Deep Learning

So far: Image Classification



Computer Vision Tasks

Classification

Semantic Segmentation

Object Detection

Instance Segmentation



Object Detection: Impact of Deep Learning



Reproduced with permission.

Last Time: Object Detection Methods

"Slow" R-CNN: Run CNN independently for each region



Fast R-CNN: Apply differentiable cropping to shared image features



Recap: Slow R-CNN Training

"Slow" R-CNN: Run CNN independently for each region



Fast R-CNN: Apply differentiable cropping to shared image features



Input Image



Ground-Truth boxes

Input Image



Ground-Truth boxes

Region Proposals

Input Image



Categorize each region proposal as positive, negative, or neutral based on overlap with ground-truth boxes

Input Image











Crop pixels from each positive and negative proposal, resize to 224 x 224

Run each region through CNN. For positive boxes predict class and box offset; for negative boxes just predict background class

Input Image







Class target: Dog Box target: ——>





Class target: Cat Box target: _____;





Class target: Dog Box target: ——>





Class target: Background Box target: None

Fast R-CNN Training

Crop features for each region, use them to predict class and box targets per region



Faster R-CNN: Learnable Region Proposals

Classification

loss

Jointly train with 4 losses:

- **1. RPN classification**: anchor box is object / not an object
- **2. RPN regression**: predict transform from anchor box to proposal box
- 3. Object classification: classify proposals as background / object class
- **4. Object regression**: predict transform from proposal box to object box
- Anchor -> Region Proposal -> Object Box (Stage 1) (Stage 2)



Bounding-box

regression loss

proposals

Classification

loss

Bounding-box

regression loss

Rol pooling

RPN predicts Object / Background for each anchor, as well as regresses from Faster R-CNN Training: RPN Training anchor to object box







RPN gives lots of **anchors** which we classify as pos / neg / neutral by matching with ground-truth

CNN

Crop features for each proposal, use them to predict class and box targets per region

Faster R-CNN Training: Stage 2

Input Image





Image Features



Now proposals come from RPN rather than selective search, but otherwise this works the same as Fast R-CNN training



Recap: Fast R-CNN Feature Cropping

"Slow" R-CNN: Run CNN independently for each region



Fast R-CNN: Apply differentiable cropping to shared image features



Cropping Features: Rol Pool



Goal: Crop features for region proposal, and resize to a fixed size for downstream processing, in a <u>differentiable</u> way



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



regions have different sizes!

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



Divide into equal-sized subregions (may not be aligned to grid!)



He et al, "Mask R-CNN", ICCV 2017



Divide into equal-sized subregions (may not be aligned to grid!)

Sample features at regularly-spaced points in each subregion using **bilinear interpolation**

He et al, "Mask R-CNN", ICCV 2017

Divide into equal-sized subregions (may not be aligned to grid!)

> Sample features at regularly-spaced points in each subregion using **bilinear interpolation**



Feature f_{xy} for point (x, y) is a linear combination of features at its four neighboring grid cells:

Input Image (e.g. 3 x 640 x 480)

<u>Align</u>



Divide into equal-sized subregions (may not be aligned to grid!)



Feature f_{xy} for point (x, y) is a linear combination of features at its four neighboring grid cells:

Project proposal

CNN

onto features

<u>Align</u>

Divide into equal-sized subregions (may not be aligned to grid!)

No "snapping"!

Sample features at regularly-spaced points in each subregion using bilinear interpolation



Feature f_{xy} for point (x, y) is a linear combination of features at its four neighboring grid cells:



Project proposal

onto features

<u>Align</u>

Divide into equal-sized subregions (may not be aligned to grid!)

> Sample features at regularly-spaced points in each subregion using **bilinear interpolation**



No "snapping"!

Project proposal

onto features

<u>Align</u>

Divide into equal-sized subregions (may not be aligned to grid!)

No "snapping"! Sample features at regularly-spaced points in each subregion using bilinear interpolation



Project proposal

onto features

<u>Align</u>

Divide into equal-sized subregions (may not be aligned to grid!)

> Sample features at regularly-spaced points in each subregion using **bilinear interpolation**



No "snapping"!

+ $(f_{6.6} * 0.5 * 0.8) + (f_{7.6} * 0.5 * 0.8)$

<u>Align</u>

Divide into equal-sized subregions (may not be aligned to grid!)

No "snapping"! Sample features at **Project proposal** regularly-spaced points onto features in each subregion using bilinear interpolation f_{6,5} **†**_{7,5} **CNN** 6.5,5.8 0.2 0.5 f_{7.6} $f_{6.6}$ $f_{xy} = \sum_{i,j=1}^{2} f_{i,j} \max(0, 1 - |x - x_i|) \max(0, 1)$ $-y_j|)$ $\mathbf{f}_{6.5,5.8} = (\mathbf{f}_{6,5} * 0.5 * 0.2) + (\mathbf{f}_{7,5} * 0.5 * 0.2)$

Feature f_{xv} for point (x, y) is a linear combination of features at its four neighboring grid cells:

Divide into equal-sized subregions (may not be aligned to grid!)

No "snapping"! Sample features at <u>Align</u> **Project proposal** regularly-spaced points onto features in each subregion using bilinear interpolation f_{6,5} $f_{7,5}$ **CNN** 6.5,5.8 $f_{6.6}$ $f_{7,6}$ $f_{xy} = \sum_{i,j=1}^{2} f_{i,j} \max(0, 1 - |x - x_i|) \max(0, 1 - |y - y_j|)$

This is differentiable! Upstream gradient for sampled feature will flow backward into each of the four nearest-neighbor gridpoints

Feature f_{xy} for point (x, y) is a linear combination of features at its four neighboring grid cells:



Divide into equal-sized subregions (may not be aligned to grid!)

> Sample features at regularly-spaced points in each subregion using **bilinear interpolation**

After sampling, maxpool in each subregion



Region features (here 512 x 2 x 2; In practice e.g 512 x 7 x 7)

Computer Vision Tasks: Object Detection

Classification

Semantic Segmentation Object Detection Instance Segmentation



Computer Vision Tasks: Semantic Segmentation



Semantic Segmentation

Label each pixel in the image with a category label

Don't differentiate instances, only care about pixels



Semantic Segmentation Idea: Sliding Window



Farabet et al, "Learning Hierarchical Features for Scene Labeling," TPAMI 2013 Pinheiro and Collobert, "Recurrent Convolutional Neural Networks for Scene Labeling", ICML 2014

Semantic Segmentation Idea: Sliding Window



between overlapping patches

Farabet et al, "Learning Hierarchical Features for Scene Labeling," TPAMI 2013 Pinheiro and Collobert, "Recurrent Convolutional Neural Networks for Scene Labeling", ICML 2014

Semantic Segmentation: Fully Convolutional Network

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



Long et al, "Fully convolutional networks for semantic segmentation", CVPR 2015

Loss function: Per-Pixel cross-entropy
Design a network as a bunch of convolutional layers to make predictions for pixels all at once



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



Input:Problem #1: Effective receptive3 x H x Wfield size is linear in number of
conv layers: With L 3x3 conv
layers, receptive field is 1+2L

Problem #2: Convolution on high res images is expensive. Recall ResNet stem aggressively downsamples

Long et al, "Fully convolutional networks for semantic segmentation", CVPR 2015

Design network as a bunch of convolutional layers, with **downsampling** and **upsampling** inside the network



Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR 2015 Noh et al, "Learning Deconvolution Network for Semantic Segmentation", ICCV 2015

Downsampling: Pooling, strided convolution



Input: 3 x H x W Design network as a bunch of convolutional layers, with **downsampling** and **upsampling** inside the network





Upsampling:

Predictions:

H x W

Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR 2015 Noh et al, "Learning Deconvolution Network for Semantic Segmentation", ICCV 2015

In-Network Upsampling: "Unpooling"

0

0

 \mathbf{O}

 \mathbf{O}

Bed of Nails



Input $C \times 2 \times 2$

Output $C \times 4 \times 4$ In-Network Upsampling: "Unpooling"

Bed of Nails

Nearest Neighbor



In-Network Upsampling: Bilinear Interpolation



Input: C x 2 x 2 Output: C x 4 x 4

 $f_{x,y} = \sum_{i,j} f_{i,j} \max(0, 1 - |x - i|) \max(0, 1 - |y - j|) \quad i \in \{\lfloor x \rfloor - 1, \dots, \lceil x \rceil + 1\}$ Use two closest neighbors in x and y to construct linear approximations $j \in \{\lfloor y \rfloor - 1, \dots, \lceil y \rceil + 1\}$

In-Network Upsampling: Bicubic Interpolation



Input: C x 2 x 2

Output: C x 4 x 4

Use **three** closest neighbors in x and y to construct **cubic** approximations

In-Network Upsampling: "Max Unpooling"

Max Pooling: Remember which position had the max Max Unpooling: Place into remembered positions





Pair each downsampling layer with an upsampling layer

Noh et al, "Learning Deconvolution Network for Semantic Segmentation", ICCV 2015

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

Input: 4 x 4

Output: 4 x 4

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



Input: 4 x 4

Output: 4 x 4

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



Input: 4 x 4

Output: 4 x 4

Recall: Normal 3 x 3 convolution, <u>stride 2</u>, pad 1

Input: 4 x 4

Output: 2 x 2

Recall: Normal 3 x 3 convolution, <u>stride 2</u>, pad 1

Input: 4 x 4

Output: 2 x 2

Recall: Normal 3 x 3 convolution, <u>stride 2</u>, pad 1

Dot product between input and filter

Input: 4 x 4

Output: 2 x 2

Recall: Normal 3 x 3 convolution, <u>stride 2</u>, pad 1

Convolution with stride > 1 is "Learnable Downsampling" Can we use stride < 1 for "Learnable Upsampling"?

Dot product between input and filter

Input: 4 x 4

Output: 2 x 2

3 x 3 convolution transpose, stride 2

Input: 2 x 2

Output: 4 x 4

3 x 3 convolution transpose, stride 2

Weight filter by input value and copy to output

Output: 4 x 4

3 x 3 convolution transpose, stride 2

Filter moves 2 pixels in <u>output</u> for every 1 pixel in <u>input</u>

Weight filter by input value and copy to output

Output: 4 x 4

3 x 3 convolution transpose, stride 2

Filter moves 2 pixels in <u>output</u> for every 1 pixel in <u>input</u>

Weight filter by input value and copy to output

Sum where

Output: 4 x 4

3 x 3 convolution transpose, stride 2

This gives 5x5 output – need to trim one pixel from top and left to give 4x4 output

Weight filter by input value and copy to output

Sum where

Output: 4 x 4

Transposed Convolution: 1D example

Filter Input Output ах ay Х a bx az+ Y b by Ζ bz

Output has copies of filter weighted by input

Stride 2: Move 2 pixels output for each pixel in input

Sum at overlaps

Transposed Convolution: 1D example

This has many names:

- Deconvolution (bad)!
- Upconvolution
- Fractionally strided convolution
- Backward strided convolution
- Transposed Convolution

(best name)

We can express convolution in terms of a matrix multiplication

$$\vec{x} * \vec{a} = X\vec{a}$$

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

Example: 1D conv, kernel size=3, stride=1, padding=1

We can express convolution in terms of a matrix multiplication

Transposed convolution multiplies by the transpose of the same matrix:

 $T \rightarrow -T \rightarrow$

$$\vec{x} * \vec{a} = X\vec{a} \qquad \qquad \vec{x} *^{T} \vec{a} = X^{T} \vec{a}$$

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

Example: 1D conv, kernel size=3, stride=1, padding=1

When stride=1, transposed conv is just a regular conv (with different padding rules)

We can express convolution in terms of a matrix multiplication

$$\vec{x} * \vec{a} = X\vec{a}$$

$$\begin{bmatrix} x & y & x & 0 & 0 & 0 \\ 0 & 0 & x & y & x & 0 \end{bmatrix} \begin{bmatrix} 0 \\ a \\ b \\ c \\ d \\ 0 \end{bmatrix} = \begin{bmatrix} ay + bz \\ bx + cy + dz \end{bmatrix}$$

Transposed convolution multiplies by the transpose of the same matrix:

$$\vec{x} *^T \vec{a} = X^T \vec{a}$$

Example: 1D conv, kernel size=3, <u>stride=2</u>, padding=1

We can express convolution in terms of a matrix multiplication

$$\vec{x} * \vec{a} = X\vec{a}$$

$$\begin{bmatrix} x & y & x & 0 & 0 & 0 \\ 0 & 0 & x & y & x & 0 \end{bmatrix} \begin{bmatrix} 0 \\ a \\ b \\ c \\ d \\ 0 \end{bmatrix} = \begin{bmatrix} ay + bz \\ bx + cy + dz \end{bmatrix}$$

Transposed convolution multiplies by the transpose of the same matrix:

$$\vec{x} *^{T} \vec{a} = X^{T} \vec{a}$$

$$\begin{bmatrix} x & 0 \\ y & 0 \\ z & x \\ 0 & y \\ 0 & z \\ 0 & 0 \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = \begin{bmatrix} ax \\ ay \\ az + bx \\ by \\ bz \\ 0 \end{bmatrix}$$

Example: 1D conv, kernel size=3, <u>stride=2</u>, padding=1

Downsampling: Pooling, strided convolution

Input: 3 x H x W Design network as a bunch of convolutional layers, with downsampling and upsampling inside the network

Upsampling: linterpolation, transposed conv

Predictions: H x W

Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR 2015 Noh et al, "Learning Deconvolution Network for Semantic Segmentation", ICCV 2015

Loss function: Per-Pixel cross-entropy

Computer Vision Tasks

Object Detection: Detects individual object instances, but only gives box

Semantic Segmentation: Gives perpixel labels, but merges instances

Things and Stuff

Things: Object categories that can be separated into object instances (e.g. cats, cars, person)

Stuff: Object categories that cannot be separated into instances (e.g. sky, grass, water, trees)

Computer Vision Tasks

Object Detection: Detects individual object instances, but only gives box (only things)

Semantic Segmentation: Gives perpixel labels, but merges instances (both things and stuff)

Computer Vision Tasks: Instance Segmentation

Computer Vision Tasks: Instance Segmentation

Instance Segmentation:

Detect all objects in the image, and identify the pixels that belong to each object (only things)

Computer Vision Tasks: Instance Segmentation

Instance Segmentation:

Detect all objects in the image, and identify the pixels that belong to each object (only things)

Approach: Perform object detection, then predict a segmentation mask for each object!

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



Predict a mask for each of C classes: C x 28 x 28









Mask R-CNN: Very Good Results!



Beyond Instance Segmentation

Instance Segmentation: Separate object instances, but only things



Semantic Segmentation: Identify both things and stuff, but doesn't separate instances



Beyond Instance Segmentation: Panoptic Segmentation

Label all pixels in the image (both things and stuff)

For "thing" categories also separate into instances



Kirillov et al, "Panoptic Segmentation", CVPR 2019 Kirillov et al, "Panoptic Feature Pyramid Networks", CVPR 2019

Beyond Instance Segmentation: Panoptic Segmentation



Kirillov et al, "Panoptic Feature Pyramid Networks", CVPR 2019

Beyond Instance Segmentation: Human Keypoints

Represent the pose of a human by locating a set of **keypoints**

e.g. 17 keypoints:

- Nose
- Left / Right eye
- Left / Right ear
- Left / Right shoulder
- Left / Right elbow
- Left / Right wrist
- Left / Right hip
- Left / Right knee
- Left / Right ankle











Ground-truth has one "pixel" turned on per keypoint. Train with softmax loss

Joint Instance Segmentation and Pose Estimation







CHARK P

image

Dense Captioning





Summary: Many Computer Vision Tasks

Classification

Semantic Segmentation

Object Detection

Instance Segmentation

