

Image Segmentation 2018

with Mask R-CNN

Gesina Schwalbe

7th May 2018

- 1 Definition and Goals
- 2 Architecture Components
- 3 Model
- 4 Training

Section 1

Definition and Goals

- 1 Definition and Goals
 - Problem

General Goals

- bounding boxes
- classification of each box
- pixel-mask for each box

Datasets

General Goals

- bounding boxes
- classification of each box
- pixel-mask for each box

Datasets

General Goals

- bounding boxes
- classification of each box
- pixel-mask for each box

Datasets

- COCO
- ImageNet

General Goals

- bounding boxes
- classification of each box
- pixel-mask for each box

Datasets

- COCO
- ImageNet

General Goals

- bounding boxes
- classification of each box
- pixel-mask for each box

Datasets

- COCO
- ImageNet

General Goals

- bounding boxes
- classification of each box
- pixel-mask for each box

Datasets

- COCO
- ImageNet

Applications

- live detection of signs, obstacles in traffic (example video)

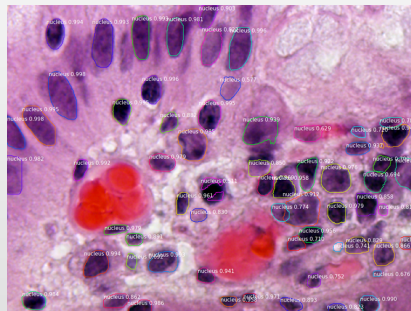
Applications

- automatic street map enhancement



Applications

- automatic evaluation of microscopy images



Architectual Requirements

- Master natural images: large nets
- Learn and process fast/efficiently
(currently best: 32h wt. 8GPU learning; 5fps inference):
 - share many features (FPN, RPN)
 - parallelize tasks/components

Architectual Requirements

- Master natural images: large nets
- Learn and process fast/efficiently
(currently best: 32h wt. 8GPU learning; 5fps inference):
 - special network components
(CNN, ResNe(X)t, FPN)
 - share many features (FPN, RPN)
 - parallelize tasks/components

Architectual Requirements

- Master natural images: large nets
- Learn and process fast/efficiently
(currently best: 32h wt. 8GPU learning; 5fps inference):
 - special network components
(CNN, ResNe(X)t, FPN)
 - share many features (FPN, RPN)
 - parallelize tasks/components

Architectual Requirements

- Master natural images: large nets
- Learn and process fast/efficiently
(currently best: 32h wt. 8GPU learning; 5fps inference):
 - special network components
(CNN, ResNe(X)t, FPN)
 - share many features (FPN, RPN)
 - parallelize tasks/components

Architectual Requirements

- Master natural images: large nets
- Learn and process fast/efficiently
(currently best: 32h wt. 8GPU learning; 5fps inference):
 - special network components
(CNN, ResNe(X)t, FPN)
 - share many features (FPN, RPN)
 - parallelize tasks/components

Section 2

Architecture Components

- 2 Architecture Components
 - Convolutional Networks
 - Deep CNNs: ResNet and ResNeXt

What

- Feedforward neural network with only local connections
- Massive weight sharing
- (often:) Downscaling of feature space:
 - Translation invariance
 - Multiple scale feature spaces available

What

- Feedforward neural network with only local connections
- Massive weight sharing
- (often:) Downscaling of feature space:
 - Multiple scale feature spaces available

What

- Feedforward neural network with only local connections
- Massive weight sharing
- (often:) Downscaling of feature space:
 - Translation invariance
 - Multiple scale feature spaces available

What

- Feedforward neural network with only local connections
- Massive weight sharing
- (often:) Downscaling of feature space:
 - Translation invariance
 - Multiple scale feature spaces available

What

- Feedforward neural network with only local connections
- Massive weight sharing
- (often:) Downscaling of feature space:
 - Translation invariance
 - Multiple scale feature spaces available

Convolution Operation

linear map $\mathbb{R}^{n_1 \times n_2 \times \dots \times n_3} \rightarrow \mathbb{R}^{n_1 \times n_2 \times \dots \times n_3}$ described by

- a fixed sliding window shape
- a window of that shape with a weight value at each coordinate, called the *kernel*.

(Animation)

Convolution Operation

linear map $\mathbb{R}^{n_1 \times n_2 \times \dots \times n_3} \rightarrow \mathbb{R}^{n_1 \times n_2 \times \dots \times n_3}$ described by

- a fixed sliding window shape
- a window of that shape with a weight value at each coordinate, called the *kernel*.

(Animation)

Convolutional Layer

Convolution Hyperparameters:

- size of kernel (= weight-window)
- padding variant (= border treatment)
- number of filters (= number of parallel convolutions)
- stride (= downsampling rate)
- dilation (=upsampling rate)
- activation function

Convolutional Layer

Convolution Hyperparameters:

- size of kernel (= weight-window)
- padding variant (= border treatment)
- number of filters (= number of parallel convolutions)
- stride (= downsampling rate)
- dilation (=upsampling rate)
- activation function

Convolutional Layer

Convolution Hyperparameters:

- size of kernel (= weight-window)
- padding variant (= border treatment)
- number of filters (= number of parallel convolutions)
- stride (= downsampling rate)
- dilation (=upsampling rate)
- activation function

Convolutional Layer

Convolution Hyperparameters:

- size of kernel (= weight-window)
- padding variant (= border treatment)
- number of filters (= number of parallel convolutions)
- stride (= downsampling rate)

- initialization and normalization
- activation function

Convolutional Layer

Convolution Hyperparameters:

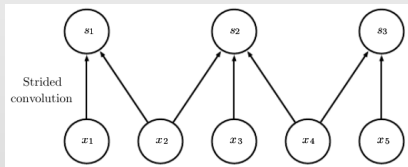
- size of kernel (= weight-window)
- padding variant (= border treatment)
- number of filters (= number of parallel convolutions)
- stride (= downsampling rate)

- dilation (=upsampling rate)
- activation function

Convolutional Layer

Convolution Hyperparameters:

- size of kernel (= weight-window)
- padding variant (= border treatment)
- number of filters (= number of parallel convolutions)
- stride (= downsampling rate)

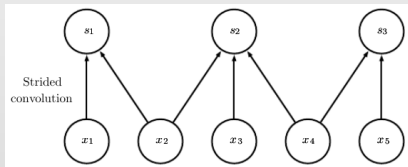


- dilation (=upsampling rate)
- activation function

Convolutional Layer

Convolution Hyperparameters:

- size of kernel (= weight-window)
- padding variant (= border treatment)
- number of filters (= number of parallel convolutions)
- stride (= downsampling rate)

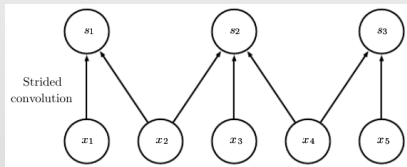


- dilation (=upsampling rate)
- activation function

Convolutional Layer

Convolution Hyperparameters:

- size of kernel (= weight-window)
- padding variant (= border treatment)
- number of filters (= number of parallel convolutions)
- stride (= downsampling rate)



- dilation (=upsampling rate)
- activation function

Convolutional Layer

Weights: kernel values, bias

Pooling Layer

Usual pooling functions

- Average Pooling (linear)
- MaxPooling (non-linear)

Pooling Layer

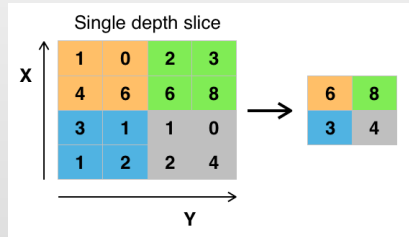
Usual pooling functions

- Average Pooling (linear)
- MaxPooling (non-linear)

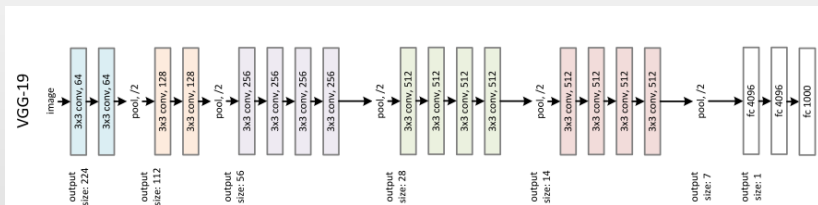
Pooling Layer

Usual pooling functions

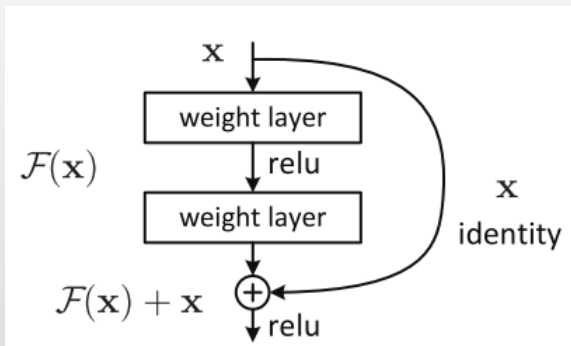
- Average Pooling (linear)
- MaxPooling (non-linear)



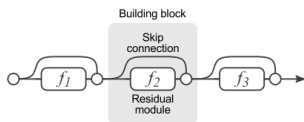
(Deep) CNN



Residual CNNs: ResNet

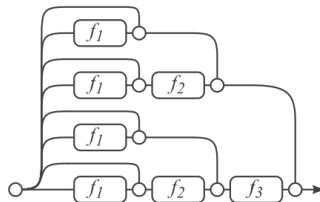


Residual CNNs: ResNet

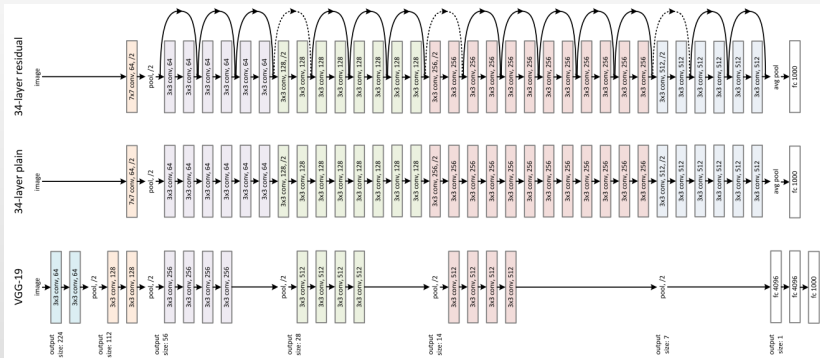


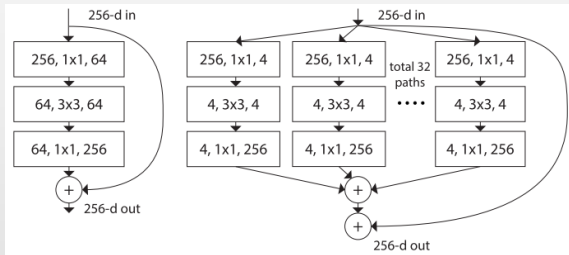
(a) Conventional 3-block residual network

=

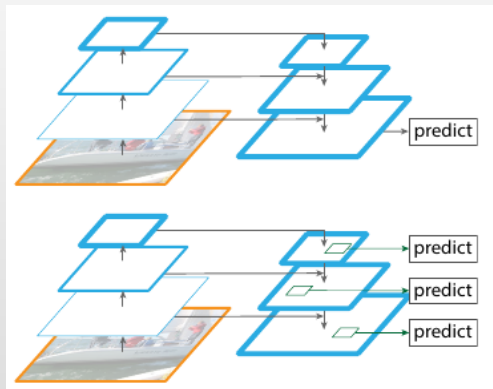


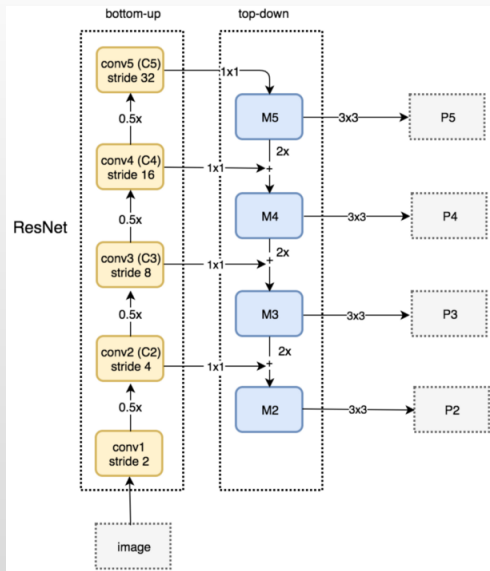
(b) Unraveled view of (a)





More Feature sharing: FPNs





Section 3

Model

3 Model

- Overview
- Convolutional Backbone
- Region Proposal Network
- RoI-Pooling
- Frontend

Predecessor Problems

Object Classification one object per image to one class per image

Object Detection

• find all bounding boxes in the image

- identify each window
- maybe enhance the window selection

Predecessor Problems

Object Classification one object per image to one class per image

Object Detection

- (intelligently) choose windows of the image
- classify each window
- maybe enhance the window selection

Predecessor Problems

Object Classification one object per image to one class per image

Object Detection

- (intelligently) choose windows of the image
- classify each window
- maybe enhance the window selection

Predecessor Problems

Object Classification one object per image to one class per image

Object Detection

- (intelligently) choose windows of the image
- classify each window
- maybe enhance the window selection

Predecessor Problems

Object Classification one object per image to one class per image

Object Detection

- (intelligently) choose windows of the image
- classify each window
- maybe enhance the window selection

Components

- 1 Convolutional backbone: extract important features
- 2 Region proposal; *in parallel*:
 - Window objectness scoring
- 3 Region of Interest (RoI) Pooling
 - Classification
 - Masking
 - Bounding Box Optimization

Components

- ① Convolutional backbone: extract important features
- ② Region proposal; *in parallel*:
 - Window objectness scoring
 - Window correction
- ③ Frontend; *in parallel*:
 - Classification
 - Masking
 - Bounding Box Optimization

Components

- ① Convolutional backbone: extract important features
- ② Region proposal; *in parallel*:
 - Window objectness scoring
 - Window correction
- ③ Frontend; *in parallel*:
 - Classification
 - Masking
 - Bounding Box Optimization

Components

- ① Convolutional backbone: extract important features
- ② Region proposal; *in parallel*:
 - Window objectness scoring
 - Window correction
- ③ Frontend; *in parallel*:
 - Masking
 - Masking
 - Bounding Box Optimization

Components

- ① Convolutional backbone: extract important features
- ② Region proposal; *in parallel*:
 - Window objectness scoring
 - Window correction
- ③ Frontend; *in parallel*:
 - Classification
 - Masking
 - Bounding Box Optimization

Components

- ① Convolutional backbone: extract important features
- ② Region proposal; *in parallel*:
 - Window objectness scoring
 - Window correction
- ③ Frontend; *in parallel*:
 - Classification
 - Masking
 - Bounding Box Optimization

Components

- ① Convolutional backbone: extract important features
- ② Region proposal; *in parallel*:
 - Window objectness scoring
 - Window correction
- ③ Frontend; *in parallel*:
 - Classification
 - Masking
 - Bounding Box Optimization

Components

- ① Convolutional backbone: extract important features
- ② Region proposal; *in parallel*:
 - Window objectness scoring
 - Window correction
- ③ Frontend; *in parallel*:
 - Classification
 - Masking
 - Bounding Box Optimization

Convolutional Backbone

Goal extract features

Architecture same as for Object Detection

Predecessors/Alternatives

Pixel Merging (e.g. SelectiveSearch)

Window scoring (e.g. Objectness)

Separate NN (e.g. Multibox)

Predecessors/Alternatives

Pixel Merging (e.g. SelectiveSearch)

Window scoring (e.g. Objectness)

Separate NN (e.g. Multibox)

Predecessors/Alternatives

Pixel Merging (e.g. SelectiveSearch)

Window scoring (e.g. Objectness)

Separate NN (e.g. Multibox)

Main Ideas of the Mask-RCNN Approach

Window scoring with enhancements:

Decoupling of classification and window proposals

All Scales at once using different candidate window shapes

Bounding box correction in *parallel* to scoring

Excessive weight sharing amongst same shapes

RoI-Pooling — Feature sharing
— one convolutional layer from backbone and pool each proposal window to fixed size

Main Ideas of the Mask-RCNN Approach

Window scoring with enhancements:

Decoupling of classification and window proposals

All Scales at once using different candidate window shapes

Bounding box correction in *parallel* to scoring

Excessive weight sharing amongst same shapes

RoI-Pooling = Feature sharing, i.e. use convolutional features from backbone and pool each proposal window to fixed size

Main Ideas of the Mask-RCNN Approach

Window scoring with enhancements:

Decoupling of classification and window proposals

All Scales at once using different candidate window shapes

Bounding box correction in *parallel* to scoring

Excessive weight sharing amongst same shapes

RoI-Pooling = Feature sharing i.e. use convolutional features from backbone and pool each proposal window to fixed size

Main Ideas of the Mask-RCNN Approach

Window scoring with enhancements:

Decoupling of classification and window proposals

All Scales at once using different candidate window shapes

Bounding box correction in *parallel* to scoring

Excessive weight sharing amongst same shapes

RoI-Pooling = Feature sharing i.e. use convolutional features from backbone and pool each proposal window to fixed size

Main Ideas of the Mask-RCNN Approach

Window scoring with enhancements:

Decoupling of classification and window proposals

All Scales at once using different candidate window shapes

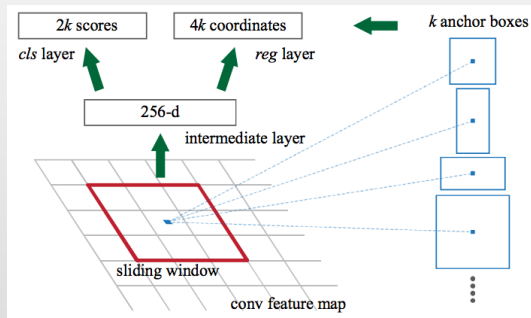
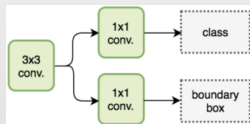
Bounding box correction in *parallel* to scoring

Excessive weight sharing amongst same shapes

RoI-Pooling = Feature sharing i.e. use convolutional features from backbone and pool each proposal window to fixed size

Architecture Overview

Shared Conv Layer Sliding window on feature space with fixed set of anchor shapes per window



Architecture Overview

Shared Conv Layer Sliding window on feature space with fixed set of anchor shapes per window

cls, reg Per window and anchor shape do in parallel

Architecture Overview

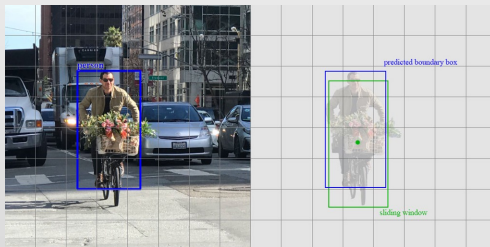
Shared Conv Layer Sliding window on feature space with fixed set of anchor shapes per window

cls, reg Per window and anchor shape do in parallel
objectness score (*cls*)

Architecture Overview

Shared Conv Layer Sliding window on feature space with fixed set of anchor shapes per window

cls, reg Per window and anchor shape do in parallel
objectness score (*cls*)
coordinate correction (*reg*)

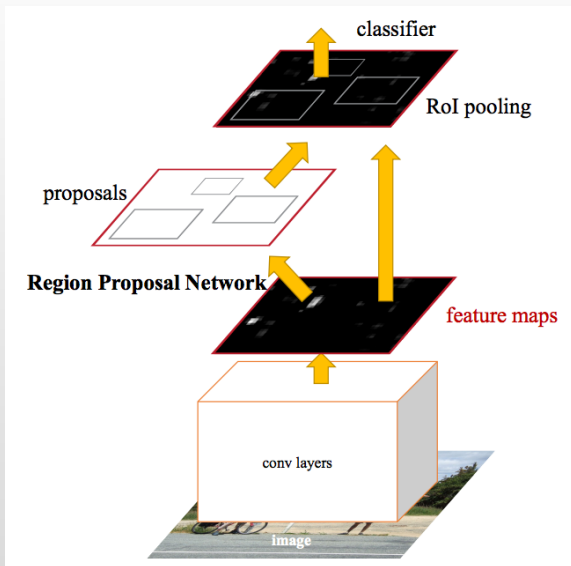


Architecture Overview

Shared Conv Layer Sliding window on feature space with fixed set of anchor shapes per window

cls, reg Per window and anchor shape do in parallel
objectness score (*cls*)
coordinate correction (*reg*)

Proposal Layer Select best proposals



Coordinate and RPN *reg* output encoding

Box coordinate encoding (all normalized):

(x_1, y_1) # upper left corner

(x_2, y_2) # lower right corner

Coordinate and RPN *reg* output encoding

Coordinate correction (= *reg* output) encoding:

$$\left(dx, \# \frac{\text{centerpoint x-difference}}{\text{anchor width}}\right)$$

$$dy, \# \frac{\text{centerpoint y-difference}}{\text{anchor height}}$$

$$\log(dw), \# \frac{\text{ground-truth box width}}{\text{anchor width}}$$

$$\log(dh) \# \frac{\text{ground-truth box height}}{\text{anchor height}}$$

Metrics

Metric for matching:

$$\text{IoU}(A, B) := \frac{\text{Intersection Area}}{\text{Union Area}}$$

Labels

object=1 with ground-truth box b if

- best **IoU** for b , else if
- **IoU** with $b \geq 0.3$

no object \Rightarrow 0 if not pos. and **IoU** ≤ 0.3

negative IoU

Balance positive and negative boxes for training!

Labels

object=1 with ground-truth box b if

- best **IoU** for b , else if
- **IoU** with $b \geq 0.3$

no **object=-1** if not pos. and **IoU** ≤ 0.3

neutral=0 else (excluded from training)

Balance positive and negative boxes for training!

Labels

object=1 with ground-truth box b if

- best **IoU** for b , else if
- **IoU** with $b \geq 0.3$

no **object=-1** if not pos. and **IoU** ≤ 0.3

neutral=0 else (excluded from training)

Balance positive and negative boxes for training!

Labels

object=1 with ground-truth box b if

- best **IoU** for b , else if
- **IoU** with $b \geq 0.3$

no object=-1 if not pos. and **IoU** ≤ 0.3

neutral=0 else (excluded from training)

Balance positive and negative boxes for training!

Labels

object=1 with ground-truth box b if

- best **IoU** for b , else if
- **IoU** with $b \geq 0.3$

no object=-1 if not pos. and **IoU** ≤ 0.3

neutral=0 else (excluded from training)

Balance positive and negative boxes for training!

Architecture Details

Sliding window Shared Conv layer with “valid”-padding

Architecture Details

Objectness classification Conv layer with

- 1×1 -sized kernel
- 2-class softmax activation: (non-object score, object score)

Loss: crossentropy for non-neutral anchors

Architecture Details

Coordinate correction Conv layer with

- 1×1 -sized kernel
- $4 \times$ (number of anchors) filters: (dx, dy, dw, dh)
coordinate correction for each anchor

Loss: smooth L_1 -loss

for positive anchors that do not cross image bounds

Architecture Details

Proposal selection

- 1 Trim to N best-object-scored anchors.
- 2 Apply coordinate correction.
- 3 Clip boxes.
- 4 Non-maximum suppression:
 - Sort anchors
 - Reject boxes which have high IoU with better scored ones
 - Trim to best *num_proposals*

Architecture Details

Proposal selection

- 1 Trim to N best-object-scored anchors.
- 2 Apply coordinate correction.
- 3 Clip boxes.
- 4 Non-maximum suppression:
 - Sort anchors
 - Reject boxes which have high IoU with better scored ones
 - Trim to best *num_proposals*

Architecture Details

Proposal selection

- 1 Trim to N best-object-scored anchors.
- 2 Apply coordinate correction.
- 3 Clip boxes.
- 4 Non-maximum suppression:
 - Remove boxes which have high IoU with better scored ones
 - Trim to best *num_proposals*

Architecture Details

Proposal selection

- 1 Trim to N best-object-scored anchors.
- 2 Apply coordinate correction.
- 3 Clip boxes.
- 4 Non-maximum suppression:
 - Sort by object score
 - Iterate over object-scored ones
 - Trim to best *num_proposals*

Architecture Details

Proposal selection

- 1 Trim to N best-object-scored anchors.
- 2 Apply coordinate correction.
- 3 Clip boxes.
- 4 Non-maximum suppression:
 - 1 Sort by score
 - 2 Reject boxes which have high IoU with better scored ones
 - 3 Trim to best *num_proposals*

Architecture Details

Proposal selection

- ① Trim to N best-object-scored anchors.
- ② Apply coordinate correction.
- ③ Clip boxes.
- ④ Non-maximum suppression:
 - ① Sort by score
 - ② Reject boxes which have high IoU with better scored ones
 - ③ Trim to best *num_proposals*

Architecture Details

Proposal selection

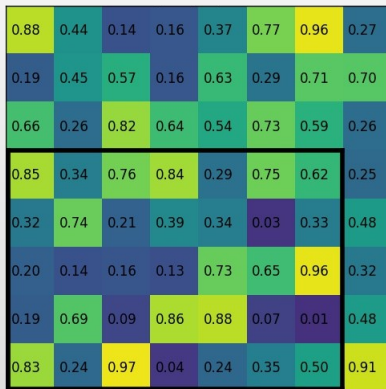
- ① Trim to N best-object-scored anchors.
- ② Apply coordinate correction.
- ③ Clip boxes.
- ④ Non-maximum suppression:
 - ① Sort by score
 - ② Reject boxes which have high IoU with better scored ones
 - ③ Trim to best *num_proposals*

Architecture Details

Proposal selection

- ① Trim to N best-object-scored anchors.
- ② Apply coordinate correction.
- ③ Clip boxes.
- ④ Non-maximum suppression:
 - ① Sort by score
 - ② Reject boxes which have high IoU with better scored ones
 - ③ Trim to best *num_proposals*

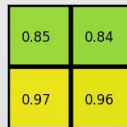
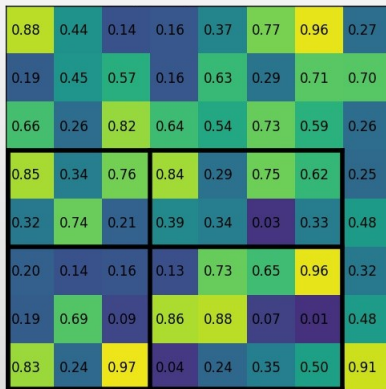
Rol-Pooling



Rol-Pooling

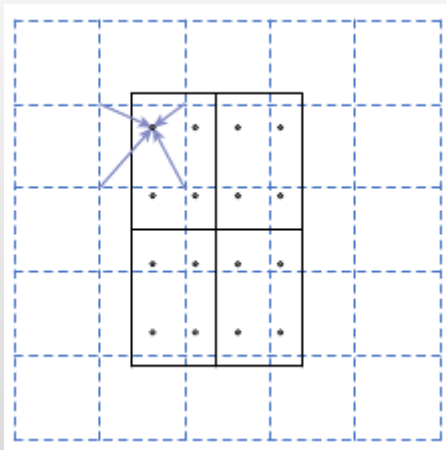
0.88	0.44	0.14	0.16	0.37	0.77	0.96	0.27
0.19	0.45	0.57	0.16	0.63	0.29	0.71	0.70
0.66	0.26	0.82	0.64	0.54	0.73	0.59	0.26
0.85	0.34	0.76	0.84	0.29	0.75	0.62	0.25
0.32	0.74	0.21	0.39	0.34	0.03	0.33	0.48
0.20	0.14	0.16	0.13	0.73	0.65	0.96	0.32
0.19	0.69	0.09	0.86	0.88	0.07	0.01	0.48
0.83	0.24	0.97	0.04	0.24	0.35	0.50	0.91

Rol-Pooling



Rol-Align

Bilinear interpolation instead of cropping:



Main Ideas

Decouple

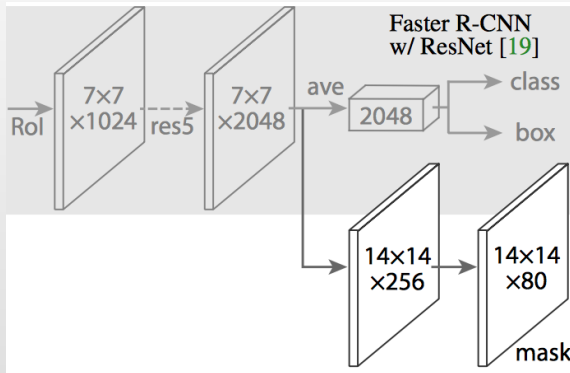
- classification, bounding box optimization, and masks;
- mask predictions for the different classes

Main Ideas

Decouple

- classification, bounding box optimization, and masks;
- mask predictions for the different classes

Architecture



Architecture

Classification fully connected layers ending in softmax (include class “No object”)

Loss: multinomial crossentropy

Architecture

Mask generation

- 1 Few (1–3) Conv Layers, maybe with upscaling parts
- 2 Conv Layer with
 - a filter for each class
 - sigmoid activation

Loss: binary cross-entropy

Architecture

Bounding box regression Linear regression

Section 4

Training

4 Training

- Overview
- Backbone Pretraining
- Alternating Training

Steps

- 1 Backbone
- 2 RPN
- 3 Alternating training of RPN and Frontend

Steps

- 1 Backbone
- 2 RPN
- 3 Alternating training of RPN and Frontend

Steps

- 1 Backbone
- 2 RPN
- 3 Alternating training of RPN and Frontend

Backbone Pretraining

Separate ConvNet Training Model

- 1 Backbone (Conv & Pooling)
- 2 Dense Layers
- 3 Classification Softmax-Layer

Training Input

Backbone Pretraining

Separate ConvNet Training Model

- 1 Backbone (Conv & Pooling)
- 2 Dense Layers
- 3 Classification Softmax-Layer

Training Input

Backbone Pretraining

Separate ConvNet Training Model

- 1 Backbone (Conv & Pooling)
- 2 Dense Layers
- 3 Classification Softmax-Layer

Training Input

Backbone Pretraining

Separate ConvNet Training Model

- 1 Backbone (Conv & Pooling)
- 2 Dense Layers
- 3 Classification Softmax-Layer

Training Input

inputs: one-object images (ca. size of later bounding boxes)
targets: the single objects' labels

Backbone Pretraining

Separate ConvNet Training Model

- 1 Backbone (Conv & Pooling)
- 2 Dense Layers
- 3 Classification Softmax-Layer

Training Input

inputs one-object images (ca. size of later bounding boxes)
labels the single objects' labels

Backbone Pretraining

Separate ConvNet Training Model

- 1 Backbone (Conv & Pooling)
- 2 Dense Layers
- 3 Classification Softmax-Layer

Training Input

inputs one-object images (ca. size of later bounding boxes)

labels the single objects' labels

Backbone Pretraining

Separate ConvNet Training Model

- 1 Backbone (Conv & Pooling)
- 2 Dense Layers
- 3 Classification Softmax-Layer

Training Input

inputs one-object images (ca. size of later bounding boxes)

labels the single objects' labels

Alternating Training

(after RPN is trained)

- 1 Frontend: RPN fixed, backbone not shared
- 2 RPN: frontend fixed, shared backbone fixed
- 3 Frontend: RPN fixed, shared backbone fixed
- 4 ...

Alternating Training

(after RPN is trained)

- 1 Frontend: RPN fixed, backbone not shared
- 2 RPN: frontend fixed, shared backbone fixed
- 3 Frontend: RPN fixed, shared backbone fixed
- 4 ...

Alternating Training

(after RPN is trained)

- 1 Frontend: RPN fixed, backbone not shared
- 2 RPN: frontend fixed, shared backbone fixed
- 3 Frontend: RPN fixed, shared backbone fixed
- 4 ...

Alternating Training

(after RPN is trained)

- 1 Frontend: RPN fixed, backbone not shared
- 2 RPN: frontend fixed, shared backbone fixed
- 3 Frontend: RPN fixed, shared backbone fixed
- 4 ...

Section 5

Implementation Review

- Source Code
- Lessons learned

Example Sources

- Keras implementation by matterport: [1]
- Easy example for handwritten number detection using autogenerated data based on MNIST: github

Keras Implementation Specialties

- Use the functional API!
- Custom layers:
 - Loss Layers: custom losses
 - Region Proposal Network: select RPN proposals
- Separate models for training and inference

Keras Implementation Specialties

- Use the functional API!
- Custom layers:
 - Loss Layers custom losses
 - Proposal Layer select RPN proposals
 - RoI-Pooling Layer reshape proposals and mask labels (can use `tf.crop_and_resize()`)
- Separate models for training and inference

Keras Implementation Specialties

- Use the functional API!
- Custom layers:
 - Loss Layers** custom losses
 - Proposal Layer** select RPN proposals
 - RoI-Pooling Layer** reshape proposals and mask labels (can use `tf.crop_and_resize()`)
- Separate models for training and inference

Lessons Learned I

- Always double-check and note down tensor/array dimensions; mind the padding for convolutions.
- Always double-check and test your algorithms.
- Always have a look at *all* input and output data:
 - Do they roughly make sense, e.g. do positive classes have different probability output than negative ones?)
 - Are there error patterns?

Visualization is your friend. But it will contain bugs, too.

- Directly document and update all input and output formats of functions. Resp. in general: Keep your code VERY clean and understandable.

Lessons Learned II

- Convolution is very geometric: Depict your sliding windows and downscaling factor(s) to check whether they make sense with your data/object sizes.
- Mind your RAM when optimizing data generation/tagging
- Check your loss and validation values; Does the model actually get better?
- The deeper, the much more time;
- Make data as easy (small) as possible; if you can still classify, the network can do, too (e.g. grayscale instead of several color channels).