

THE UNIVERSITY OF BRITISH COLUMBIA

CPSC 425: Computer Vision



Lecture 23: Neural Network Vision Applications

Menu for Today

Topics:

- Object Detection
- Segmentation

Redings:

- Today's Lecture: N/A

Reminders:

- Assignment 6: Deep Learning is out and due next Thursday
- No lecture on Monday (holiday)
- No office hours on Friday (holiday)
- Additional **Office Hours** week after next (for Final Prep)

Image Generation — Quiz 5 and 6 **online** (see Canvas) **Note**: grading same 0.5 point for participation



Please fill out **Student Evaluations** (on Canvas)



3 depth

Filters always extend the full depth of the input volume

5 x 5 x 3 filter

Convolve the filter with the image (i.e., "slide over the image spatially, computing dot products"





Review Lecture 22: Convolutional Layer

If we have 6 5x5 filter, we'll get 6 separate activation maps:





this results in the "new image" of size 28 x 28 x 6!

activation map





3 depth





3 depth









3 depth



28 width





3 depth





3 depth





Review Lecture 22: Pooling Layer

- Makes representation smaller, more manageable and spatially invariant Operates over each activation map independently





Review Lecture 22: Deep Learning Terminology



generally kept fixed, requires some knowledge of the problem and NN to sensibly set

requires knowledge of the nature of the problem

• **Network structure:** number and types of layers, forms of activation functions, dimensionality of each layer and connections (defines computational graph)

deeper = better

• Loss function: objective function being optimized (softmax, cross entropy, etc.)

• Parameters: trainable parameters of the network, including weights/biases of linear/fc layers, parameters of the activation functions, etc. optimized using SGD or variants

• Hyper-parameters: parameters, including for optimization, that are not optimized

directly as part of training (e.g., learning rate, batch size, drop-out rate) grid search



Review Lecture 22: Deep Learning Terminology

generally kept fixed, requires some knowledge of the problem and NN to sensibly set deeper = better

requires knowledge of the nature of the problem

Specification of neural architecture will define a **computational** graph.



• **Network structure:** number and types of layers, forms of activation functions, dimensionality of each layer and connections (defines computational graph)

• Loss function: objective function being optimized (softmax, cross entropy, etc.)

Computer Vision Problems

Computer Vision Problems

Categorization



Multi-class: Horse Church Toothbrush Person IM GENET

Multi-**label**: Horse

Church Toothbrush Person



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



Computer Vision Problems (no language for now)

Categorization

Detection





Horse Multi-class: 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)





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



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



Region Proposals (older idea in vision)

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



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

[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)





[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 (Regions with CNN features) algorithm:

- Extract promising candidate regions using an object proposals algorithm
- Resize each proposal window to the size of the input layer of a trained convolutional neural network
- Input each resized image patch to the convolutional neural network

Implementation detail: Instead of using the classification scores of the input feature to a trained support vector machine (SVM)

network directly, the output of the final fully-connected layer can be used as an

R-CNN vs. SPP



R-CNN 2000 nets on image regions

[He et al, ECCV 2014]



R-CNN vs. SPP



R-CNN 2000 nets on image regions

[He et al, ECCV 2014]



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




* image from Ross Girshick

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**



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

[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



Rol Align



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**





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**



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]



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



YOLO: You Only Look Once





Input image $3 \times H \times W$

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)

Divide image into grid 7 x 7

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





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]



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!



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

(in terms of storage)



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

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 2 pad 1



Dot product between filter and input

Input: 4 × 4



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 × 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

U-Net Architecture

ResNet-like Fully convolutional CNN



[Ronneberger et al, CVPR 2015]
Image Colorization



[Zhang et al. 2016]

Depth Estimation







Direct supervision via Kinect RGB+D

U-Net with skip connections

Loss, e.g., L2



Depth Estimation





[DenseNet Huang et al 2018] [Alhashim Wonka 2019]

Super Resolution

bicubic (21.59dB/0.6423)







SRResNet (23.53dB/0.7832)

original





A state of the art super-res network trained with L2 loss is good at sharpening edges, but results lack realistic texture [Ledig et al 2017]

Image-to-Image Translation





• • •

Figure 1 in the original paper.

[lsola et al., 2016]

Generative Adversarial Networks

Architecture: DCGAN-based

Training is conditioned on the **images** from the source domain





Figure 2 in the original paper.

[Isola et al., 2016]

Image-to-Image Translation

Style transfer: change the style of an image while preserving the content



Data: two unrelated collections of image, one for each style

[Zhu et al., 2017]

Denoising **Diffusion** Models

Key Idea: Learning to generate by denoising

Denoising diffusion models consist of two processes:

- Forward diffusion process that gradually adds noise to input
- Reverse denoising process that learns to generate data by denoising

Forward diffusion process (fixed)

Data



Sohl-Dickstein et al., Deep Unsupervised Learning using Nonequilibrium Thermodynamics, ICML 2015 Ho et al., Denoising Diffusion Probabilistic Models, NeurIPS 2020 Song et al., Score-Based Generative Modeling through Stochastic Differential Equations, ICLR 2021

Slide from: https://www.cs.unc.edu/~ronisen/teaching/fall_2022/pdf_lectures/lecture7-8_diffusion_model.pdf

Noise

Reverse denoising process (generative)

Forward Diffusion Process

The formal definition of the forward process in T steps:



Slide from: https://www.cs.unc.edu/~ronisen/teaching/fall_2022/pdf_lectures/lecture7-8_diffusion_model.pdf

Forward diffusion process (fixed)

Noise

Reverse Denoising Process

Data

Formal definition of forward and reverse processes in T steps:



Slide from: https://www.cs.unc.edu/~ronisen/teaching/fall_2022/pdf_lectures/lecture7-8_diffusion_model.pdf

Noise

the denoised image x(t-1) given x(t).

Implementation: Network Architecture

Diffusion models often use U-Net architectures with ResNet blocks and self-attention layers to represent $\epsilon_{\theta}(\mathbf{x}_t, t)$



Time representation: sinusoidal positional embeddings or random Fourier features.

Time features are fed to the residual blocks using either simple spatial addition or using adaptive group normalization layers. (see <u>Dharivwal and Nichol NeurIPS 2021</u>)

Slide from: https://www.cs.unc.edu/~ronisen/teaching/fall_2022/pdf_lectures/lecture7-8_diffusion_model.pdf

28

4. Visual Imagination



A brain riding a rocketship heading towards the moon.

A dragon fruit wearing karate belt in the snow.

A marble statue of a Koala DJ in front of a marble statue of a turntable. The Koala has wearing large marble headphones.



A Pomeranian is sitting on the Kings throne wearing a crown. Two tiger soldiers are standing next to the throne.

An extremely angry bird.

Android Mascot made from bamboo.

Three spheres made of glass falling into ocean. Water is splashing. Sun is setting.

A single beam of light enter the room from the ceiling. The beam of light is illuminating an easel. On the easel there is a Rembrandt painting of a raccoon.

- imagen.research.google
- Text to image generation
- Uses diffusion process, training using large dataset of text (web scale) and image-text (400M) pairs



"A Golden Retriever dog wearing a blue checkered beret and red dotted turtleneck."





Computer Vision Problems (no language for now)



Instance Segmentation



Horse1 Horse2 Person1 Person2

Mask R-CNN



[He et al, 2017]

Summary

Common types of layers:

- 1. Convolutional Layer - Parameters define a set of learnable filters
- 2. **Pooling** Layer - Performs a downsampling along the spatial dimensions
- 3. Fully-Connected Layer As in a regular neural network

Each layer accepts an input 3D volume and transforms it to an output 3D volume through a differentiable function

Summary

The parameters of a neural network are learned using **backpropagation**, which computes gradients via recursive application of the chain rule

the network architecture to reduce the number of parameters

A convolutional layer applies a set of learnable filters

A **pooling layer** performs spatial downsampling

A fully-connected layer is the same as in a regular neural network

- A convolutional neural network assumes inputs are images, and constrains
- Convolutional neural networks can be seen as learning a hierarchy of filters