

THE UNIVERSITY OF BRITISH COLUMBIA

Topics in AI (CPSC 532L): **Multimodal Learning with Vision, Language and Sound**

Lecture 5: Convolutional Neural Networks (Part 2)



Course Logistics

Assignment 2 questions Assignment 1

Computational Graph

- Had to break functions to the elementary units
- Expectation is a single graph

Computational Graph

- Had to break functions to the elementary units
- Expectation is a single graph

Forward Pass

- Show work
- Had to do computations right

Computational Graph

- Had to break functions to the elementary units
- Expectation is a single graph

AutoDiff Forward Mode

- Had to have two passes (one for each **input** variable)
- Needed to show enough work, that I was fairly certain you understood the process

Forward Pass

- Show work
- Had to do computations right

Computational Graph

- Had to break functions to the elementary units
- Expectation is a single graph

AutoDiff Forward Mode

- Had to have two passes (one for each **input** variable)
- Needed to show enough work, that I was fairly certain you understood the process

AutoDiff Backwards Mode / BackProp

- Had to have two passes (one for each output variable)
- Needed to show enough work, that I was fairly certain you understood the process

Forward Pass

- Show work
- Had to do computations right



Categorization





Categorization



Multi-class: Horse Church Toothbrush Person **IM** GENET



Categorization



Horse Multi-class: Church Toothbrush Person **IM** GENET

Multi-label: Horse

Church Toothbrush Person



Categorization

Detection





Multi-**class:** Horse Church Toothbrush Person IM GENET

Multi-label: Horse

Church Toothbrush Person

Horse (x, y, w, h) Horse (x, y, w, h) Person (x, y, w, h) Person (x, y, w, h)





Categorization

Detection





Multi-**class:** Horse Church Toothbrush Person **IM** GENET

Multi-label: Horse

Church Toothbrush Person

Horse (x, y, w, h) Horse (x, y, w, h) Person (x, y, w, h) Person (x, y, w, h)





Segmentation

Horse Person



Categorization

Detection





Multi-class: Horse Church Toothbrush Person **IM** GENET

Multi-label: Horse

Church Toothbrush Person

Horse (x, y, w, h) Horse (x, y, w, h) Person (x, y, w, h) Person (x, y, w, h)





Segmentation

Horse Person



Instance Segmentation

Horse1 Horse₂ Person1 Person2



Object Classification



Problem: For each image predict which category it belongs to out of a fixed set





Object Classification





Category	Predictio
Dog	No
Cat	No
Couch	No
Flowers	No
Leopard	Yes

Problem: For each image predict which category it belongs to out of a fixed set







Object Classification







Problem: For each image predict which category it belongs to out of a fixed set

 \mathbf{x}^t

CNN Architectures: LeNet-5



Architecture: $CONV \longrightarrow POOL \longrightarrow CONV \longrightarrow POOL \longrightarrow FC \longrightarrow FC$ Conv filters: 5x5, Stride: 1 **Pooling:** 2x2, Stride: 2

[LeCun et al., 1998]

ImageNet Dataset

Over **14 million** (high resolution) web **images** Roughly labeled with **22K synset** categories Labeled on Amazon Mechanical Turk (AMT)



Popular Synsets

Animal

fish bird mammal invertebrate

Plant

tree flower vegetable

Activity

sport

Material

fabric

Instrumentation

utensil appliance tool musical instrument

Scene

room geological formation

Food

beverage



ImageNet Competition (ILSVRC)

Annual competition of image classification at scale Focuses on a subset of **1K synset** categories **Scoring:** need to predict true label within top K (K=5)





Architecture: CONV1 MAX POOL1 NORM1 CONV2 MAX POOL2 NORM2 CONV3 CONV4 CONV5 Max POOL3 FC6 FC7 FC8

Input: 227 x 227 x 3 images



[Krizhevsky et al., 2012]

Architecture: CONV1 MAX POOL1 NORM1 CONV2 MAX POOL2 NORM2 CONV3 CONV4 CONV5 Max POOL3 FC6 FC7 FC8

Input: 227 x **CONV1:** 96 1



[Krizhevsky et al., 2012]

Input: 227 x 227 x 3 images

CONV1: 96 11 x 11 filters applied at stride 4

Architecture: CONV1 MAX POOL1 NORM1 CONV2 MAX POOL2 NORM2 CONV3 CONV4 CONV5 Max POOL3 FC6 FC7 FC8

Output: 55 x 55 x 96



[Krizhevsky et al., 2012]

Input: 227 x 227 x 3 images

CONV1: 96 11 x 11 filters applied at stride 4

Architecture: CONV1 MAX POOL1 NORM1 CONV2 MAX POOL2 NORM2 CONV3 CONV4 CONV5 Max POOL3 FC6 FC7 FC8

Output: 55 x 55 x 96



[Krizhevsky et al., 2012]

Input: 227 x 227 x 3 images

CONV1: 96 11 x 11 filters applied at stride 4 **Parameters:** 35K

Architecture: CONV1 MAX POOL1 NORM1 CONV2 MAX POOL2 NORM2 CONV3 CONV4 CONV5 Max POOL3 FC6 FC7 FC8

Output: 55 x 55 x 96



[Krizhevsky et al., 2012]

Input: 227 x 227 x 3 images

CONV1: 96 11 x 11 filters applied at stride 4

Parameters: 35K

MAX POOL1: 96 11 x 11 filters applied at stride 4

Architecture: CONV1 MAX POOL1 NORM1 CONV2 MAX POOL2 NORM2 CONV3 CONV4 CONV5 Max POOL3 FC6 FC7 FC8

Output: 55 x 55 x 96



[Krizhevsky et al., 2012]

Input: 227 x 227 x 3 images

CONV1: 96 11 x 11 filters applied at stride 4

Parameters: 35K

MAX POOL1: 96 11 x 11 filters applied at stride 4 **Output:** 27 x 27 x 96

Architecture: CONV1 MAX POOL1 NORM1 CONV2 MAX POOL2 NORM2 CONV3 CONV4 CONV5 Max POOL3 FC6 FC7 FC8

Output: 55 x 55 x 96



[Krizhevsky et al., 2012]

Input: 227 x 227 x 3 images

CONV1: 96 11 x 11 filters applied at stride 4

Parameters: 35K

MAX POOL1: 96 11 x 11 filters applied at stride 4 **Output:** 27 x 27 x 96 **Parameters:** 0

Full (simplified) AlexNet architecture: [227x227x3] INPUT [55x55x96] CONV1: 96 11x11 filters at stride 4, pad 0 [27x27x96] MAX POOL1: 3x3 filters at stride 2 [27x27x96] NORM1: Normalization layer [27x27x256] CONV2: 256 5x5 filters at stride 1, pad 2 [13x13x256] MAX POOL2: 3x3 filters at stride 2 [13x13x256] NORM2: Normalization layer [13x13x384] CONV3: 384 3x3 filters at stride 1, pad 1 [13x13x384] CONV4: 384 3x3 filters at stride 1, pad 1 [13x13x256] CONV5: 256 3x3 filters at stride 1, pad 1 [6x6x256] MAX POOL3: 3x3 filters at stride 2 [4096] FC6: 4096 neurons [4096] FC7: 4096 neurons [1000] FC8: 1000 neurons (class scores)



Krizhevsky et al., 2012

Details / Comments

- First use of ReLU
- Used contrast normalization layers
- Heavy data augmentation
- Dropout of 0.5
- Batch size of 128
- SGD Momentum (0.9)
- Learning rate (1e-2) reduced by 10 manually when validation accuracy plateaus
- L2 weight decay
- -7 CNN ensample: 18.2% -> 15.4%



ILSVRC winner 2012



ZF Net

AlexNet with small modifications:

- CONV1 (11 x 11 stride 4) to (7 x 7 stride 2)
- CONV3 # of filters 384 -> 512
- CONV4 # of filters 384 -> 1024
- CONV5 # of filters 256 -> 512

[Zeiler and Fergus, 2013]

ILSVRC winner 2012

Trend:

- -smaller filters (3 x 3)
- -deeper network (16 or 19 vs. 8 in AlexNet)

[Simonyan and Zisserman, 2014]

			Solumax
			FC 1000
		Softmax	FC 4096
		FC 1000	FC 4096
		FC 4096	Pool
		FC 4096	3x3 conv, 512
		Pool	3x3 conv, 512
-)		3x3 conv, 512	3x3 conv, 512
-)		3x3 conv, 512	3x3 conv, 512
		3x3 conv, 512	Pool
		Pool	3x3 conv, 512
	Softmax	3x3 conv, 512	3x3 conv, 512
	FC 1000	3x3 conv, 512	3x3 conv, 512
	FC 4096	3x3 conv, 512	3x3 conv, 512
	FC 4096	Pool	Pool
	Pool	3x3 conv, 256	3x3 conv, 256
	3x3 conv, 256	3x3 conv, 256	3x3 conv, 256
	3x3 conv, 384	Pool	Pool
	Pool	3x3 conv, 128	3x3 conv, 128
	3x3 conv, 384	3x3 conv, 128	3x3 conv, 128
	Pool	Pool	Pool
	5x5 conv, 256	3x3 conv, 64	3x3 conv, 64
	11x11 conv, 96	3x3 conv, 64	3x3 conv, 64
	Input	Input	Input
	AlexNet	VGG16	VGG19

Trend:

-smaller filters (3 x 3)

-deeper network (16 or 19 vs. 8 in AlexNet)

Why?

[Simonyan and Zisserman, 2014]

			Solumax
			FC 1000
		Softmax	FC 4096
		FC 1000	FC 4096
		FC 4096	Pool
		FC 4096	3x3 conv, 512
		Pool	3x3 conv, 512
-)		3x3 conv, 512	3x3 conv, 512
-)		3x3 conv, 512	3x3 conv, 512
		3x3 conv, 512	Pool
		Pool	3x3 conv, 512
	Softmax	3x3 conv, 512	3x3 conv, 512
	FC 1000	3x3 conv, 512	3x3 conv, 512
	FC 4096	3x3 conv, 512	3x3 conv, 512
	FC 4096	Pool	Pool
	Pool	3x3 conv, 256	3x3 conv, 256
	3x3 conv, 256	3x3 conv, 256	3x3 conv, 256
	3x3 conv, 384	Pool	Pool
	Pool	3x3 conv, 128	3x3 conv, 128
	3x3 conv, 384	3x3 conv, 128	3x3 conv, 128
	Pool	Pool	Pool
	5x5 conv, 256	3x3 conv, 64	3x3 conv, 64
	11x11 conv, 96	3x3 conv, 64	3x3 conv, 64
	Input	Input	Input
	AlexNet	VGG16	VGG19

Trend:

- -smaller filters (3 x 3)
- -deeper network (16 or 19 vs. 8 in Alex

Why?

- receptive field of a 3 layer ConvNet with filter size = is the same as 1 layer ConvNet with filter 7x7 (at stride

- deeper = more non-linearity (so richer filters)
- fewer parameters

[Simonyan and Zisserman, 2014]

			EC 1000
		Softmax	FC 1000
		FC 1000	FC 4096
		FC 4096	Pool
		FC 4096	3x3 conv, 512
		Pool	3x3 conv, 512
xNet)		3x3 conv, 512	3x3 conv, 512
		3x3 conv, 512	3x3 conv, 512
		3x3 conv, 512	Pool
		Pool	3x3 conv, 512
	Softmax	3x3 conv, 512	3x3 conv, 512
	FC 1000	3x3 conv, 512	3x3 conv, 512
	FC 4096	3x3 conv, 512	3x3 conv, 512
	FC 4096	Pool	Pool
= 3x3	Pool	3x3 conv, 256	3x3 conv, 256
1)	3x3 conv, 256	3x3 conv, 256	3x3 conv, 256
• /	3x3 conv, 384	Pool	Pool
	Pool	3x3 conv, 128	3x3 conv, 128
	3x3 conv, 384	3x3 conv, 128	3x3 conv, 128
	Pool	Pool	Pool
	5x5 conv, 256	3x3 conv, 64	3x3 conv, 64
	11x11 conv, 96	3x3 conv, 64	3x3 conv, 64
	Input	Input	Input
	AlexNet	VGG16	VGG19

INPUT: [224x224x3] memory: 224*224*3=150K params: 0 CONV3-64: [224x224x64] memory: 224*224*64=3.2M params: (3*3*3)*64 = 1,728 CONV3-64: [224x224x64] memory: 224*224*64=3.2M params: (3*3*64)*64 = 36,864 POOL2: [112x112x64] memory: 112*112*64=800K params: 0 CONV3-128: [112x112x128] memory: 112*112*128=1.6M params: (3*3*64)*128 = 73,728 CONV3-128: [112x112x128] memory: 112*112*128=1.6M params: (3*3*128)*128 = 147,456 POOL2: [56x56x128] memory: 56*56*128=400K params: 0 CONV3-256: [56x56x256] memory: 56*56*256=800K params: (3*3*128)*256 = 294,912 CONV3-256: [56x56x256] memory: 56*56*256=800K params: (3*3*256)*256 = 589,824 CONV3-256: [56x56x256] memory: 56*56*256=800K params: (3*3*256)*256 = 589,824 POOL2: [28x28x256] memory: 28*28*256=200K params: 0 CONV3-512: [28x28x512] memory: 28*28*512=400K params: (3*3*256)*512 = 1,179,648 CONV3-512: [28x28x512] memory: 28*28*512=400K params: (3*3*512)*512 = 2,359,296 CONV3-512: [28x28x512] memory: 28*28*512=400K params: (3*3*512)*512 = 2,359,296 POOL2: [14x14x512] memory: 14*14*512=100K params: 0 CONV3-512: [14x14x512] memory: 14*14*512=100K params: (3*3*512)*512 = 2,359,296 CONV3-512: [14x14x512] memory: 14*14*512=100K params: (3*3*512)*512 = 2,359,296 CONV3-512: [14x14x512] memory: 14*14*512=100K params: (3*3*512)*512 = 2,359,296 POOL2: [7x7x512] memory: 7*7*512=25K params: 0 FC: [1x1x4096] memory: 4096 params: 7*7*512*4096 = 102,760,448 FC: [1x1x4096] memory: 4096 params: 4096*4096 = 16,777,216 FC: [1x1x1000] memory: 1000 params: 4096*1000 = 4,096,000

TOTAL memory: 24M * 4 bytes ~= 96MB / image (only forward! ~*2 for bwd) TOTAL params: 138M parameters

[Simonyan and Zisserman, 2014]

```
(not counting biases)
```

* slide from Fei-Dei Li, Justin Johnson, Serena Yeung, cs231n Stanford

VGG16

ILSVRC winner 2012

GoogleLeNet

even deeper network with computational efficiency

- -22 layers
- Efficient "Inception" module
- No FC layers
- Only 5 million parameters!
- (12x less than AlexNet!)
- Better performance (@6.7 top 5 error)

[Szegedy et al., 2014]

these modules

Szegedy et al., 2014]

Idea: design good local topology ("network within network") and then stack



these modules

Apply parallel filter operations on the input from previous layer

 Multiple receptive field sizes for convolution (1x1, 3x3, 5x5)

- Pooling operation (3x3)

Concatenate all filter outputs together at output depth-wise Szegedy et al., 2014

Idea: design good local topology ("network within network") and then stack



Naive Inception module

these modules

Apply parallel filter operations on the input from previous layer

 Multiple receptive field sizes for convolution (1x1, 3x3, 5x5)

- Pooling operation (3x3)

Concatenate all filter outputs together at output depth-wise Szegedy et al., 2014

Idea: design good local topology ("network within network") and then stack

What's the problem?



Naive Inception module

these modules

Apply parallel filter operations on the input from previous layer

 Multiple receptive field sizes for convolution (1x1, 3x3, 5x5)

- Pooling operation (3x3)

Concatenate all filter outputs together at output depth-wise

Szegedy et al., 2014]

Idea: design good local topology ("network within network") and then stack



these modules

Apply parallel filter operations on the input from previous layer

 Multiple receptive field sizes for convolution (1x1, 3x3, 5x5)

- Pooling operation (3x3)

Concatenate all filter outputs together at output depth-wise

Szegedy et al., 2014]

Idea: design good local topology ("network within network") and then stack

28x28x192 28x28x128 28x28x96 28x28x256



these modules

Apply parallel filter operations on the input from previous layer

 Multiple receptive field sizes for convolution (1x1, 3x3, 5x5)

- Pooling operation (3x3)

Concatenate all filter outputs together at output depth-wise

Szegedy et al., 2014

Idea: design good local topology ("network within network") and then stack



Convolutional Layer: 1x1 convolutions

56 x 56 x 64 **image**



Idea: design good local topology ("network within network") and then stack these modules



Naive Inception module

Szegedy et al., 2014]

1x1 "bottleneck" layers



Inception module with dimension reduction

saves approximately 60% of computations

GoogleLeNet

even deeper network with computational efficiency

- -22 layers
- Efficient "Inception" module
- No FC layers
- Only 5 million parameters!
- (12x less than AlexNet!)
- Better performance (@6.7 top 5 error)

[Szegedy et al., 2014]



ILSVRC winner 2012



ResNet

even deeper — **152 layers**! using residual connections

[He et al., 2015]







[He et al., 2015]

What happens when we continue to stack ing deeper layers on a "plain" CNN





[He et al., 2015]

What happens when we continue to stack ing deeper layers on a "plain" CNN



Whats the **problem**?





[He et al., 2015]

What happens when we continue to stack ing deeper layers on a "plain" CNN



Whats the **problem**?



Hypothesis: deeper models are harder to optimize (optimization problem)

[He et al., 2015]



Hypothesis: deeper models are harder to optimize (optimization problem)

Observation: the deeper model should (conceptually) perform just as well (e.g., take shallower model and use identity for all remaining layers)

[He et al., 2015]



Hypothesis: deeper models are harder to optimize (optimization problem)

Observation: the deeper model should (conceptually) perform just as well (e.g., take shallower model and use identity for all remaining layers)

How do we implement this idea in practice

[He et al., 2015]



ResNet

Solution: use network to fit residual mapping instead of directly trying to fit a desired underlying mapping





[He et al., 2015]



ResNet

Full details

- Stacked **residual blocks**
- Every residual block consists of two 3x3 filters
- Periodically double # of filters and downsample spatially using stride of 2
- Additional convolutional layer in the beginning
- No FC layers at the end (only FC to output 1000 classes)

[He et al., 2015]







ILSVRC winner 2012



ILSVRC winner 2012



Regularization: Stochastic Depth

Effectively "dropout" but for layers

some layer (for each batch)



Huang et al., ECCV 2016]



Comparing **Complexity**



An Analysis of Deep Neural Network Models for Practical Applications, 2017.



Computer Vision Problems (no language for now)

Categorization

Detection





Multi-class: Horse Church Toothbrush Person **IM** GENET

Multi-label: Horse

Church Toothbrush Person

Horse (x, y, w, h) Horse (x, y, w, h) Person (x, y, w, h) Person (x, y, w, h)





Segmentation

Horse Person



Instance Segmentation

Horse1 Horse₂ Person1 Person2



Computer Vision Problems (no language for now)



Segmentation



Horse Person



Semantic Segmentation

Label every pixel with a category label (without differentiating instances)







Sky





Semantic Segmentation: Sliding Window

Extract **patches**



[Farabet et al, TPAMI 2013] [Pinheiro et al, ICML 2014]

Classify center pixel with CNN







Semantic Segmentation: Sliding Window

Extract **patches**



Problem: VERY inefficient, no reuse of computations for overlapping patches

[Farabet et al, TPAMI 2013] [•] Pinheiro et al, ICML 2014]

Classify center pixel with CNN





Design a network as a number of convolutional layers to make predictions for all pixels at once!





Problem: Convolutions at the original image scale will be very expensive

Design a network as a number of convolutional layers to make predictions for all pixels at once!





Input **Image**

 $3 \times H \times W$



 $D_1 \times H/2 \times W/2$

Design a network as a number of convolutional layers with downsampling and upsampling inside the network!



Predicted Labels

HxW

[Long et al, CVPR 2015] [Noh et al, ICCV 2015]



Input **Image**

 $3 \times H \times W$

High-res: $D_1 \times H/2 \times W/2$

Downsampling = Pooling

Design a network as a number of convolutional layers with downsampling and upsampling inside the network!





Predicted Labels

HxW

Upsampling = ???

[Long et al, CVPR 2015] [Noh et al, ICCV 2015]

In-network Up Sampling (a.k.a "Unpooling")

Nearest Neighbor



Input: 2 x 2

Output: 4 × 4

In-network Up Sampling (a.k.a "Unpooling")

Nearest Neighbor



Input: 2 x 2

Output: 4 × 4

"Bed of Nails"



In-network Up Sampling: Max Unpooling

Max Pooling

Remember which element was max!



Corresponding pairs of downsampling and upsampling layers

Max Unpooling Use positions from pooling layer

In-network Up Sampling: Transpose Convolution

Recall: Normal 3 x 3 convolution, stride 1 pad 1



Input: 4 × 4

Dot product between filter and input



Output: 4 × 4
Recall: Normal 3 x 3 convolution, stride 1 pad 1



Dot product between filter and input

Input: 4 × 4



Output: 4 × 4

Recall: Normal 3 x 3 convolution, stride 2 pad 1



Input: 4 × 4

Dot product between filter and input



Output: 2 x 2

Recall: Normal 3 x 3 convolution, stride 1 pad 1





Input: 4 × 4

Dot product between filter and input



Output: 2 × 2

Filter moves 2 pixels in the **input** for every one pixel in the **output**

Stride gives ratio in movement in input vs output

3 x 3 transpose convolution, stride 2 pad 1

Input: 2 x 2

Output: 4 × 4

3 x 3 transpose convolution, stride 2 pad 1



Input gives weight for filter

Input: 2 x 2



Output: 4 × 4

3 x 3 transpose convolution, stride 2 pad 1



Input gives weight for filter

Input: 2 x 2



Output: 4 × 4

Filter moves 2 pixels in the **output** for every one pixel in the **input**

Stride gives ratio in movement in output vs input

Transpose Convolution: 1-D Example



Output contains copies of the filter weighted multiplied by the input, summing at overlaps in the output

Computer Vision Problems (no language for now)

Categorization

Detection





Multi-class: Horse Church Toothbrush Person **IM** GENET

Multi-label: Horse

Church Toothbrush Person

Horse (x, y, w, h) Horse (x, y, w, h) Person (x, y, w, h) Person (x, y, w, h)





Segmentation

Horse Person



Instance Segmentation

Horse1 Horse₂ Person1 Person2



Computer Vision Problems (no language for now)

Detection



Horse (x, y, w, h) Horse (x, y, w, h) Person (x, y, w, h) Person (x, y, w, h)





Datasets: Pascal VOC

20 classes: aeroplane, bicycle, boat, bottle, bus, car, cat, chair, cow, dining table, dog, horse, motorbike, person, potted plant, sheep, train, TV



Real images downloaded from flickr, not filtered for "quality"

* slide from Andrew Zisserman

Datasets: Pascal VOC

20 classes: aeroplane, bicycle, boat, bottle, bus, car, cat, chair, cow, dining table, dog, horse, motorbike, person, potted plant, sheep, train, TV



Real images downloaded from flickr, not filtered for "quality"

* slide from Andrew Zisserman

Datasets: COCO



Object segmentation
Recognition in context
Superpixel stuff segmentation
330K images (>200K labeled)
1.5 million object instances
80 object categories
91 stuff categories
5 captions per image
250,000 people with keypoints

Object **Detection**



* plot from Ross Girshick, 2015



Object **Detection** as Regression Problem





Object **Detection** as Regression Problem









Object **Detection** as Regression Problem





Problem: each image needs a different number of outputs











Apply CNN to many different crops in the image and (classification) CNN classifies each patch as object or background





Apply CNN to many different crops in the image and (classification) CNN classifies each patch as object or background





Apply CNN to many different crops in the image and (classification) CNN classifies each patch as object or background





Apply CNN to many different crops in the image and (classification) CNN classifies each patch as object or background





Apply CNN to many different crops in the image and (classification) CNN classifies each patch as object or background





Problem: Need to apply CNN to **many** patches in each image

Apply CNN to many different crops in the image and (classification) CNN classifies each patch as object or background

Region Proposals (older idea in vision)

Find image regions that are likely contain objects (any object at all)



[Alexe et al, TPAMI 2012] [Uijkings et al, IJCV 2013] [Cheng et al, CVPR 2014] [Zitnick and Dollar, ECCV 2014]

- typically works by looking at histogram distributions, region aspect ratio, closed contours, coherent color

Relatively fast to run (Selective Search gives 1000 region proposals in a few seconds on a CPU)



Goal: Get "true" object regions to be in as few top K proposals as possible



[Girshick et al, CVPR 2014]





[Girshick et al, CVPR 2014]





[Girshick et al, CVPR 2014]

Warped image regions

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





[Girshick et al, CVPR 2014]

Forward each region through a CNN

Warped image regions

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





[Girshick et al, CVPR 2014]

Classify regions with SVM

Forward each region through a **CNN**

Warped image regions

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



Linear Regression for bounding box offsets



[Girshick et al, CVPR 2014]

Classify regions with SVM

Forward each region through a **CNN**

Warped image regions

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



R-CNN: Training

Fine-tuning ImageNet CNN on object proposal patches

- > 50% Intersection-over-Union overlap with GT considered "object" others "background"
- batches of 128 (**32 positives, 96 negatives**)

[Girshick et al, CVPR 2014]





R-CNN: Issues

Ad-hoc training objectives

- Fine-tune network with softmax objective (**log** loss)
- Train post-hoc linear SVM (**hinge** loss)
- Train post-hoc bounding-box regression (least squares)

Training is slow

84 hours and takes a lot of disk space

Inference / Detection is slow

- 47 sec / image with VGG16 [Simonyan et al, ICLR 2015]

[Girshick et al, CVPR 2014]





R-CNN vs. SPP



R-CNN 2000 nets on image regions

[He et al, ECCV 2014]



SPP-net **1 net on full image**





[Girshick et al, ICCV 2015]

Input Image





[Girshick et al, ICCV 2015]

Input Image





[Girshick et al, ICCV 2015]

"conv5" feature map

Forward prop the **whole image** through CNN

Input **Image**



Regions of Interest "conv5" feature map from the Forward prop the **whole image** through CNN proposal method ConvNet

[Girshick et al, ICCV 2015]



Input **Image**


Fast **R-CNN**

Regions of $\overline{}$ Interest from the proposal method ConvNet

R/

[Girshick et al, ICCV 2015]

- "Rol Pooling" layer
- "conv5" feature map
 - Forward prop the whole image through CNN



Input Image

Girshick, "Fast R-C Figure copyright Re

* image from Ross Girshick



Fast **R-CNN**

Object classification

Regions of Interest from the proposal method



Multi-task loss

[Girshick et al, ICCV 2015]

Bounding box regression

- "Rol Pooling" layer
- "conv5" feature map
 - Forward prop the **whole image** through CNN

Input **Image**

* image from Ross Girshick





Multi-task loss

[Girshick et al, ICCV 2015]

Bounding box regression

- "Rol Pooling" layer
- "conv5" feature map
 - Forward prop the **whole image** through CNN

Input Image

* image from Ross Girshick



R-CNN vs. SPP vs. Fast R-CNN



[Girshick et al, CVPR 2014] [Girshick et al, ICCV 2015] [He et al, ECCV 2014]



R-CNN vs. SPP vs. Fast R-CNN



Observation: Performance dominated by the region proposals at this point!

Girshick et al, CVPR 2014 [Girshick et al, ICCV 2015] [He et al, ECCV 2014]



Faster **R-CNN**

Make CNN do proposals!

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



Jointly train with 4 losses:

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

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



LSDA: Large Scale Detection through Adaptation





 $+\delta \mathbf{W}_{cat}$

[Hoffman et al, NIPS 2014]

YOLO: You Only Look Once





Input image 3 x H x W

Divide image into grid 7 x 7

Image a set of **base boxes** centered at each grid cell Here B = 3

Redmon et al, CVPR 2016]

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: $7 \times 7 \times (5 * B + C)$







YOLO: You Only Look Once





Input image 3 x H x W

Divide image into grid 7 x 7

Image a set of **base boxes** centered at each grid cell Here B = 3

[Redmon et al, CVPR 2016]



http://pureddie.com/yolo





http://pureddie.com/yolo



