



Module 8: Supervised learning

5LSM0: Convolutional neural networks for computer vision

Fons van der Sommen

Electrical Engineering / VCA research group



Administrative

We're half-way!

Assignments

- Assignment 1 due
- Assignment 2 online [due Mar-25]
- Start on final assignment as soon as you've finished assignment 2

Staff changes

- Joost → Chengyang
 - *Pytorch expert & smart guy*

You are here!



WK	Date	Module
6	Tue Feb-5	1. Introduction
	Thu Feb-7	2. Data-driven image classification
7	Tue Feb-12	3. Loss functions and optimization
	Thu Feb-14	4. Neural networks and backpropagation
8	Tue Feb-19	5. Convolutional neural networks
	Thu Feb-21	6. CNN architectures
9	Tue Feb-26	7. Training neural networks (part I)
	Thu Feb-28	7. Training neural networks (part II)
10		<u>NO LECTURES</u>
11	Tue Mar-12	8. Supervised learning
	Thu Mar-14	9. Unsupervised learning
12	Tue Mar-19	10. Semi/self-supervised learning
	Thu Mar-21	11. Temporal modeling & reinforcement learning
13	Tue Mar-26	12. Visualization and understanding
	Thu Mar-28	<u>PAPER SESSIONS</u>
14	Tue Apr-2	13. Efficient deep learning
	Thu Apr-4	Recap + question hours
15+16		<u>EXAMS</u>



Last time(s): training neural networks

Weight initialization

Layer mean

Layer standard deviation

Output distribution over the layers

Batch norm

$$\hat{x}^{(k)} = \frac{x^{(k)} - E[x^{(k)}]}{\sqrt{\text{Var}[x^{(k)}]}}$$

Original data

Zero-mean data

Normalized data

Data normalization

3 SLSMO Module 8: Supervised learning

Last time(s): training neural networks

Saddle points

Local minima

Poor conditioning

Optimization

SGD

$$x_{t+1} = x_t - \alpha \nabla f(x_t)$$

SGD+ momentum

$$x_{t+1} = x_t - \alpha v_{t+1}$$

$$v_{t+1} = \rho v_t + \nabla f(x_t)$$

AdaGrad

$$x_{t+1} = x_t - \frac{\alpha \nabla f(x_t)}{g_{\cdot}^{-1/2}}$$

$$g_{t+1,i} = \gamma g_{t,i} + (1 - \gamma) \left(\frac{\partial L}{\partial w_i} \right)_t^2$$

ADAM

$$x_{t+1} = x_t - \alpha \frac{m_{t+1}}{\sqrt{v_{t+1} + \epsilon}}$$

$$m_{t+1} = \beta_1 m_t + (1 - \beta_1) \nabla f(x_t)$$

$$v_{t+1} = \beta_2 v_t + (1 - \beta_2) (\nabla f(x_t))^2$$

4 SLSMO Module 8: Supervised learning

Last time(s): training neural networks

Step decay learning rate

$$J(\theta) \approx J(\theta_0) + (\theta - \theta_0)^T \nabla_{\theta} J(\theta_0) + \frac{1}{2} (\theta - \theta_0)^T H (\theta - \theta_0)$$

$$\theta^* = \theta_0 - H^{-1} \nabla_{\theta} J(\theta_0)$$

First-order approximation

Second-order approximation

More optimization...

5 SLSM0 Module 8: Supervised learning

Last time(s): training neural networks

Re-use

Regularization

How can we reduce this gap?

Data augmentation

Transfer learning

CNN codes

Conventional classifier (e.g. SVM)

Dropout

6 SLSM0 Module 8: Supervised learning

This time

Supervised learning!

Classification (recap)

- What **image** belongs to what **category**?

Segmentation

- What **pixels** belong to what **category**?

Detection

- What **pixels** belong to what **instance** of what **category**?



7

5LSM0 Module 8: Supervised learning

TU/e

Supervised learning

So far: **just classification**

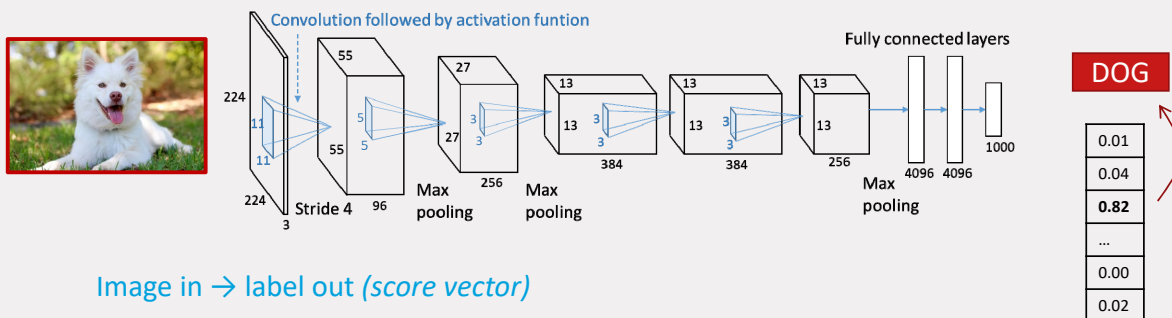


Image in → label out (*score vector*)

Typically an image shows more than a single class though...



8

5LSM0 Module 8: Supervised learning

TU/e

Images license free (www.pexels.com)

Supervised learning

DOG 2

DOG

DOG

DOG 1 DOG 2

DOG DOG STICK

9 5LSMO Module 8: Supervised learning

TU/e

Supervised learning

Classification+localization	Object detection	Instance segmentation	Segmentic segmentation
Car	Pedestrian Car Pedestrian	Pedestrian Car Pedestrian	Truck Road Vegetation Pedestrian Sidewalk Car
Single object	Multiple objects		No objects, just pixels

10 5LSMO Module 8: Supervised learning

TU/e

Semantic segmentation

- Assign a class to each pixel in an image
- Do not distinguish between several instances of the same category

Pedestrian

Road

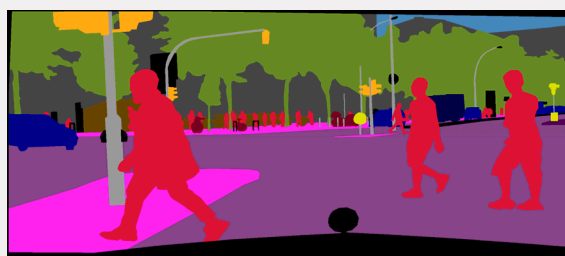
Truck

Car

Vegetation

Sidewalk

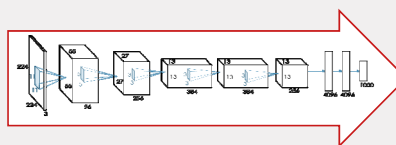
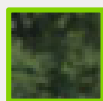
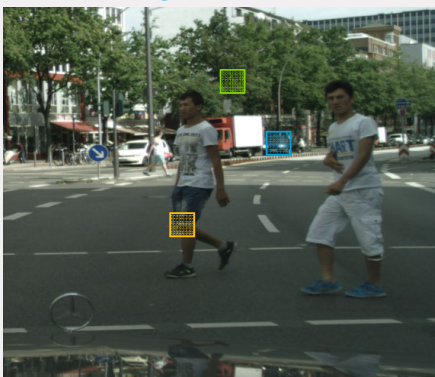
Q: Main difference with instance segmentation?



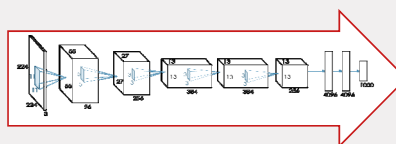
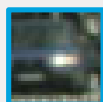
Semantic segmentation

Q: Why is this not a good idea?

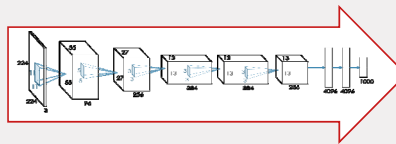
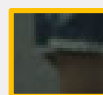
Idea: sliding window



Car



Car





Pedestrian



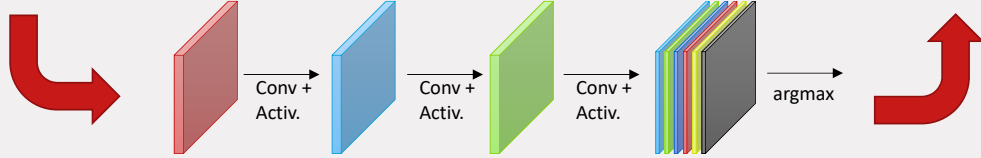
Semantic segmentation


Idea: fully convolutional


Q1: How could we characterize the loss in this case?

Q2: Problem in this case?




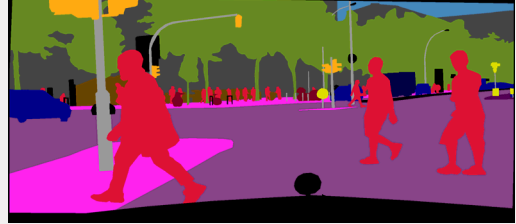


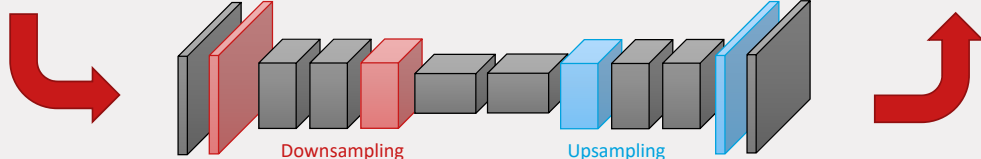
13 5LSMO Module 8: Supervised learning




Semantic segmentation


Idea: fully convolutional





14 5LSMO Module 8: Supervised learning



Semantic segmentation

Upsampling... but how?

- Unpooling

6	8
3	4

→

6	6	8	8
6	6	8	8
3	3	4	4
3	3	4	4

Nearest neighbor

6	8
3	4

→

6	0	8	0
0	0	0	0
3	0	4	0
0	0	0	0

"Bed of nails"

15 5LSM0 Module 8: Supervised learning

Semantic segmentation

Upsampling... but how?

- Unpooling

1	1	2	4
5	6	8	7
3	2	1	0
1	2	3	4

→

6	8
3	4

→

0	0	0	0
0	6	8	0
3	0	0	0
0	0	0	4

Max unpooling

Max positions

→

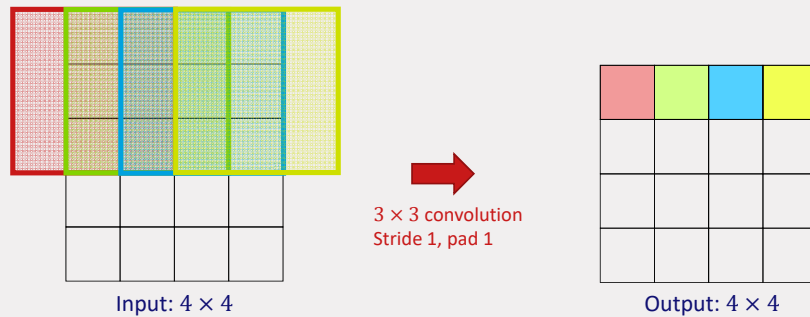
16 5LSM0 Module 8: Supervised learning

Semantic segmentation

Upsampling... but how?

- Unpooling
- **Transpose convolution**

Recall: normal convolution



17 5LSM0 Module 8: Supervised learning

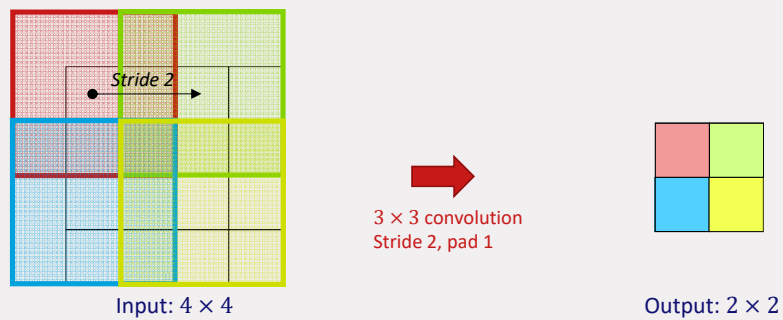
TU/e

Semantic segmentation

Upsampling... but how?

- Unpooling
- **Transpose convolution**

Strided convolution



18 5LSM0 Module 8: Supervised learning

TU/e

Semantic segmentation

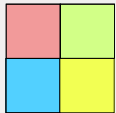
Upsampling... but how?

- Unpooling
- **Transpose convolution**

Q: Sometimes called *deconvolution*, why is this inadequate nomenclature?

Transpose convolution

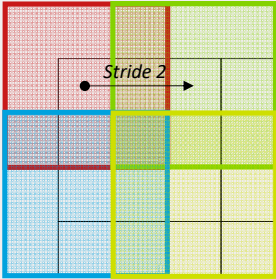
Multiply input value with filter kernel
Sum where overlaps




Output: 2 × 2

➔


3 × 3 convolution
Stride 2, pad 1



Input: 4 × 4



19 5LSM0 Module 8: Supervised learning




Semantic segmentation

Upsampling... but how?

- Unpooling
- **Transpose convolution**

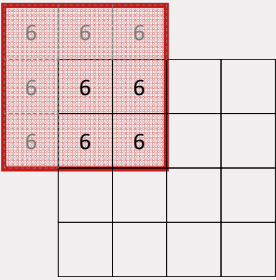
Transpose convolution




Output: 2 × 2

➔


transpose convolution
stride 2, pad 1
3 × 3 kernel of ones




Input: 4 × 4



3 × 3 kernel of ones



20 5LSM0 Module 8: Supervised learning



Semantic segmentation

Upsampling... but how?

- Unpooling
- **Transpose convolution**

6	8
3	4

Output: 2 × 2

→

transpose convolution
stride 2, pad 1
3 × 3 kernel of ones

Transpose convolution

6	6	14	8	8
6	6	14	8	8
6	6	14	8	8

Input: 4 × 4

1	1	1
1	1	1
1	1	1

3 × 3 kernel of ones

+8

21 5LSM0 Module 8: Supervised learning

Semantic segmentation

Upsampling... but how?

- Unpooling
- **Transpose convolution**

6	8
3	4

Output: 2 × 2

→

transpose convolution
stride 2, pad 1
3 × 3 kernel of ones

Transpose convolution

6	6	14	8	8
6	6	14	8	8
9	9	17	8	8
3	3	3		
3	3	3		

Input: 4 × 4

1	1	1
1	1	1
1	1	1

3 × 3 kernel of ones

+3

22 5LSM0 Module 8: Supervised learning

Semantic segmentation

Upsampling... but how?

- Unpooling
- **Transpose convolution**

1	1	1
1	1	1
1	1	1

3 × 3 kernel of ones

Transpose convolution

6	8
3	4

Output: 2 × 2

transpose convolution
stride 2, pad 1
3 × 3 kernel of ones

6	6	14	8	8
6	6	14	8	8
9	9	21	12	12
3	3	7	4	4
3	3	7	4	4

Input: 4 × 4 +4

23 5LSM0 Module 8: Supervised learning

Semantic segmentation

Q: Output size of this convolution?

Convolution as matrix multiplication

- We can write a convolution as a matrix multiplication

$$p = \begin{bmatrix} x \\ y \\ z \end{bmatrix}$$

$$q = \begin{bmatrix} a \\ b \\ c \\ d \end{bmatrix}$$

$$y = p * q$$

$$y_n = \sum_{k=-\infty}^{\infty} p_k q_{n-k}$$

ax

bx + ay

cx + by + az

dz

24 5LSM0 Module 8: Supervised learning

Semantic segmentation

Convolution as matrix multiplication

- We can write a convolution as a matrix multiplication

$$\mathbf{p} * \mathbf{q} = \mathbf{Qp}_{00} \quad \leftarrow \text{Zero-padded}$$

$\mathbf{p} = \begin{bmatrix} x \\ y \\ z \end{bmatrix}$

$\mathbf{q} = \begin{bmatrix} a \\ b \\ c \\ d \end{bmatrix}$

$\mathbf{y} = \mathbf{p} * \mathbf{q} \quad y_n = \sum_{k=-\infty}^{\infty} p_k q_{n-k}$

$\begin{bmatrix} d & c & b & a & 0 & 0 & 0 & 0 & 0 \\ 0 & d & c & b & a & 0 & 0 & 0 & 0 \\ 0 & 0 & d & c & b & a & 0 & 0 & 0 \\ 0 & 0 & 0 & d & c & b & a & 0 & 0 \\ 0 & 0 & 0 & 0 & d & c & b & a & 0 \\ 0 & 0 & 0 & 0 & 0 & d & c & b & a \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ x \\ y \\ z \\ 0 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} ax \\ bx + ay \\ cx + by + az \\ dx + cy + bz \\ dy + cz \\ dz \end{bmatrix}$



Semantic segmentation

Transpose convolution as matrix multiplication

$$\mathbf{p} *^T \mathbf{q} = \mathbf{Q}^T \mathbf{p}_{00} \quad = \mathbf{p} * \mathbf{q} \quad \text{For stride} = 1$$

$$\begin{bmatrix} d & 0 & 0 & 0 & 0 & 0 \\ c & d & 0 & 0 & 0 & 0 \\ b & c & d & 0 & 0 & 0 \\ a & b & c & d & 0 & 0 \\ 0 & a & b & c & d & 0 \\ 0 & 0 & a & b & c & d \\ 0 & 0 & 0 & a & b & c \\ 0 & 0 & 0 & 0 & a & b \\ 0 & 0 & 0 & 0 & 0 & a \end{bmatrix} \begin{bmatrix} x \\ y \\ z \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} dx \\ cx + dy \\ bx + cy + dz \\ ax + by + cz \\ ay + bz \\ az \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

Q: Difference between convolution and cross-correlation?



Semantic segmentation

Transpose convolution as matrix multiplication

$$\begin{bmatrix} d & c & b & a & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & d & c & b & a & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & d & c & b & a & 0 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 0 \\ x \\ y \\ z \\ 0 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} ax \\ cx + by + az \\ dy + cz \end{bmatrix}$$

Strided convolution

$p * q = Qp_{00}$

Stride 2

Q: How is this different from a convolution?

No convolution (cross-correlation) anymore!!

$$\begin{bmatrix} d & 0 & 0 \\ c & 0 & 0 \\ b & d & 0 \\ a & c & 0 \\ 0 & b & d \\ 0 & a & c \\ 0 & 0 & b \\ 0 & 0 & a \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \end{bmatrix} = \begin{bmatrix} dx \\ cy \\ bx + dz \\ ax + cz \\ by + dz \\ ay + cz \\ bz \\ az \\ 0 \end{bmatrix}$$

Strided transpose convolution

$p *^T q = Q^T p_{00}$

Stride 2

27 5LSM0 Module 8: Supervised learning

Semantic segmentation

Fully convolutional network (FCN)

➤ Pooling

➤ Strided convolution

➤ Unpooling

➤ Strided transpose convolution

Q: Fundamental difference between unpooling and transpose convolution?

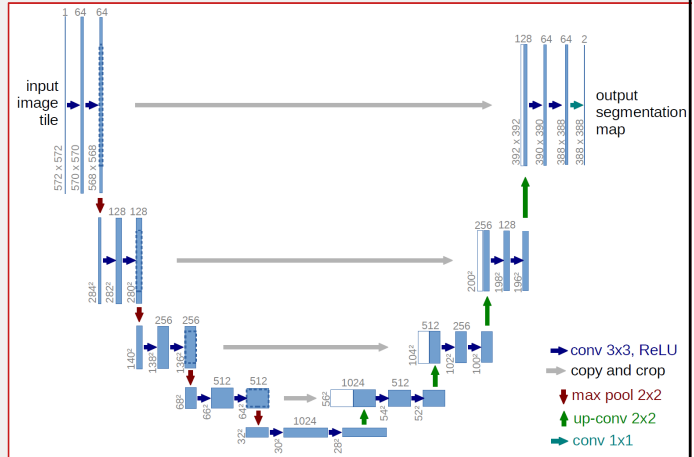
28 5LSM0 Module 8: Supervised learning

Semantic segmentation

The U-net architecture

- Problem with conventional FCNs:
 - *Sharp edge information is lost*
- Solution:
 - *Append upstream convolution maps with downstream convolution maps of the same scale / size*
 - *Fine-grained edge information preserved*

Ronneberger et al., MICCAI 2015
 -> Most cited MICCAI paper, runner up from 1998...
 (4,568 vs 3,310 @ March 10, 2019)



Classification + localization

Image classification

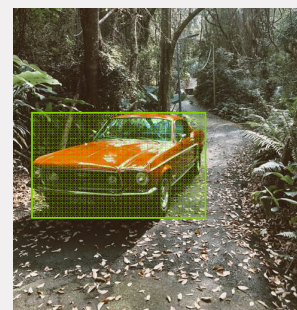
- Category label assigent to the image

Classification + localization

- Predict category label
- New: Where is the <category> in that image?

Example:

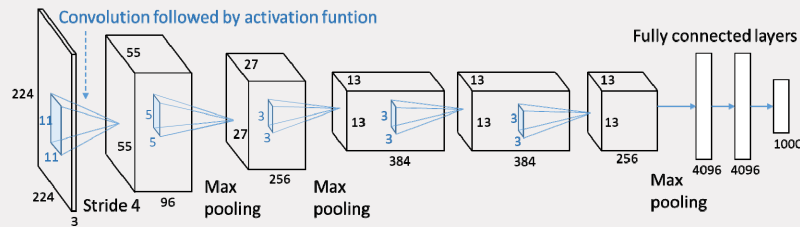
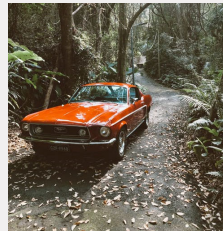
- Image labeled as CAR
- Draw a bounding box around the car within the image



Car

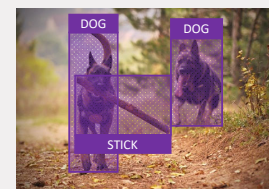


Classification + Localization



Object detection

Fixed set of categories that we're interested in

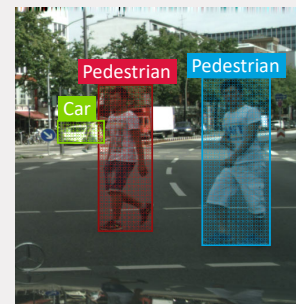


Every time that one of those categories appears in an input image:

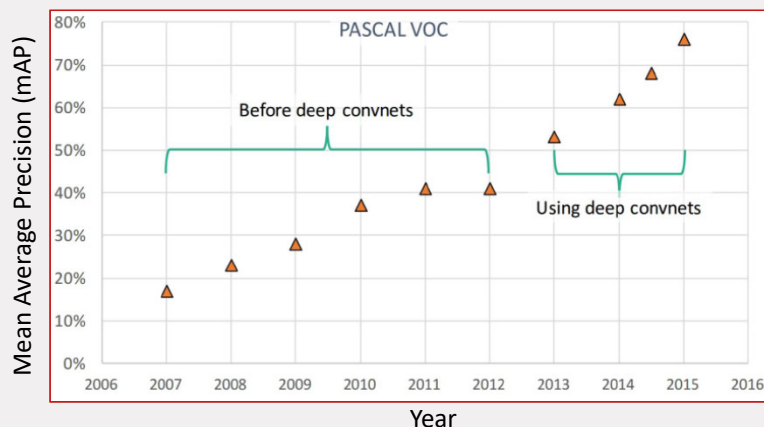
- put a bounding box around it and classify it.

Why can't we use the classification + localization approach from the previous slides?

- Problem: we don't know how many objects of interest to expect in the image
- Hence, we don't know how many bounding boxes to estimate and classify...



Object detection

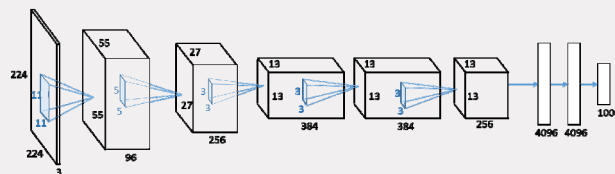
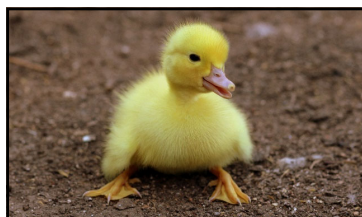


Current state-of-the-art well over 80%

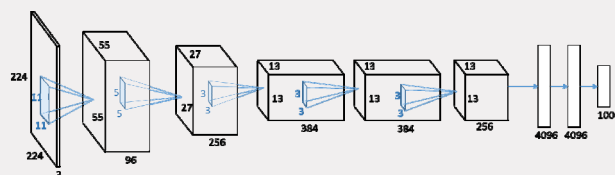


33 5LSM0 Module 8: Supervised learning

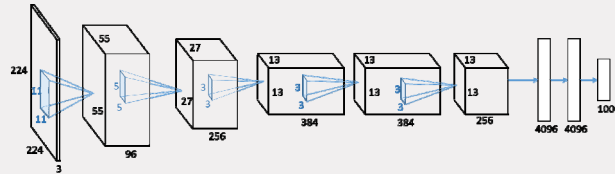
*Example from [cs231n \(2017\)](#): lecture 11 – slide 54



Duckling (x,y,w,h)



Cat (x,y,w,h)
 Cat (x,y,w,h)
 Cat (x,y,w,h)
 Cat (x,y,w,h)
 Cat (x,y,w,h)



Fish (x,y,w,h)
 Fish (x,y,w,h)
 Fish (x,y,w,h)
 ...
 Fish (x,y,w,h)



34 5LSM0 Module 8: Supervised learning



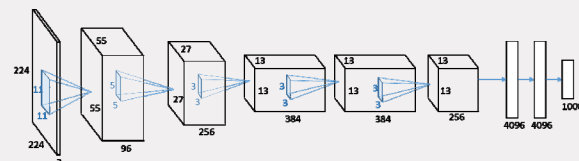
Object detection

Other approach: sliding window



Move a window over the image and classify each window independently:

- + Add category for background

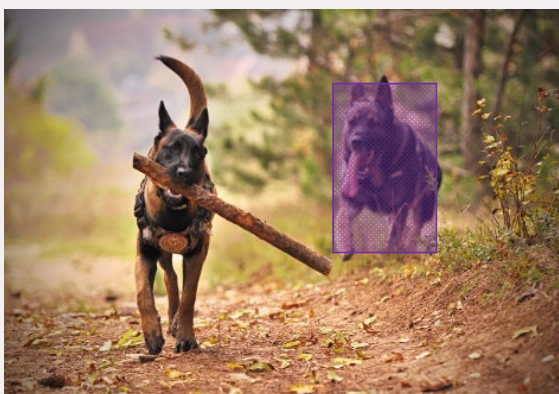


Cat? NO
 Dog? NO
 Stick? NO
 Background? YES



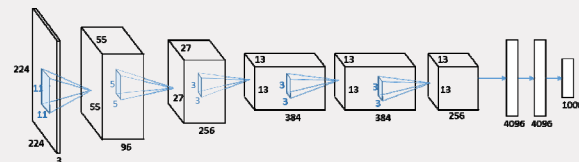
Object detection

Other approach: sliding window



Move a window over the image and classify each window independently:

- + Add category for background

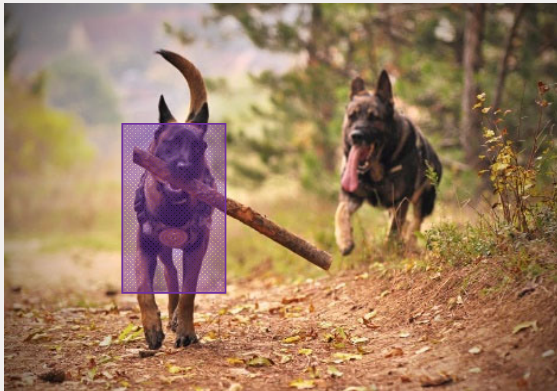


Cat? NO
 Dog? YES
 Stick? NO
 Background? NO



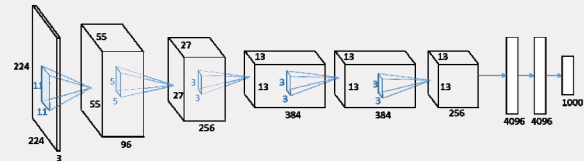
Object detection

Other approach: sliding window



Move a window over the image and classify each window independently:

- + Add category for background



Cat? NO
 Dog? YES
 Stick? MAYBE?
 Background? NO

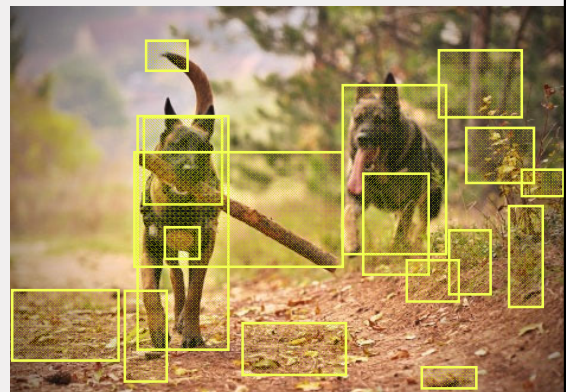
Q: Problem(s) with this approach?



Object detection

Region proposals

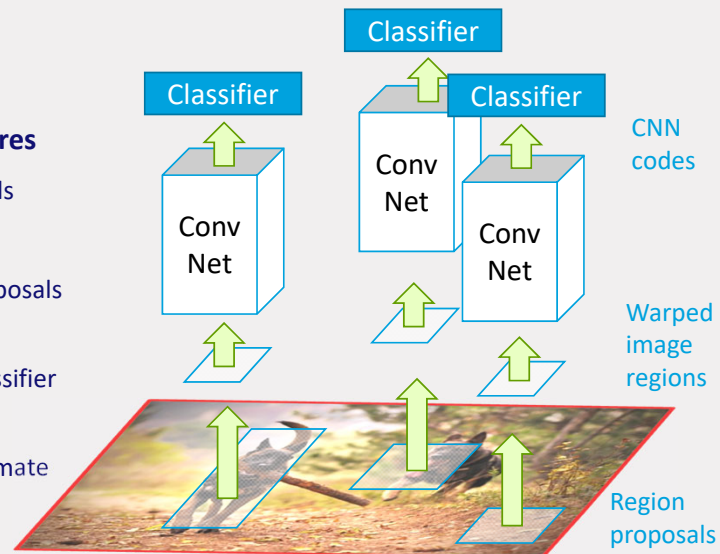
- Use conventional (fast) image processing techniques to find potential locations of objects
 - Regions in which an object may be present
 - Candidate proposal regions
- Typically a lot of proposals are generated
 - The majority of them may not actually fully contain an object
 - Very high recall → if there are objects in the image, some proposals will contain them



Object detection

R-CNN: Regions with CNN features

- Find regions using region proposals
- Warp images to normalize size
- Compute features of warped proposals with pre-trained CNN
- Feed features to conventional classifier (e.g. SVM) to classify the region
- In addition, use regression to estimate a correction for the bounding box



39 5LSM0 Module 8: Supervised learning

Girschik et al., "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014

TU/e

Object detection

Problems with R-CNN

- Training is very slow and memory heavy...
 - *Features of region candidates typically dumped to disk*
- Inference is slow...
 - *A lot of candidate regions which have to be classified*
- Different types of losses for training
 - *E.g. SVM loss for classification, Least squares regression loss, fine-tune network with softmax*
- Convolution results are not shared/re-used
- Region proposals are not learned...



40 5LSM0 Module 8: Supervised learning

TU/e

Object detection

Fast R-CNN!

- First run entire image through some convolutional layer(s).
- Region proposals for original image.
- Rather than taking crops from the image, project the proposals onto the feature map and take crops from there.

Q: Why would this strategy fail for the FC layers?

Note: now you can backprop through this!

Softmax classifier

Log loss + Smooth L1 loss

Linear + softmax

Linear

Bounding box regressors

FCs

ROI pooling layer

Region proposals on conv5 feature map

Conv Net

YCA

41 5LSM0 Module 8: Supervised learning

TU/e

Object detection

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

Fast R-CNN: how fast?

Training time (Hours)

Model	Training time (Hours)
R-CNN	84
SPP-Net	25.5
Fast R-CNN	8.75

Test time (seconds)

Model	Including Region proposals	Excluding Region Proposals
R-CNN	49	47
SPP-Net	4.3	2.3
Fast R-CNN	2.3	0.32

*Example from [cs231n \(2017\)](#): lecture 11 – slide 79

YCA

42 5LSM0 Module 8: Supervised learning

TU/e

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

Object detection

Faster R-CNN!!

- Bottleneck of Fast R-CNN: computing region proposals with external function
- Solution: insert Region Proposal Network (RPN) to come up with region proposals
- Jointly train with four losses
 - *RPN predict object? [Y/N]*
 - *RPN box coordinates?*
 - *Final class scores?*
 - *Final box coordinates?*

43 SLSM0 Module 8: Supervised learning

TU/e

Object detection

Faster R-CNN: how fast?

R-CNN Test-Time Speed

Model	Test-Time Speed
R-CNN	49
SPP-Net	4.3
Fast R-CNN	2.3
Faster R-CNN	0.2

44 SLSM0 Module 8: Supervised learning

*Example from [cs231n \(2017\)](#): lecture 11 – slide 82

TU/e

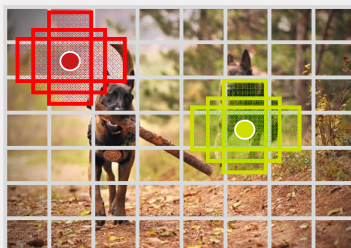
Object detection

So far: region based methods

- First propose a set of regions, then process each region independently

Alternative: all feed-forward in a single pass

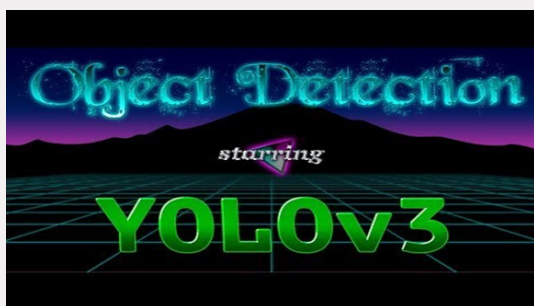
- YOLO: You Only Look Once / SSD: Single Shot Detection



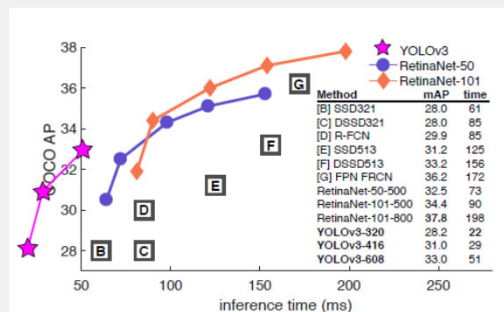
- Divide image into grid
- Use a base set of bounding boxes for each cell in the grid
- Solve one regression problem:
- For each bounding box predict:
 - True location + confidence
 - Class scores



Object detection



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

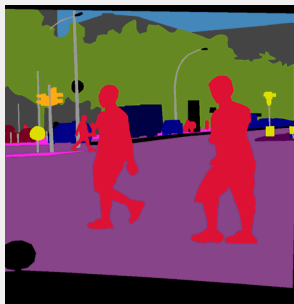


Redmon and Farhadi, "Yolov3: An incremental improvement." *arXiv preprint arXiv:1804.02767* (2018).



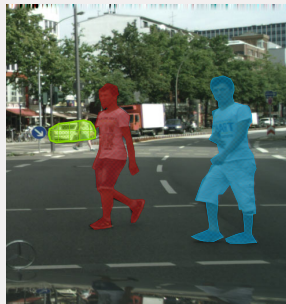
Instance segmentation

Segmentic segmentation



Pedestrian

Instance segmentation

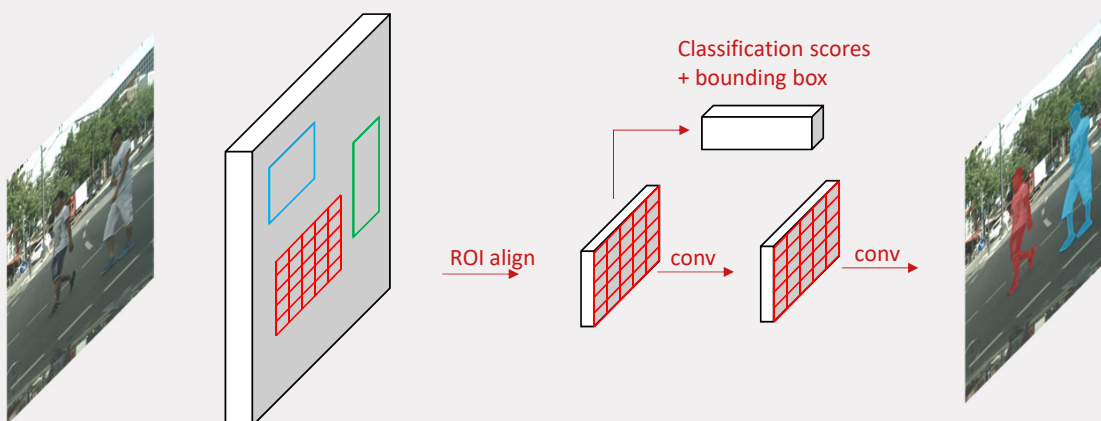


Pedestrian 1

Pedestrian 2



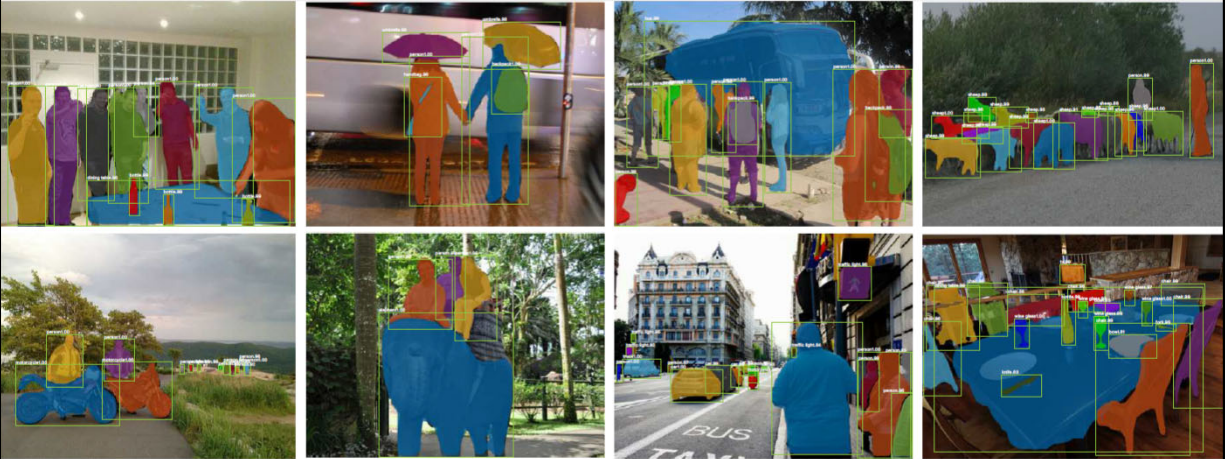
Instance segmentation



He, Kaiming, et al. "Mask r-cnn." ICCV 2017.



Instance segmentation



49 5LSM0 Module 8: Supervised learning

Kaiming He, Georgia Gkioxari, Piotr Dollar, Ross Girshick, "Mask R-CNN", The IEEE International Conference on Computer Vision (ICCV), 2017, pp. 2961-2969 ([link](#))

TU/e

Summary

Supervised learning

- Classification
 - *Predict category of an image*
- Semantic segmentation
 - *Predict category of pixels within an image*
- Classification + localization
 - *Predict category of an image + bounding box of the location of the object within the image*
- Object detection
 - *Detect multiple instances of different categories within a single image*
- Instance segmentation
 - *Detect and segment multiple instances of different categories within a single image*



50 5LSM0 Module 8: Supervised learning

TU/e

Next time

Wait a minute.. what if we ain't got no labels?

UNSUPERVISED LEARNING

