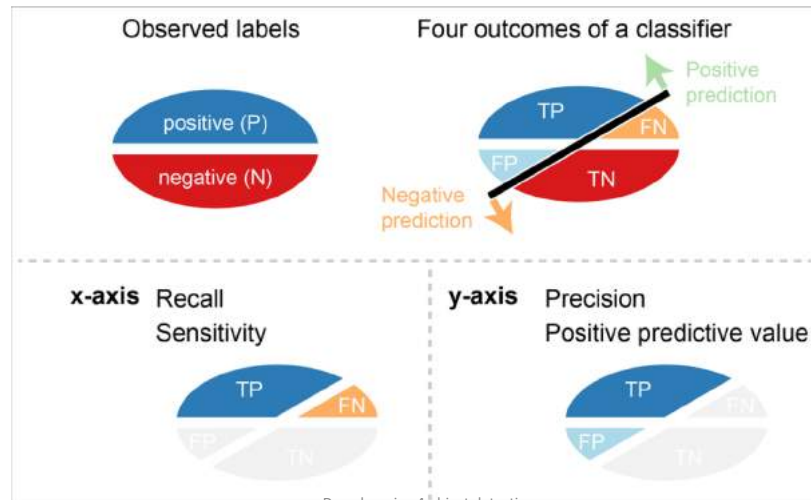
19/3/19



Deep learning 4 object detection

17

2. Precision & Recall



19/3/19

Deep learning 4 object detection

18

<https://classeeval.wordpress.com/introduction/introduction-to-the-precision-recall-plot/>

2. Precision & Recall

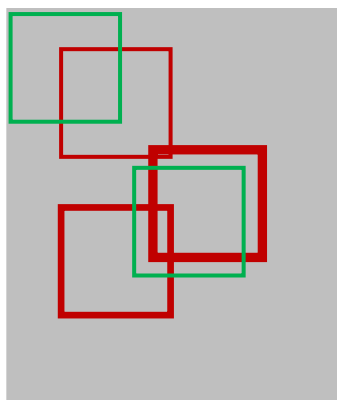
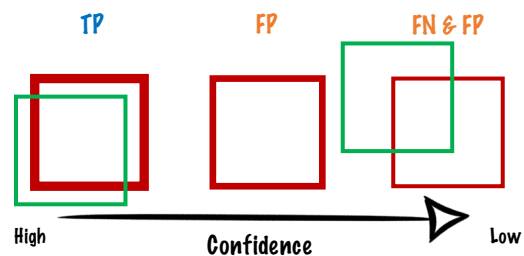


Image
Annotations
Predictions



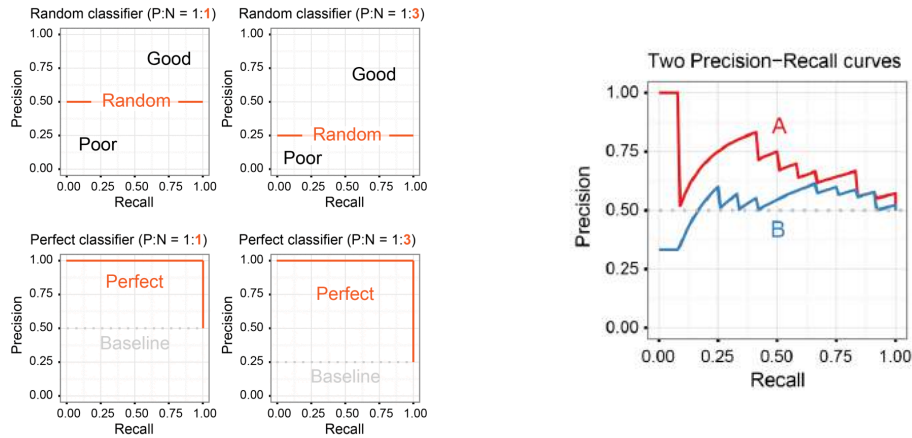
TP	1	1	1
FP	0	1	2
FN	0	0	1
Precision	1	.5	.3
Recall	.5	.5	1

19/3/19

Deep learning 4 object detection

19

3. Precision - Recall plots



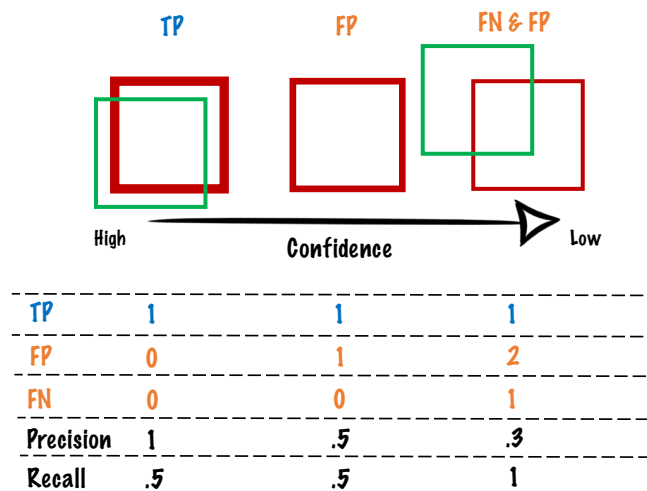
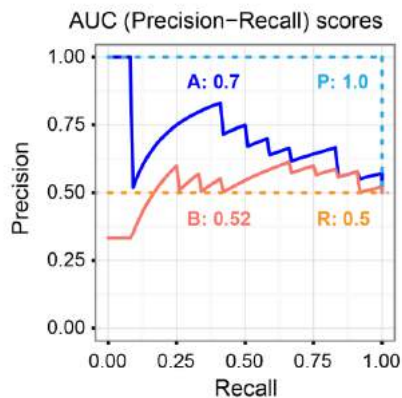
19/3/19

Deep learning 4 object detection

20

<https://classeeval.wordpress.com/introduction/introduction-to-the-precision-recall-plot/>

4. AP & AR


 $AP@(\text{IoU}=50\%) = ?$
 $R@(\text{pred}) = 0.5$

19/3/19

Deep learning 4 object detection

21

<https://classeeval.wordpress.com/introduction/introduction-to-the-precision-recall-plot/>



Average Precision (AP) : % AP at IoU=.50:.05:.95 (primary challenge metric)

AP % AP at IoU=.50 (PASCAL VOC metric)
 AP^{IoU=.75} % AP at IoU=.75 (strict metric)

AP Across Scales:

AP^{small} % AP for small objects: area < 32²
 AP^{medium} % AP for medium objects: 32² < area < 96²
 AP^{large} % AP for large objects: area > 96²

Average Recall (AR) :

AR^{max=1} % AR given 1 detection per image
 AR^{max=10} % AR given 10 detections per image
 AR^{max=100} % AR given 100 detections per image

AR Across Scales:

AR^{small} % AR for small objects: area < 32²
 AR^{medium} % AR for medium objects: 32² < area < 96²
 AR^{large} % AR for large objects: area > 96²

19/3/19

Deep learning 4 object detection

<http://mscoco.org/dataset/#detections-eval>

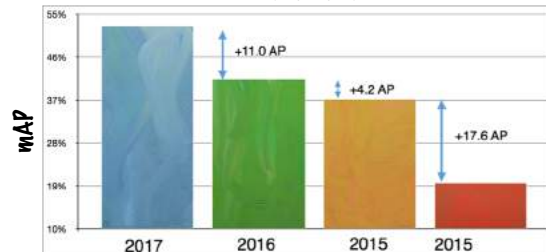

Localization



Detection



Detection



19/3/19

Deep learning 4 object detection

<http://presentations.cocodataset.org/COCO17-Detect-Overview.pdf>

Index

Part 1
Introduction

Part 2
Basic blocks & concepts

Part 3
Models

Bonus material
Weakly supervised
localization

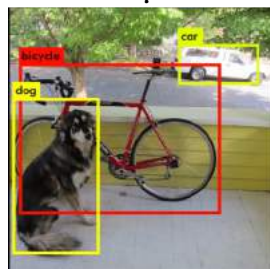
19/3/19

Deep learning 4 object detection

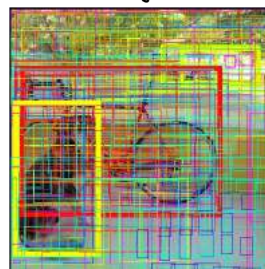
24

What is an outcome of an object detector?

We expect:



We get:



Lots of FP

19/3/19

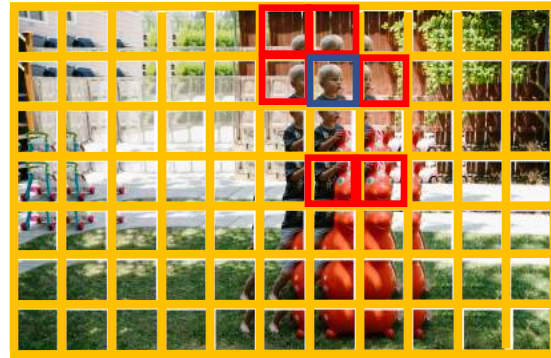
Deep learning 4 object detection

25

Image credit: <https://pjreddie.com/darknet/yolo/>

The unbalanced nature of detection. Hard Negative Mining

Build a head detector.



Positive	1
Negative	70
Hard negative	6

19/3/19

Deep learning 4 object detection

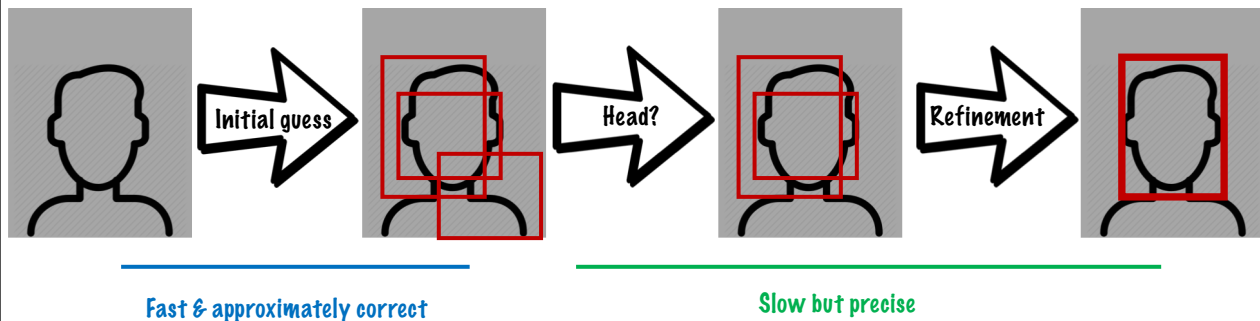
26

Image credit: <https://medium.com/@ageitgey/machine-learning-is-fun-part-3-deep-learning-and-convolutional-neural-networks-f40359318721#.dnbyjd6zg>

Object detection pipeline

Given the unbalanced nature of detection

What do we need?

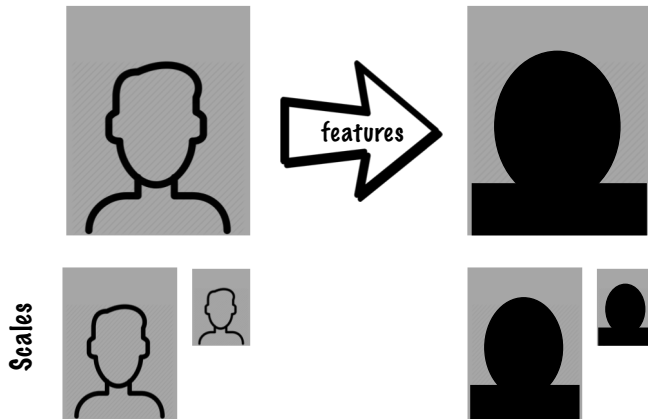


19/3/19

Deep learning 4 object detection

27

Initial guess



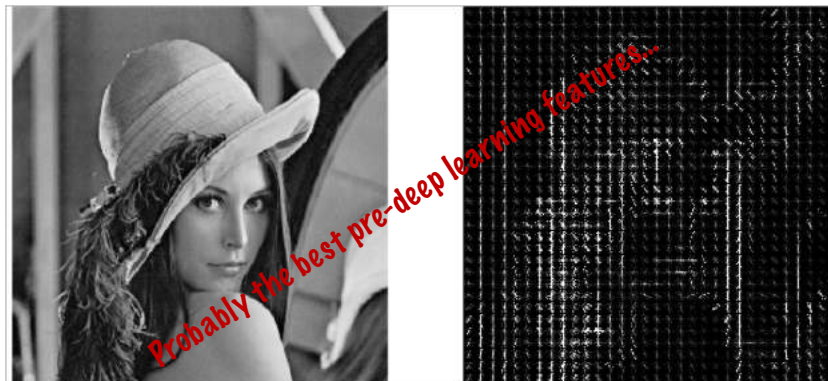
19/3/19

Deep learning 4 object detection

28

Features (in the past)

Histogram of Oriented Gradients



19/3/19

Deep learning 4 object detection

29

http://sharky93.github.io/docs/gallery/auto_examples/plot_hog.html

Features (currently)

Use pretrained neural networks!

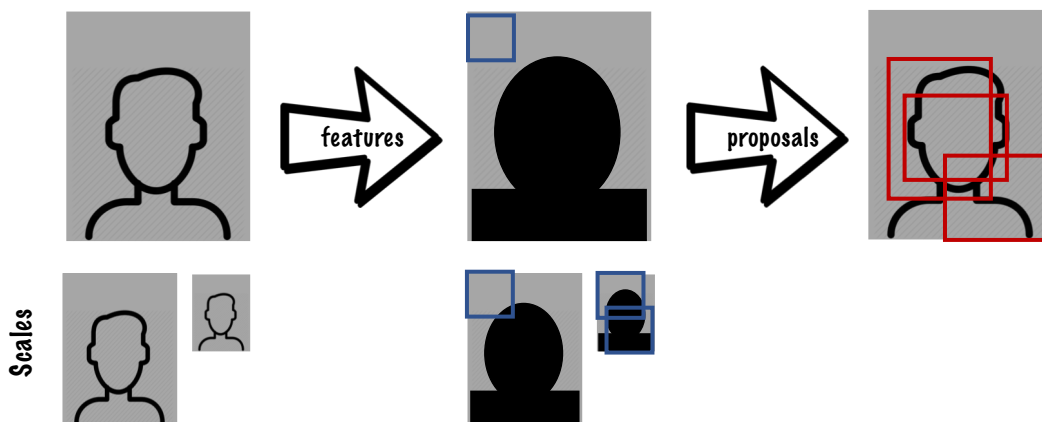
Architectures for Image Classification

19/3/19

Deep learning 4 object detection

30

Region proposals (Class agnostic classifiers)

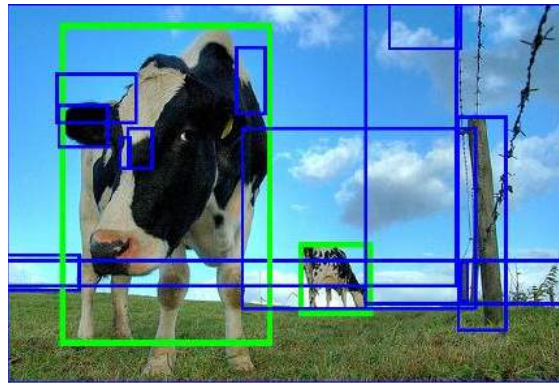


19/3/19

Deep learning 4 object detection

31

Region Proposals (Class agnostic classifiers)



Objectness [1]
 Selective search [2]
 Category-independent object proposals [3]
 Constrained parametric min-cuts (CPMC) [4]
 Multi-scale combinatorial grouping [5]

- [1] B. Alexe et al. Measuring the objectness of image windows.
- [2] J. Uijlings et al. Selective search for object recognition
- [3] I. Endres and D. Hoiem. Category independent object proposals
- [4] J. Carreira and C. Sminchisescu. CPMC: Automatic object segmentation using constrained parametric min-cuts
- [5] P. Arbelaez et al. Multiscale combinatorial grouping

19/3/19

Deep learning 4 object detection

Image: <https://ivi.fnwi.uva.nl/isis/publications/bibtexbrowser.php?key=UijlingsIJCV2013&bib=all.bib>

32

Region Proposals Selective search

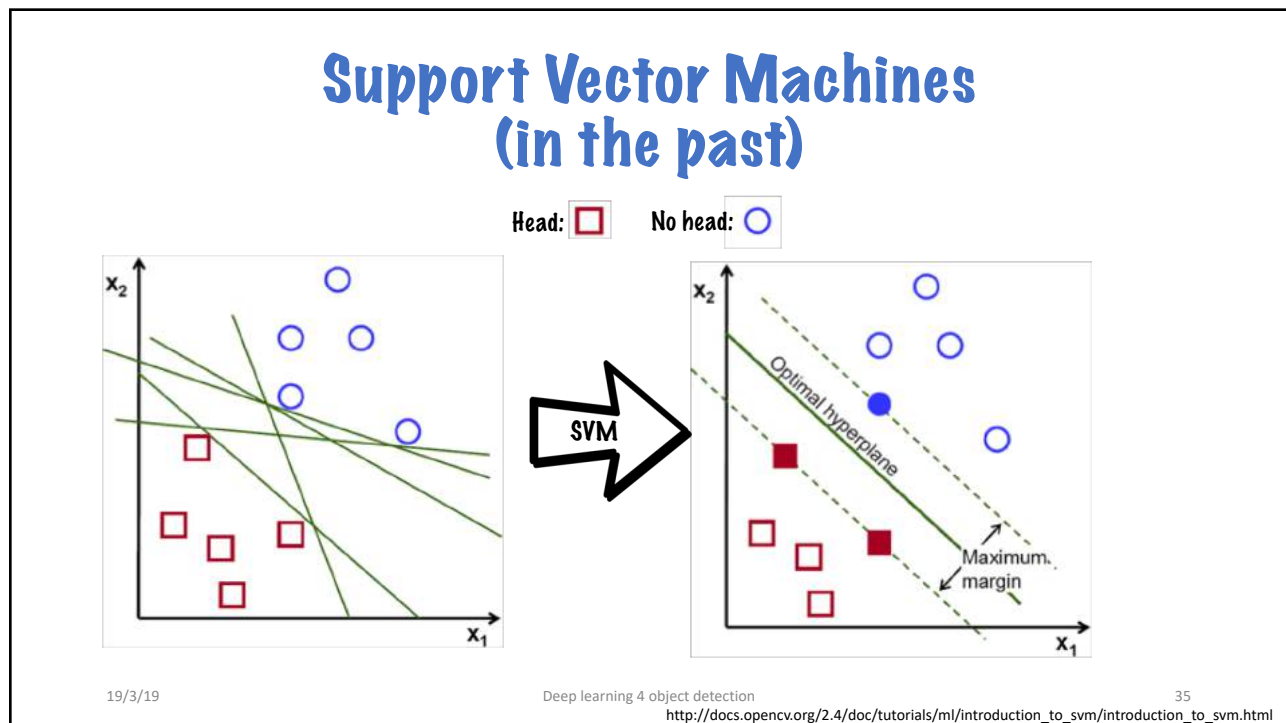
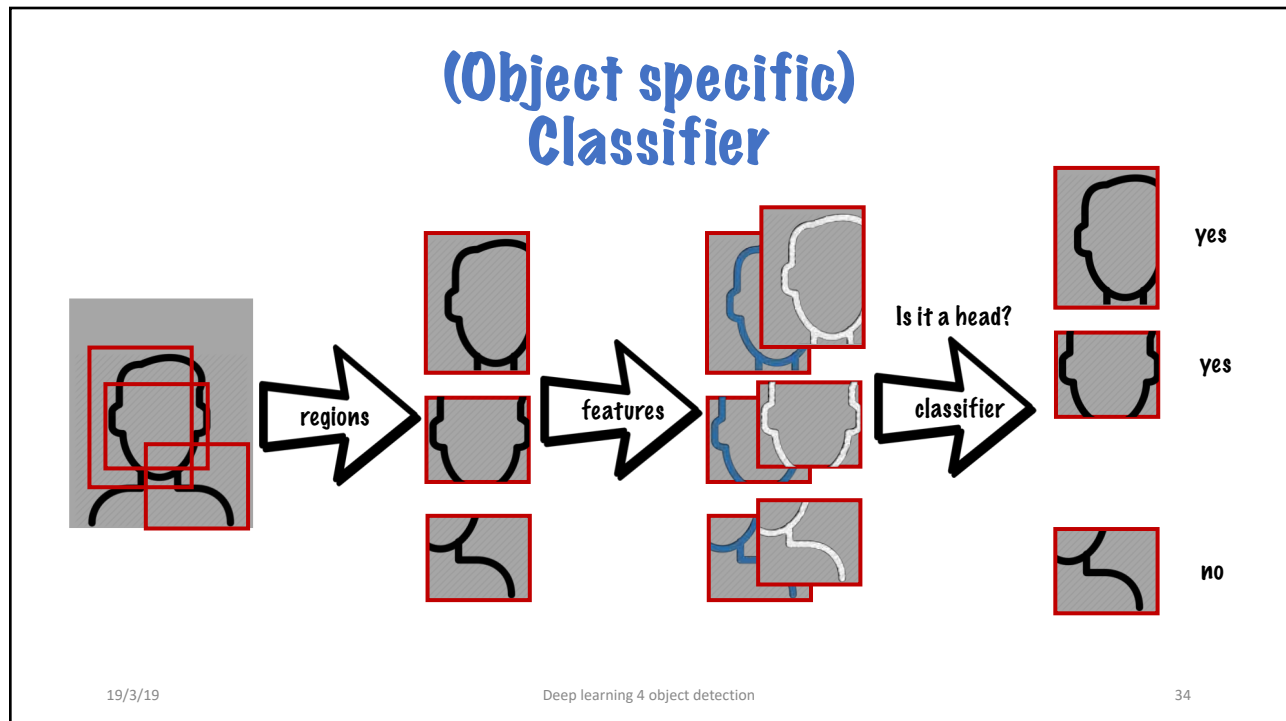


19/3/19

Deep learning 4 object detection

<http://koen.me/research/pub/vandesande-iccv2011-poster.pdf>

33



Classifier (Currently)

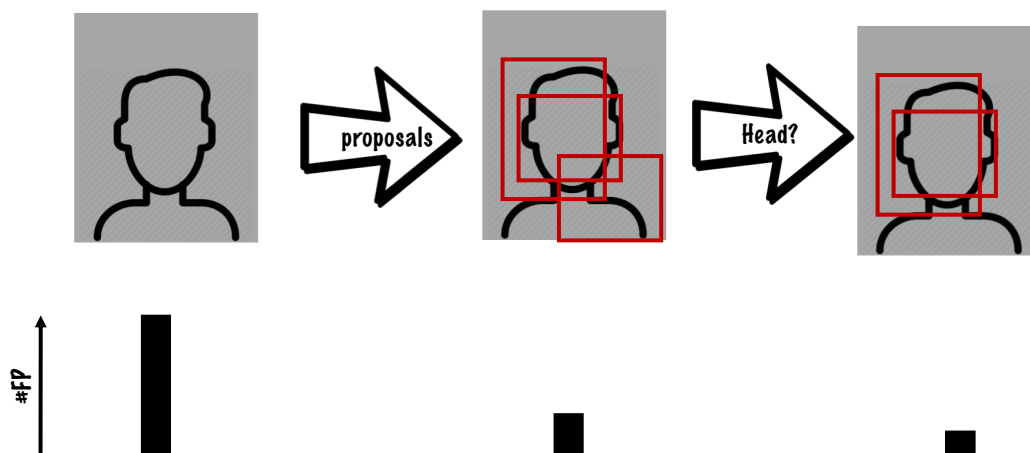
Neural Network!

19/3/19

Deep learning 4 object detection

36

What happens with FP?



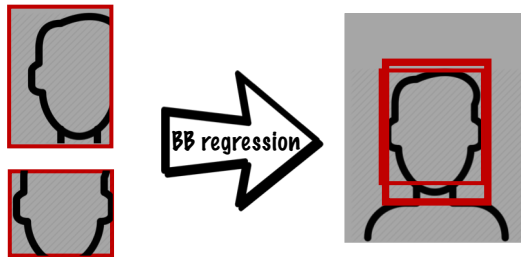
19/3/19

Deep learning 4 object detection

37

Regressor (Refinement)

Regressor is used to adjust the position of class specific bounding boxes.

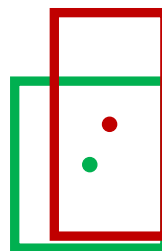
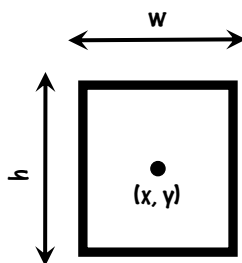


19/3/19

Deep learning 4 object detection

38

Bounding box regression



Annotation
(A_x, A_y, A_w, A_h)

Prediction
(P_x, P_y, P_w, P_h)

Regression error
(t_x, t_y, t_w, t_h)

$$\begin{aligned} t_x &= (A_x - P_x) / P_w \\ t_y &= (A_y - P_y) / P_h \\ t_w &= \log(A_w / P_w) \\ t_h &= \log(A_h / P_h) \end{aligned}$$

Why divide by P_w and P_h ?

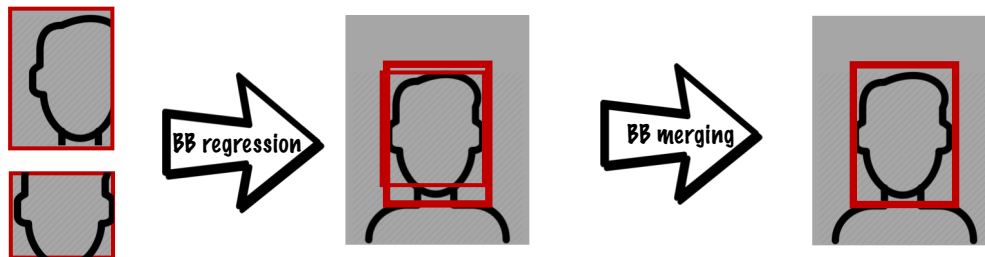
$$\text{Regression loss} = \text{loss}(t_x) + \text{loss}(t_y) + \text{loss}(t_w) + \text{loss}(t_h)$$

19/3/19

Deep learning 4 object detection

39

Regressor (Refinement)

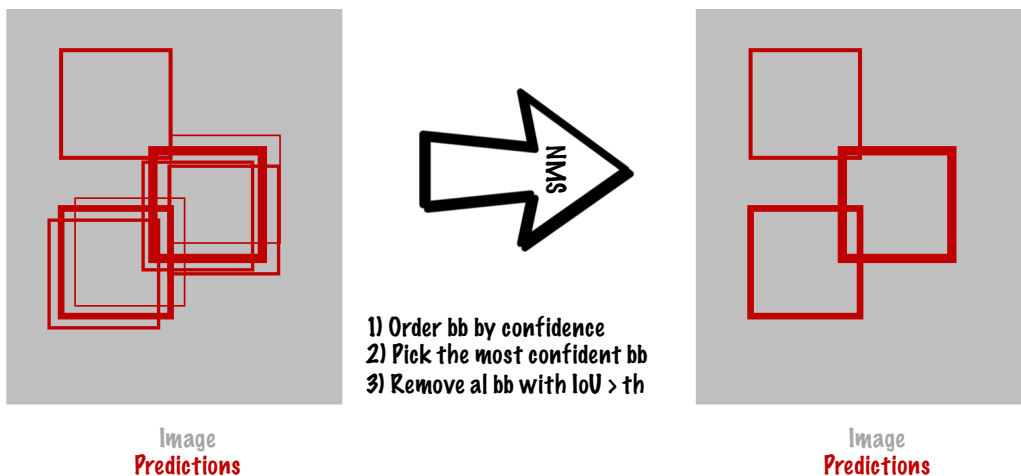


19/3/19

Deep learning 4 object detection

40

BB merging Non-Maximum Suppression (NMS)



19/3/19

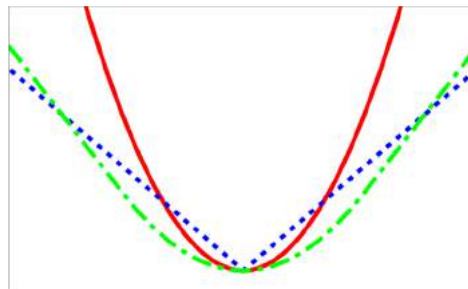
Deep learning 4 object detection

41

<http://www.computervisionblog.com/2011/08/blazing-fast-nmsm-from-exemplar-svm.html>

Detection: Loss Functions

- **Classification losses:**
 - Cross entropy (softmax)
 - Hinge loss (SVM)
- **Regression losses:**
 - **L1**
 - **Smooth L1**
 - **L2**



19/3/19

Deep learning 4 object detection

42

Index

Part 1
Introduction

Part 2
Basic blocks & concepts

Part 3
Models

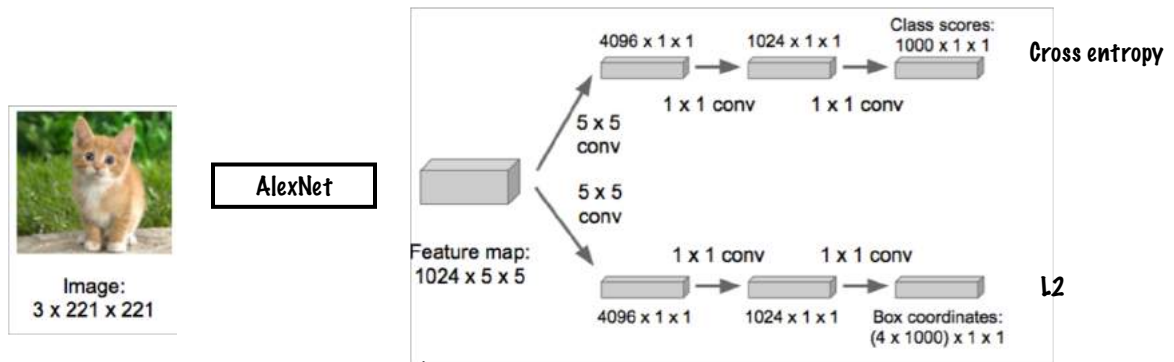
Bonus material
Weakly supervised
localization

19/3/19

Deep learning 4 object detection

43

OverFeat (2013, ICLR2014): training



- Two stage training:
1. Train the classifier (cross entropy)
 2. Train the regressor (L2)

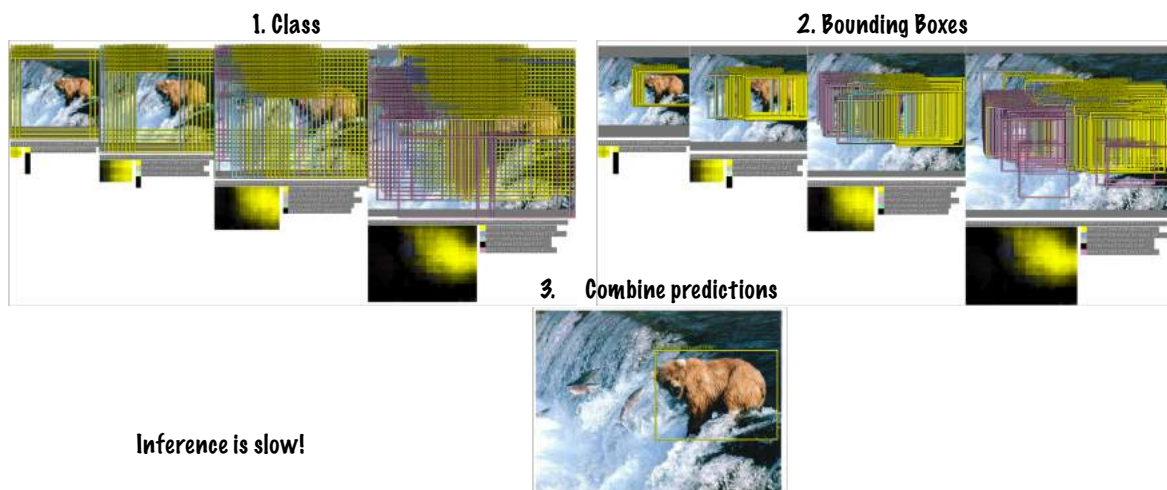
19/3/19

Deep learning 4 object detection

46

OverFeat (2013, ICLR2014): inference

Apply the network at all positions and scales and predict:

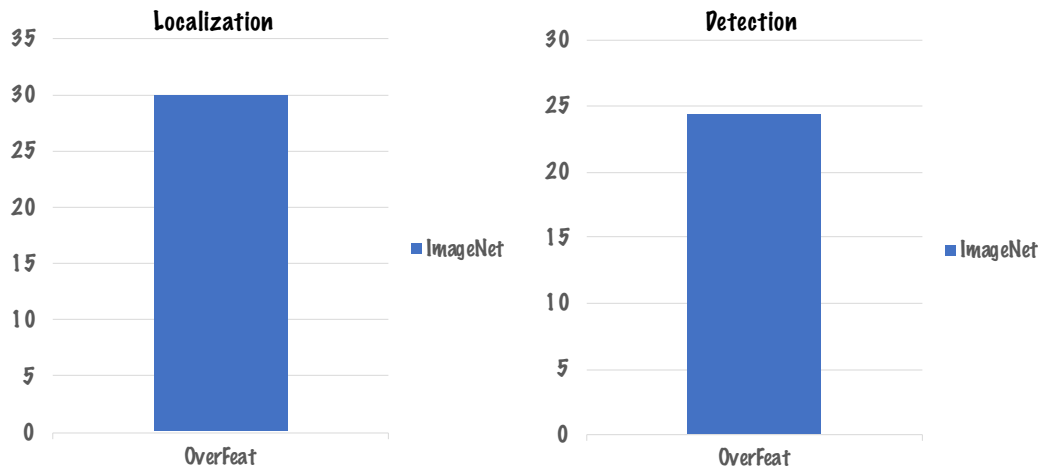


19/3/19

Deep learning 4 object detection

47

OverFeat (2013, ICLR2014): results

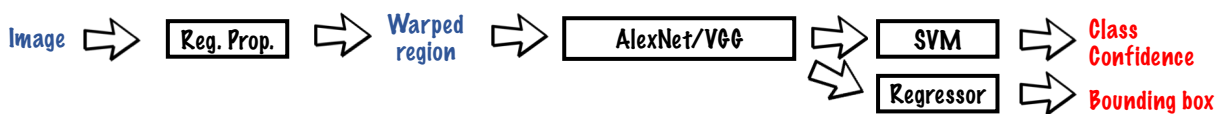
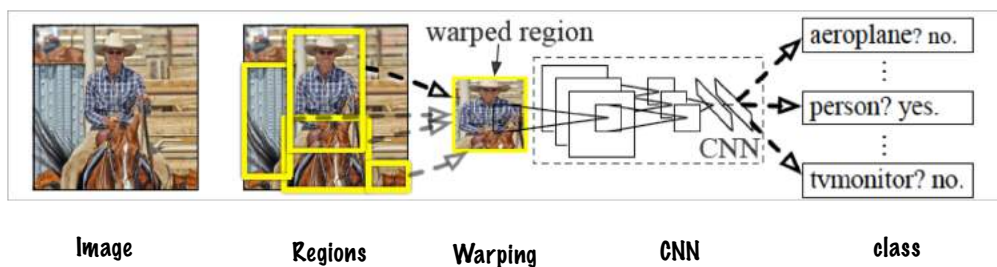


19/3/19

Deep learning 4 object detection

48

R-CNN (2013, CVPR2014)

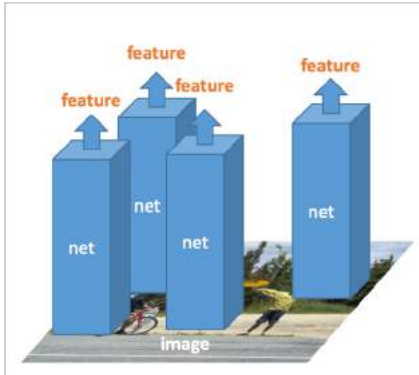
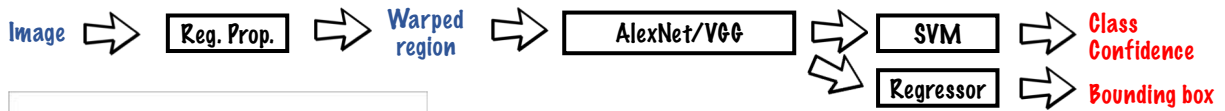


19/3/19

Deep learning 4 object detection

49

R-CNN (2013, CVPR2014): training



Training:

1. Pre-train network on Imagenet (image classification task)
2. Finetune network with softmax classifier (log loss)
3. Extract features
4. Train linear SVMs with hard negative mining (hinge loss)
5. Train bounding box regressions (least squares)

Training is slow (84h).

Why point #2?

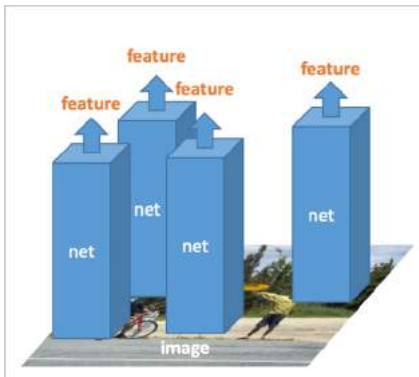
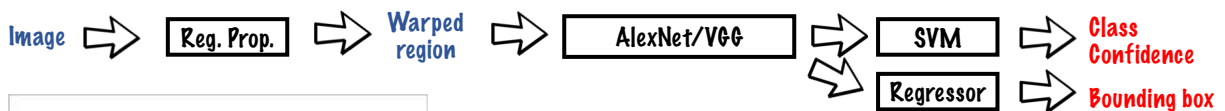
19/3/19

Deep learning 4 object detection

http://kaiminghe.com/eccv14sppnet/sppnet_ilsrc2014.pdf

50

R-CNN (2013, CVPR2014): inference



Inference:

1. Extract 2000 region proposals per image
2. Extract features for each proposal
3. Infer class, confidence and bounding box for each proposal

Inference is slow (2k passes of CNN per image).

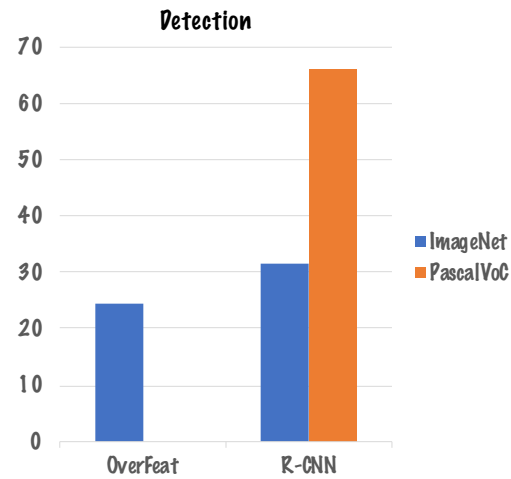
19/3/19

Deep learning 4 object detection

http://kaiminghe.com/eccv14sppnet/sppnet_ilsrc2014.pdf

51

R-CNN (2013, CVPR2014): results

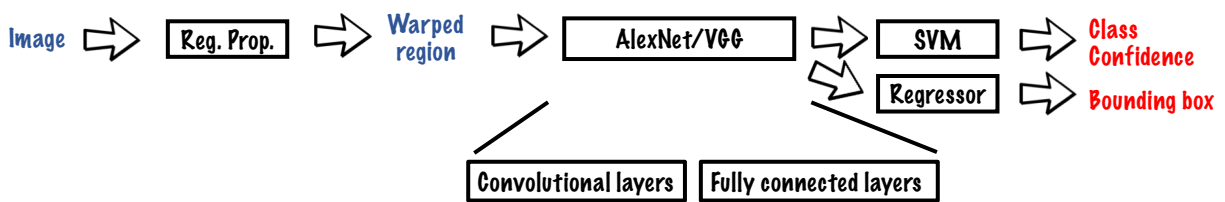


19/3/19

Deep learning 4 object detection

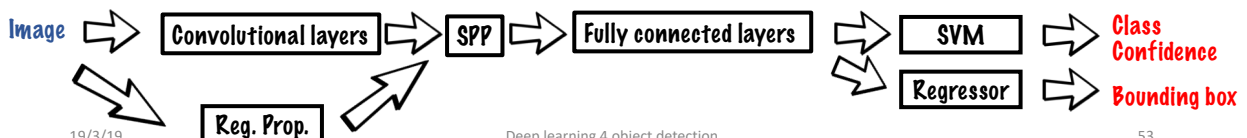
52

R-CNN (2013, CVPR2014)



Spatial Pyramid Pooling

SPP (2014, TPAMI2015)

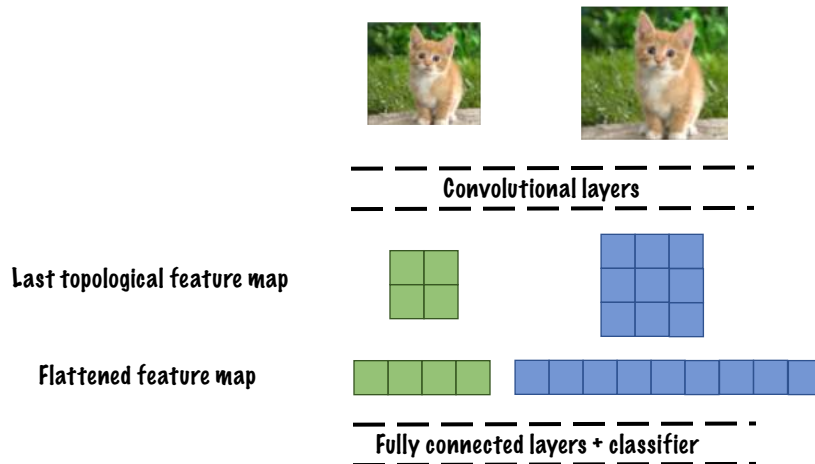


19/3/19

Deep learning 4 object detection

53

Why do we need to warp images?

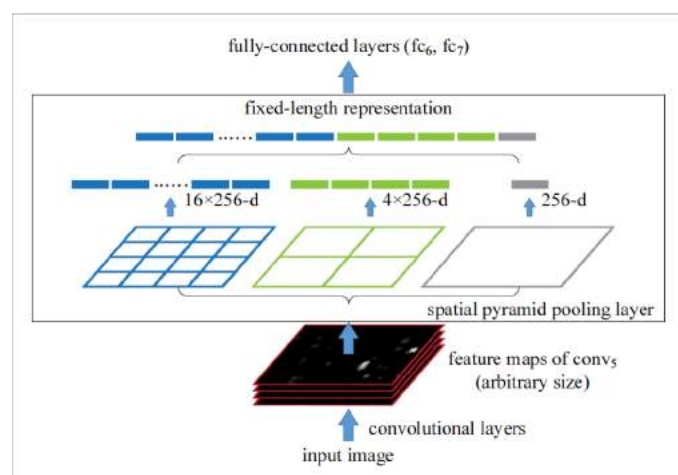


19/3/19

Deep learning 4 object detection

54

SPP (2014, TPAMI2015)

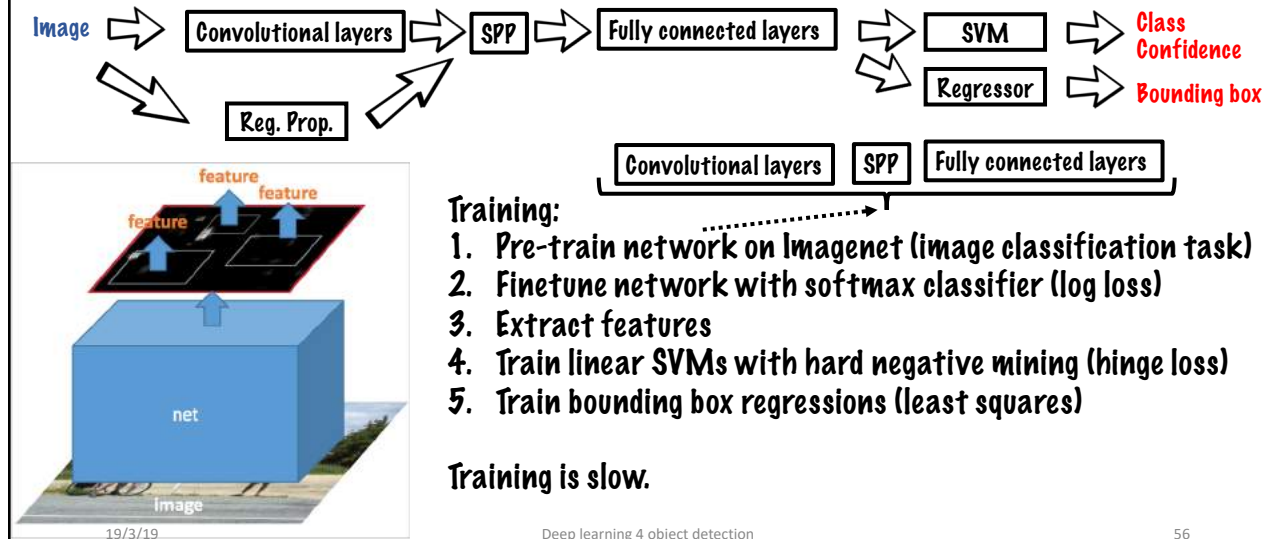


19/3/19

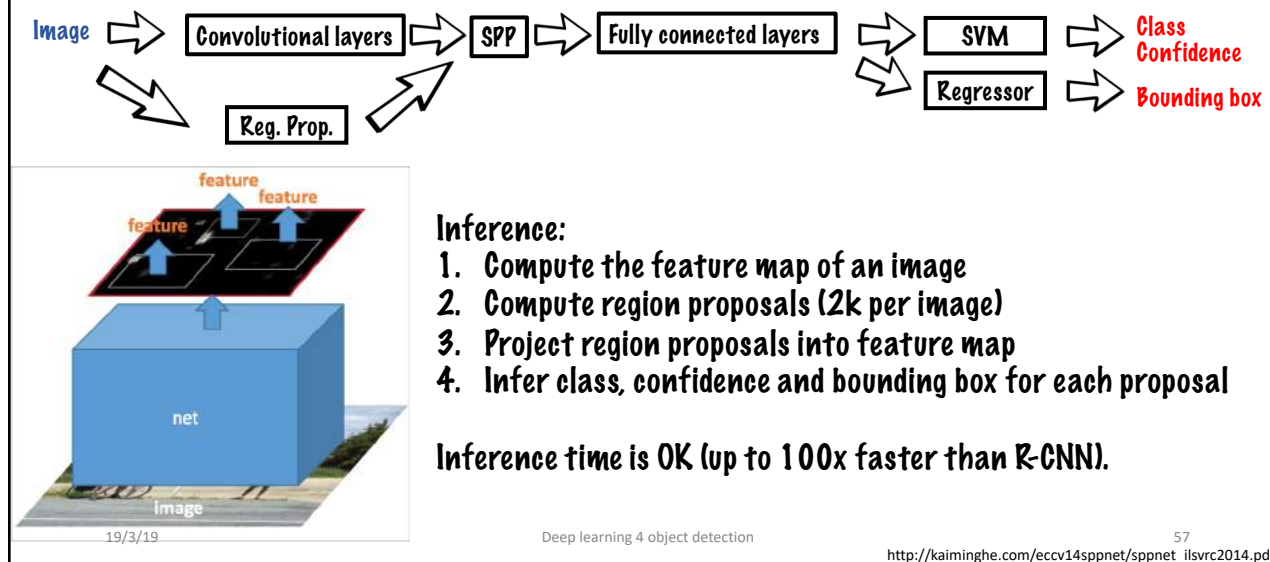
Deep learning 4 object detection

55

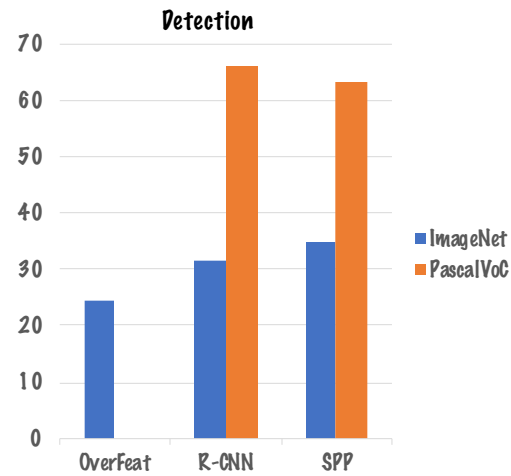
SPP (2014, TPAMI2015): training



SPP (2014, TPAMI2015): inference



SPP (2014, TPAMI2015): results

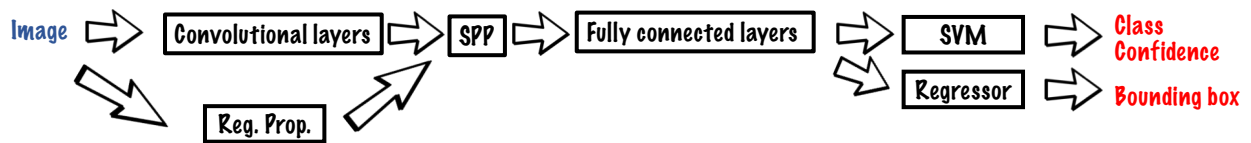


19/3/19

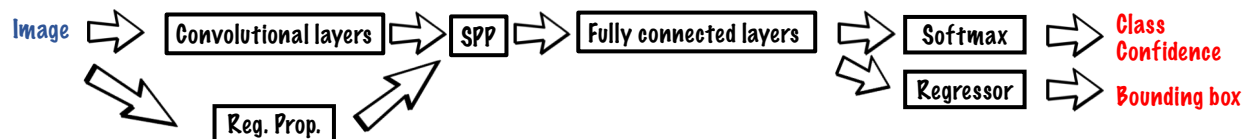
Deep learning 4 object detection

58

SPP (2014, TPAMI2015)



Fast R-CNN (2015, ICCV2015)

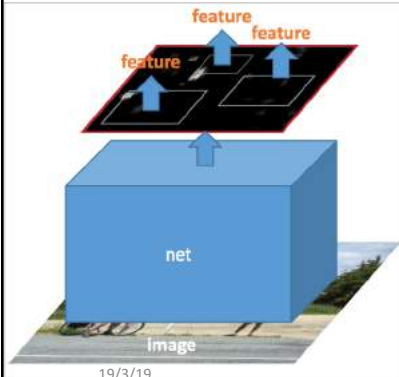
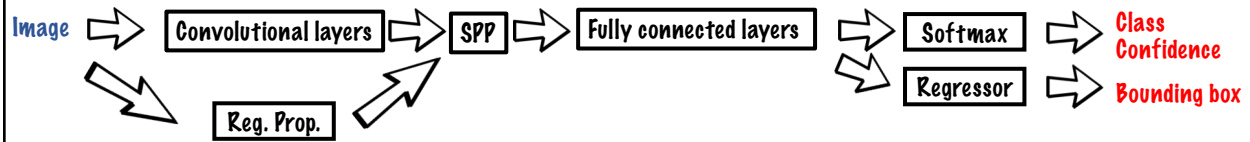


19/3/19

Deep learning 4 object detection

59

Fast R-CNN (2015, ICCV2015): training



Joint loss: log loss + smooth L1 loss

Training:

1. Pre-train network on Imagenet classification task
2. Train the model with joint loss

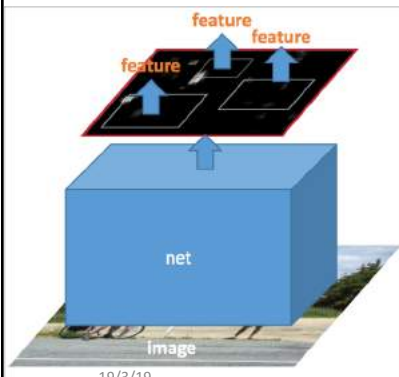
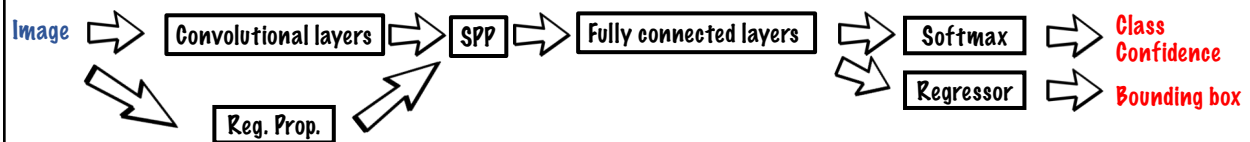
Training is elegant and fast.
Region Proposals are still required....

19/3/19

Deep learning 4 object detection

60

Fast R-CNN (2015, ICCV2015): inference



Inference:

1. Extract feature map
2. Extract region proposals
3. Infer class, confidence and bounding box for each proposal

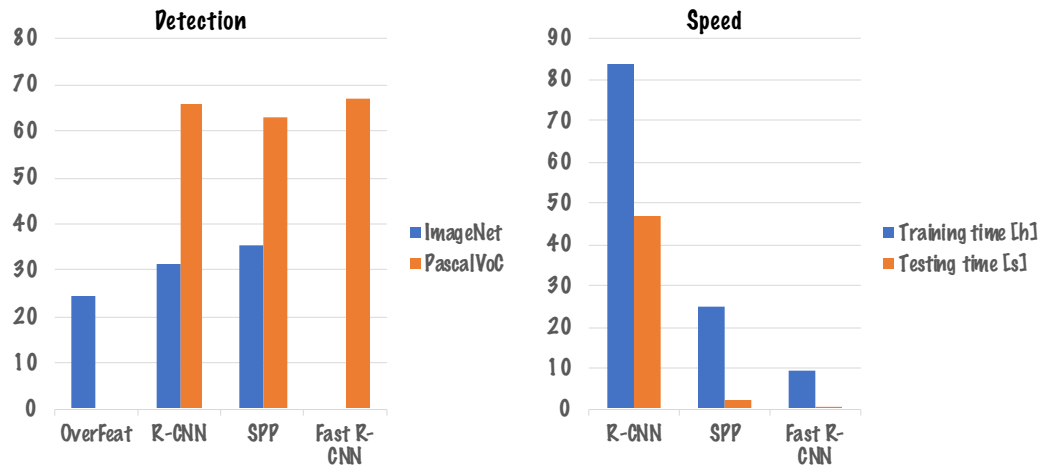
Inference is fast.

19/3/19

Deep learning 4 object detection

61

Fast R-CNN (2015, ICCV2015): results

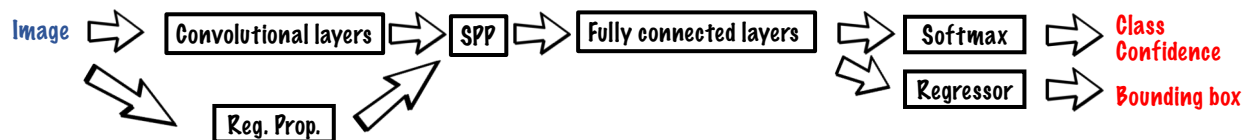


19/3/19

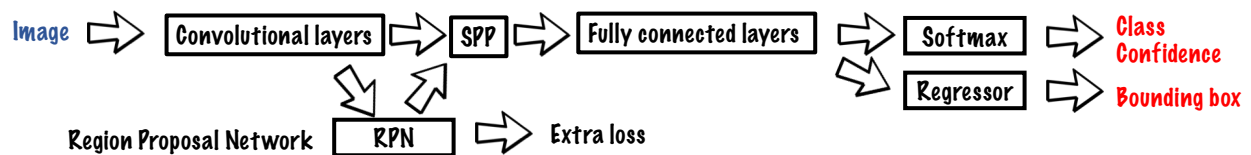
Deep learning 4 object detection

62

Fast R-CNN (2015, ICCV2015)



Faster R-CNN (2015, NIPS2015)

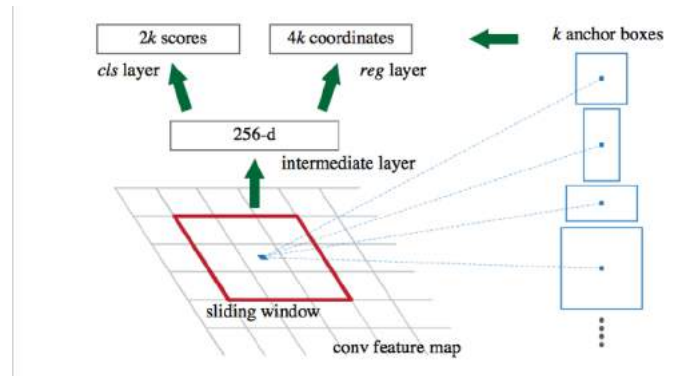


19/3/19

Deep learning 4 object detection

63

Faster R-CNN (2015, NIPS2015): Region Proposal Network (RPN)

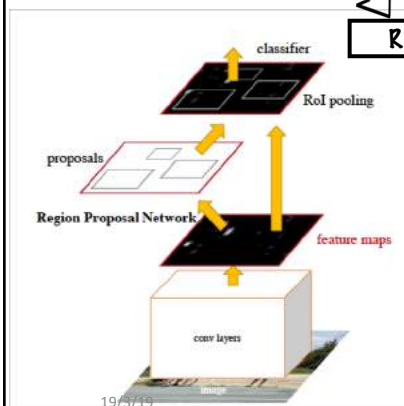


19/3/19

Deep learning 4 object detection

64

Faster R-CNN (2015, NIPS2015): training



19/3/19

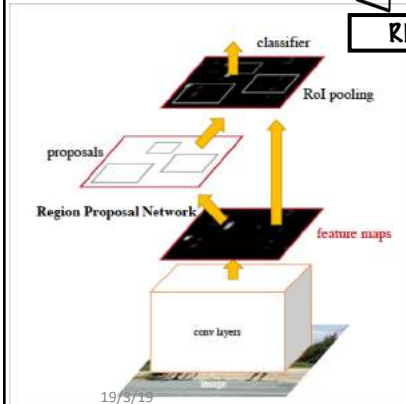
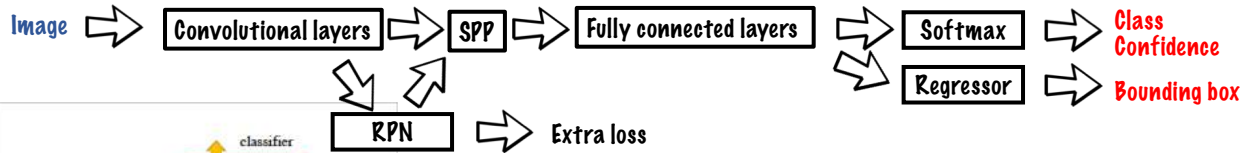
Deep learning 4 object detection

65

One network, four losses (TPAMI version)

1. RPN classification (anchor good / bad)
2. RPN regression (anchor → proposal)
3. Fast R-CNN classification (over classes)
4. Fast R-CNN regression (proposal → box)

Faster R-CNN (2015, NIPS2015): inference



Inference:

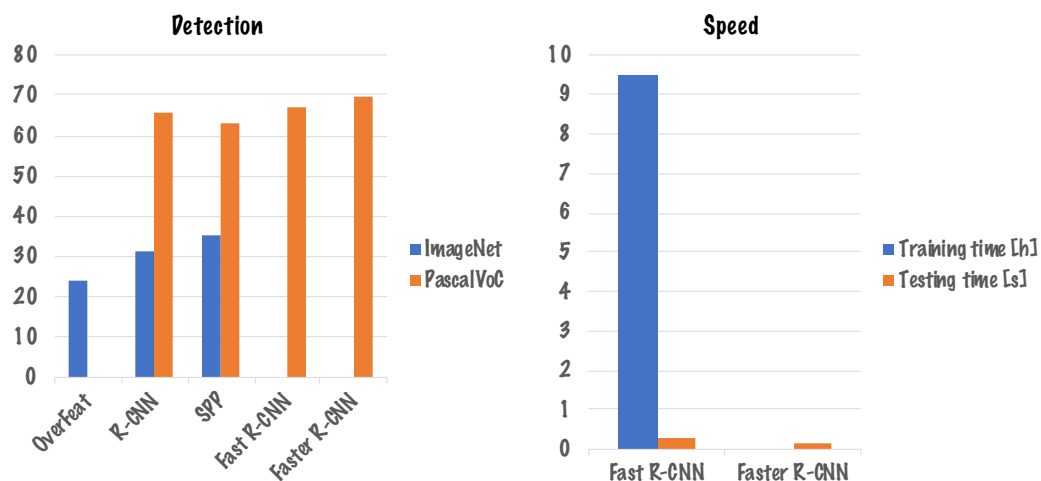
1. Extract feature map and region proposals
2. Infer class, confidence and bounding box for each proposal

19/3/19

Deep learning 4 object detection

66

Faster R-CNN (2015, NIPS2015): results

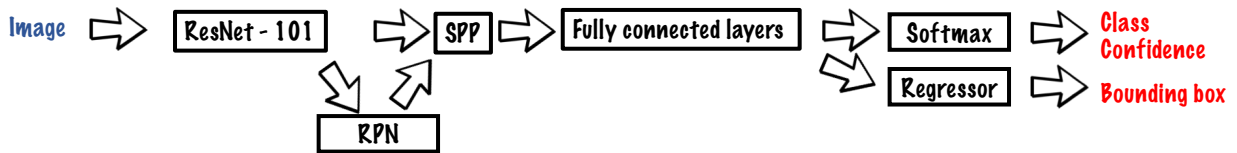


19/3/19

Deep learning 4 object detection

67

ResNets (2015, CVPR2016)



Faster R-CNN baseline	mAP@.5	mAP@.5:.95
VGG-16	41.5	21.5
ResNet-101	48.4	27.2

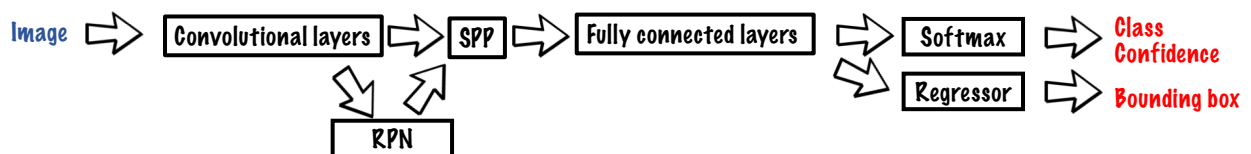
coco detection results

19/3/19

Deep learning 4 object detection

68

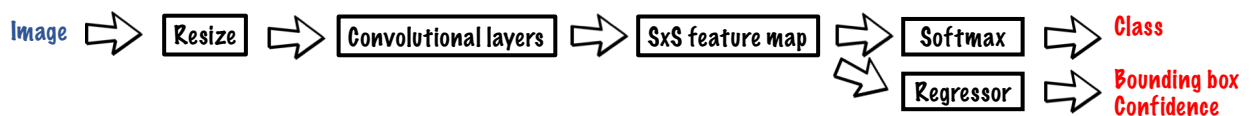
Faster R-CNN (2015, NIPS2015)



You Only Look Once

YOLO (2015, CVPR2016)

No need for region proposals.

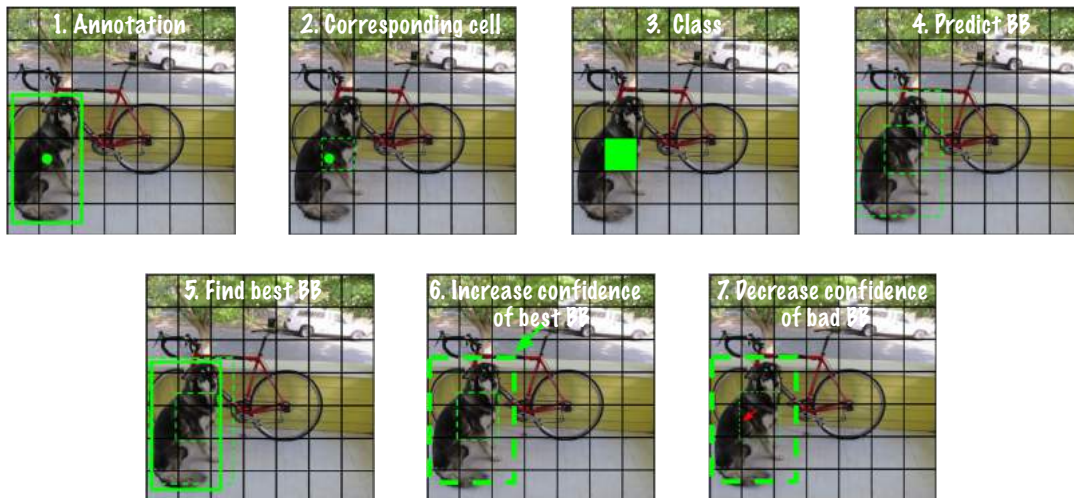


19/3/19

Deep learning 4 object detection

69

YOLO (2015, CVPR2016): key idea



19/3/19

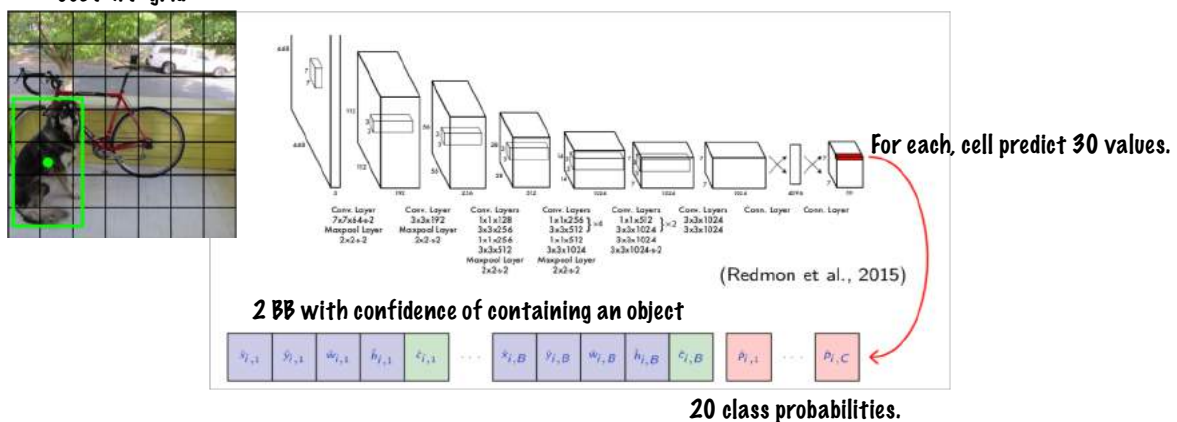
Deep learning 4 object detection 70
https://docs.google.com/presentation/d/1kAa7NOamBt4calBU9iHgT8a86RRH9Yz2oh4-GTdX6M/edit?slide=id.g151008b386_0_0

70

YOLO (2015, CVPR2016): architecture

For Pascal/VOC we have:

Use 7 x 7 grid



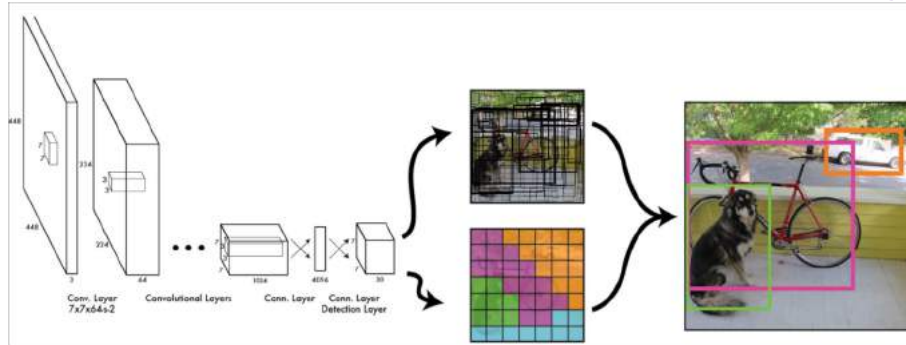
19/3/19

Deep learning 4 object detection

71

Francois Fleuret : Deep Learning Course. <https://fleuret.org/dlc/>

YOLO (2015, CVPR2016): training



Training:

1. Pre-train network on Imagenet classification task
2. Train the model with joint loss (quite engineered loss function)

19/3/19

Deep learning 4 object detection
https://docs.google.com/presentation/d/1kAa7NOamBt4calBU9IHgT8a86RRHz9Yz2oh4-GTdX6M/edit#slide=id.g151008b386_0_0

72

YOLO (2015, CVPR2016): training tricks

1. use 448×448 input for detection, instead of 224×224 ,
2. use Leaky ReLU for all layers,
3. dropout after the first fully connected layer,
4. normalize bounding boxes parameters in $[0, 1]$,
5. use a quadratic loss not only for the bounding box coordinates, but also for the confidence and the class scores,
6. reduce the weight of large bounding boxes by using the square roots of the size in the loss,
7. reduce the importance of empty cells by weighting less the confidence-related loss on them,
8. use momentum 0.9, decay $5e - 4$,
9. data augmentation with scaling, translation, and HSV transformation.

19/3/19

Deep learning 4 object detection

73

François Fleuret - Deep Learning Course, <https://fleuret.org/dlc>

SSD (2015, ECCV2016)

The diagram illustrates two deep learning architectures for object detection: SSD (Single Shot Detector) and YOLO Customized Architecture.

SSD Architecture:

- Input:** Image (300x300).
- VGG-16 through Conv5_3 layer:** The input image is processed by the VGG-16 network up to the Conv5_3 layer, resulting in feature maps of size 3x3x128.
- Extra Feature Layers:** These layers are added to the SSD architecture, including:
 - Classifier: Conv: 3x3x4x(Class+1)
 - Classifier: Conv: 3x3x6x(Class+1)
 - Conv: 3x3x128
 - Conv: 1x1x256
 - Conv: 1x1x128
 - Conv: 3x3x256-s1
 - Conv: 1x1x128
 - Conv: 3x3x256-s1
 - Conv: 1x1x128
 - Conv: 3x3x256-s1
- Detections:** 6732 per class.
- Non-Maximum Suppression:** Results in 74.3mAP and 59FPS.

YOLO Customized Architecture:

- Input:** Image (448x448).
- YOLO Customized Architecture:** The input image is processed by a customized YOLO architecture, resulting in feature maps of size 7x7x1024.
- Fully Connected:** The feature maps are passed through two fully connected layers.
- Detections:** 98 per class.
- Non-Maximum Suppression:** Results in 83.4mAP and 45FPS.

YOLOv2 (2016)

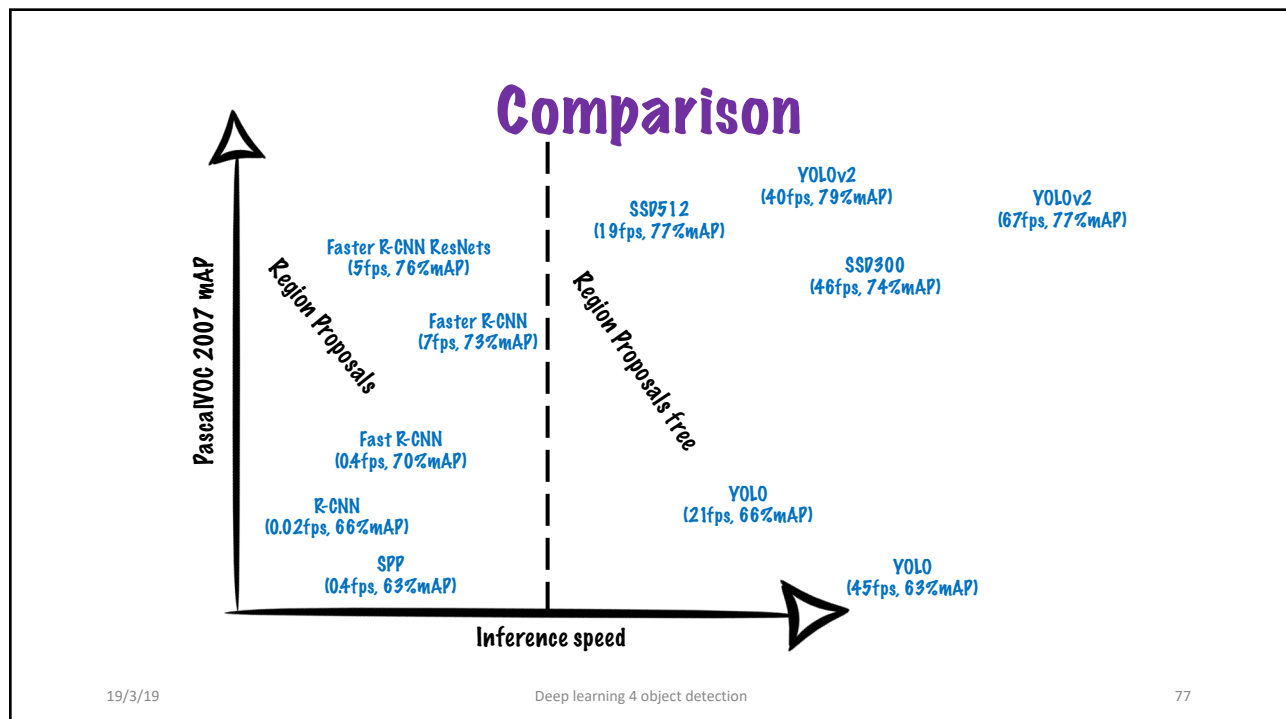
	YOLO								YOLOv2
batch norm?		✓	✓	✓	✓	✓	✓	✓	✓
hi-res classifier?			✓	✓	✓	✓	✓	✓	✓
convolutional?				✓	✓	✓	✓	✓	✓
anchor boxes?				✓	✓				
new network?					✓	✓	✓	✓	✓
dimension priors?						✓	✓	✓	✓
location prediction?						✓	✓	✓	✓
passthrough?							✓	✓	✓
multi-scale?								✓	✓
hi-res detector?									✓
VOC2007 mAP	63.4	65.8	69.5	69.2	69.6	74.4	75.4	76.8	78.6

There are a lot of tricks to get a good architecture for object detection...

19/3/19

Deep learning 4 object detection

76



19/3/19

Deep learning 4 object detection

77

Interesting papers (not covered in the class)

Feature Pyramid Networks for Object Detection

Tsung-Yi Lin^{1,2}, Piotr Dollár¹, Ross Girshick¹,
Kaiming He¹, Bharath Hariharan¹, and Serge Belongie²

¹Facebook AI Research (FAIR)

²Cornell University and Cornell Tech

Focal Loss for Dense Object Detection

Tsung-Yi Lin Priya Goyal Ross Girshick Kaiming He Piotr Dollár
Facebook AI Research (FAIR)

Mask R-CNN

Kaiming He Georgia Gkioxari Piotr Dollár Ross Girshick
Facebook AI Research (FAIR)

19/3/19

Deep learning 4 object detection

78

Index

Part 1

Introduction

Part 2

Basic blocks & concepts

Part 3

Models

Bonus material

Weakly supervised localization

19/3/19

Deep learning 4 object detection

79

Weakly supervised localization

We have:



We want:



Labeling is expensive!

Is object localization for free?

Let us have a look at one paper [1]

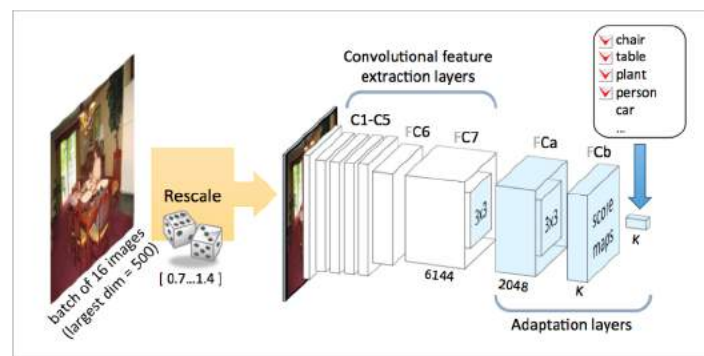
19/3/19

Deep learning 4 object detection

80

[1] <http://www.di.ens.fr/~josef/publications/Oquab15.pdf>

How is it done?



19/3/19

Deep learning 4 object detection

81

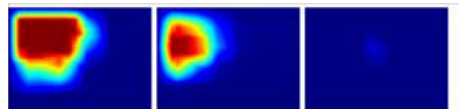
<http://www.di.ens.fr/willow/research/weakcnn/>

Is object localization for free?

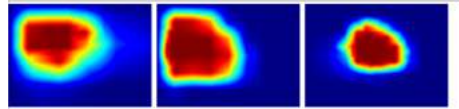
Iteration 210



Iteration 510



Iteration 4200



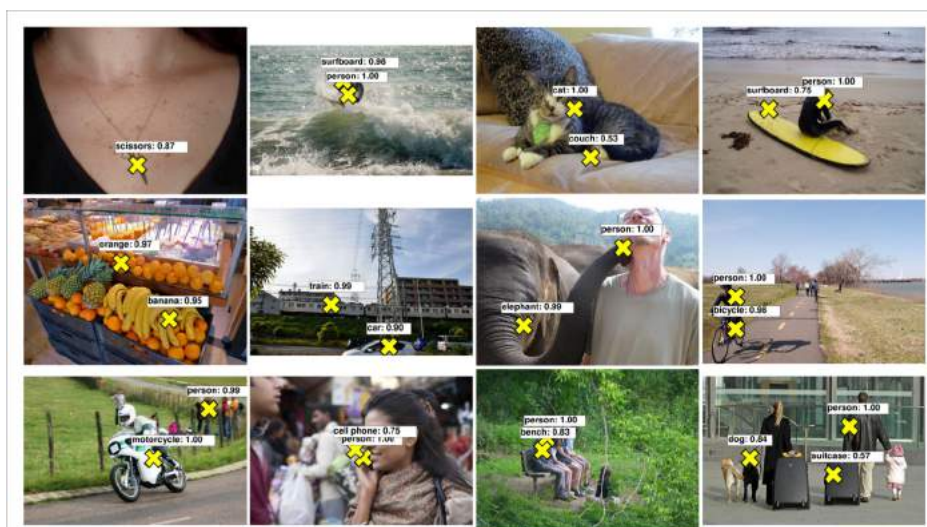
19/3/19

Deep learning 4 object detection

82

<http://www.di.ens.fr/willow/research/weakcnn/>

Results



19/3/19

Deep learning 4 object detection

83

Wrap up

Datasets
Evaluation

Part 1
Introduction

Part 2
Basic blocks & concepts

Object detection pipeline
Hard Negative Mining
Non-Maximum Suppression
Region Proposals

Part 3
Models

SPP
OverFeat
R-CNN
Fast R-CNN
Faster R-CNN
YOLO
SSD

Bonus material
Weakly supervised
localization

Is object localization for free?

19/3/19

Deep learning 4 object detection

84

References

- *Ross B. Girshick, Jeff Donahue, Trevor Darrell and Jitendra Malik*; Rich feature hierarchies for accurate object detection and semantic segmentation.
- *Pierre Sermanet, David Eigen, Xiang Zhang, Michael Mathieu, Rob Fergus, Yann LeCun*; OverFeat: Integrated Recognition, Localization and Detection using Convolutional Networks.
- *Kaiming He and Xiangyu Zhang and Shaoqing Ren and Jian Sun*; Spatial Pyramid Pooling in Deep Convolutional Networks for Visual Recognition.
- *Ross B. Girshick*; Fast R-CNN.
- *Shaoqing Ren, Kaiming He, Ross B. Girshick and Jian Sun*; Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks.
- *Joseph Redmon, Santosh Kumar Divvala, Ross B. Girshick, Ali Farhadi*; You Only Look Once: Unified, Real-Time Object Detection.
- *Kaiming He, Xiangyu Zhang, Shaoqing Ren, Jian Sun*; Deep Residual Learning for Image Recognition.
- *Joseph Redmon and Ali Farhadi*; YOLO9000: Better, Faster, Stronger.
- *M. Oquab and L. Bottou and I. Laptev and J. Sivic*; Is object localization for free? - Weakly-supervised learning with convolutional neural networks.

19/3/19

Deep learning 4 object detection

85

Deep Learning 4

Object Detection



Michal Drozdal
mdrozdal@fb.com