

Lecture 12: Optimization for Neural Networks, Convolutional Neural Network

Honglak Lee

02/19/2025



Outline

- Optimization
- CNN basics
- Examples of CNN Architectures
- Applications of CNN

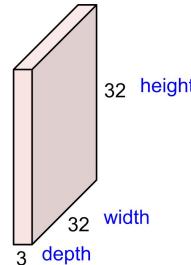
Image Classification



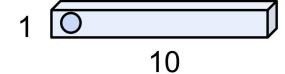
“airplane”

Image Classification

32x32x3 image



10-class output



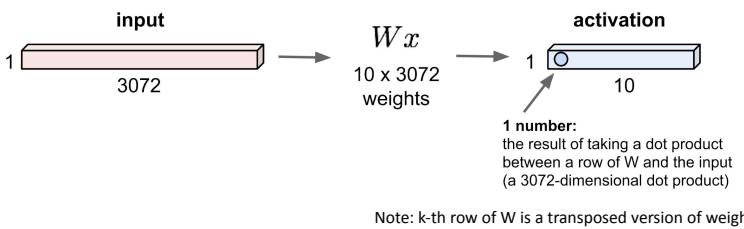
Note: in this lecture, we use vectors visualized as rows

24

Slide credit: Fei-Fei Li, Justin Johnson, Serena Yeung 25

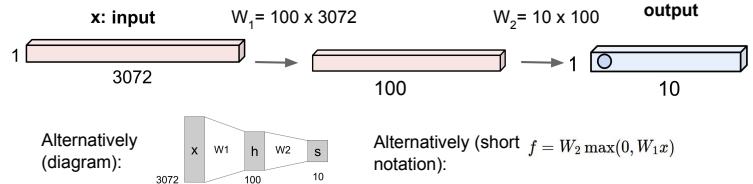
Fully Connected Layer

32x32x3 image -> stretch to 3072 x 1



Slide credit: Fei-Fei Li, Justin Johnson, Serena Yeung 27

Neural Network with 1 hidden layer (with Rectified Linear Units)



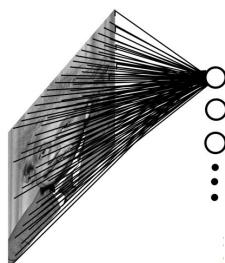
Deeper models: We can add more layers, etc.

Q. Have we solved the problem? Will this be applicable for larger-sized images?

Slide credit: Fei-Fei Li, Justin Johnson, Serena Yeung 28

Scaling to larger-sized images

Example: 200x200 image, 40k hidden units: Fully connected neural network would have ~2B parameters!



Slide credit: Russ Salakhutdinov

Motivation: Convolutional Neural Networks

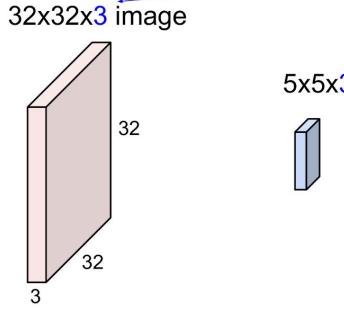
- How can we design neural networks that are specifically adapted for classifying large-sized images?
 - Must deal with very **high-dimensional inputs**: 150 x 150 pixels = 22500 inputs, or 3 x 22500 if RGB pixels
 - Can exploit the **2D topology** of pixels (or 3D for video data)
 - Can build in **invariance** to certain variations: translation, illumination, etc.
- **Convolutional networks** leverage these ideas
 - Local connectivity
 - Convolution, parameter sharing
 - Pooling / subsampling hidden units

Slide credit: Russ Salakhutdinov

29

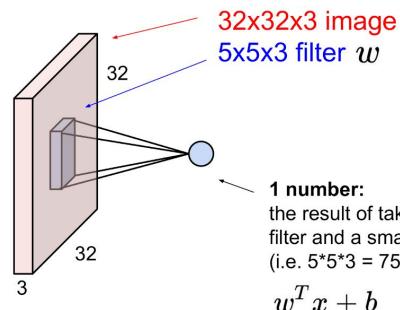
30

Convolution Layer



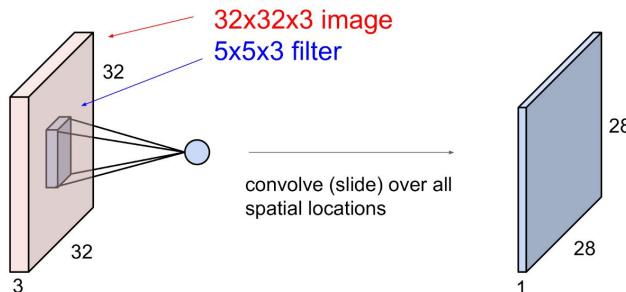
Slide credit: Fei-Fei Li, Justin Johnson, Serena Yeung 31

Convolution Layer



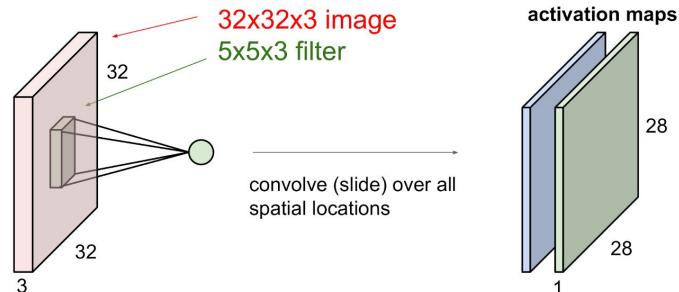
Slide credit: Fei-Fei Li, Justin Johnson, Serena Yeung 32

Convolution Layer



Slide credit: Fei-Fei Li, Justin Johnson, Serena Yeung 33

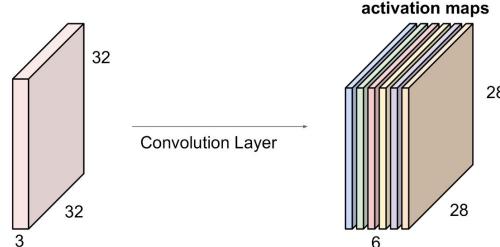
Convolution Layer



Slide credit: Fei-Fei Li, Justin Johnson, Serena Yeung 34

Convolution Layer

For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:



We stack these up to get a "new image" of size 28x28x6!

Slide credit: Fei-Fei Li, Justin Johnson, Serena Yeung 35

Discrete Convolution

- Applying filter to an image via "sliding window" can be efficiently computed as a convolution of an image x with kernel k as follows

$$(x * \text{valid } k)_{i,j} = \sum_{m=i}^{i+H_k-1} \sum_{n=j}^{j+W_k-1} x_{m,n} k_{i-m+H_k, j-n+W_k}$$

Example

1	.5	
.25	0	
		filter

0	80	40	
20	40	0	
0	0	40	image x

$$1 \times 0 + 0.5 \times 80 + 0.25 \times 20 + 0 \times 40 = 45$$

kernel is a flipped version of a filter

$$* \quad \begin{matrix} 0 & .25 \\ .5 & 1 \end{matrix} = \begin{matrix} 45 & \\ & \end{matrix}$$

kernel is a flipped version of a filter

38

Discrete Convolution

- Applying filter to an image via "sliding window" can be efficiently computed as a convolution of an image x with kernel k as follows

- Example

$$(x * \text{valid } k)_{i,j} = \sum_{m=i}^{i+H_k-1} \sum_{n=j}^{j+W_k-1} x_{m,n} k_{i-m+H_k, j-n+W_k}$$

kernel is a flipped version of a filter

1	.5	
.25	0	
		filter

0	80	40	
20	40	0	
0	0	40	image x

$$1 \times 80 + 0.5 \times 40 + 0.25 \times 40 + 0 \times 0 = 110$$

kernel is a flipped version of a filter

$$* \quad \begin{matrix} 0 & .25 \\ .5 & 1 \end{matrix} = \begin{matrix} 45 & 110 \\ & \end{matrix}$$

kernel is a flipped version of a filter

39

Discrete Convolution

- Applying filter to an image via "sliding window" can be efficiently computed as a convolution of an image x with kernel k as follows

- Example

$$(x * \text{valid } k)_{i,j} = \sum_{m=i}^{i+H_k-1} \sum_{n=j}^{j+W_k-1} x_{m,n} k_{i-m+H_k, j-n+W_k}$$

kernel is a flipped version of a filter

1	.5	
.25	0	
		filter

0	80	40	
20	40	0	
0	0	40	image x

$$1 \times 20 + 0.5 \times 40 + 0.25 \times 0 + 0 \times 0 = 40$$

kernel is a flipped version of a filter

$$* \quad \begin{matrix} 0 & .25 \\ .5 & 1 \end{matrix} = \begin{matrix} 45 & 40 \\ & \end{matrix}$$

kernel is a flipped version of a filter

40

Discrete Convolution

- Applying filter to an image via “sliding window” can be efficiently computed as a convolution of an image x with kernel k as follows

- Example

$$\begin{matrix} 1 & .5 \\ .25 & 0 \end{matrix}$$

filter

$$(x *_{\text{valid}} k)_{i,j} = \sum_{m=i}^{i+H_k-1} \sum_{n=j}^{j+W_k-1} x_{m,n} k_{i-m+H_k, j-n+W_k}$$

$$1 \times 40 + 0.5 \times 0 + 0.25 \times 0 + 0 \times 40 = 40$$

kernel is a flipped version of a filter

0	80	40
20	40	0
0	0	40

\ast

0	.25
.5	1

$=$

45	110
40	40

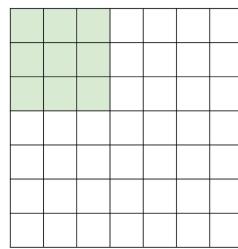
\bar{k} kernel

image x

41

Closer look at spatial dimensions

7



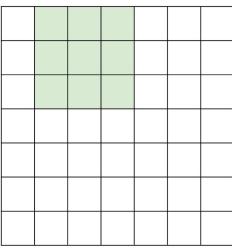
7x7 input (spatially)
assume 3x3 filter

7

Slide credit: Fei-Fei Li, Justin Johnson, Serena Yeung 42

Closer look at spatial dimensions

7



7x7 input (spatially)
assume 3x3 filter

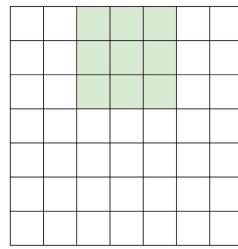
Stride is convolutional
“step size”: here, = 1

7

Slide credit: Fei-Fei Li, Justin Johnson, Serena Yeung 43

Closer look at spatial dimensions

7



7x7 input (spatially)
assume 3x3 filter

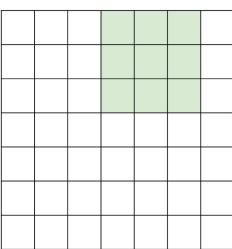
Stride is convolutional
“step size”: here, = 1

7

Slide credit: Fei-Fei Li, Justin Johnson, Serena Yeung 44

Closer look at spatial dimensions

7



7x7 input (spatially)
assume 3x3 filter

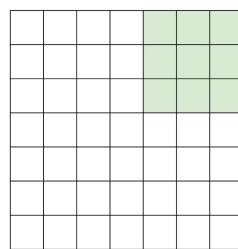
Stride is convolutional
“step size”: here, = 1

7

Slide credit: Fei-Fei Li, Justin Johnson, Serena Yeung 45

Closer look at spatial dimensions

7



7x7 input (spatially)
assume 3x3 filter

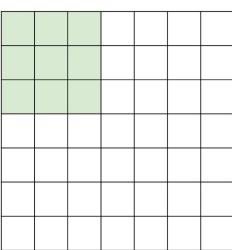
Stride is convolutional
“step size”: here, = 1

7

Slide credit: Fei-Fei Li, Justin Johnson, Serena Yeung 46

Closer look at spatial dimensions

7



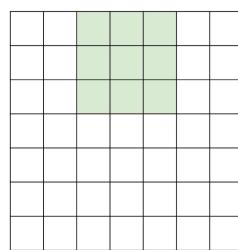
7x7 input (spatially)
assume 3x3 filter
applied with stride 2

7

Slide credit: Fei-Fei Li, Justin Johnson, Serena Yeung 47

Closer look at spatial dimensions

7

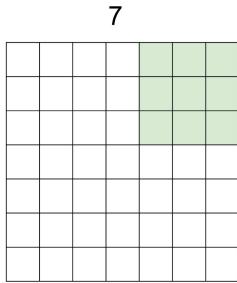


7x7 input (spatially)
assume 3x3 filter
applied with stride 2

7

Slide credit: Fei-Fei Li, Justin Johnson, Serena Yeung 48

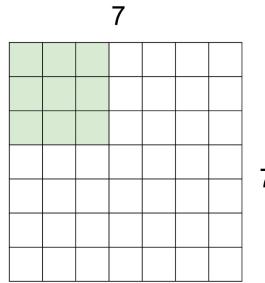
Closer look at spatial dimensions



7x7 input (spatially)
assume 3x3 filter
applied **with stride 2**
=> 3x3 output!

Slide credit: Fei-Fei Li, Justin Johnson, Serena Yeung 49

Closer look at spatial dimensions

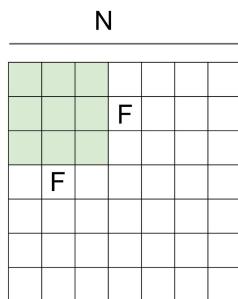


7x7 input (spatially)
assume 3x3 filter
applied **with stride 3?**

doesn't fit!
cannot apply 3x3 filter on
7x7 input with stride 3.

Slide credit: Fei-Fei Li, Justin Johnson, Serena Yeung 51

Closer look at spatial dimensions

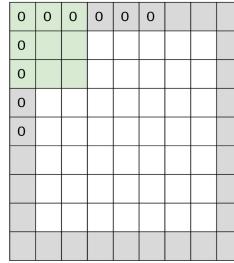


Output size:
 $(N - F) / \text{stride} + 1$

e.g. $N = 7, F = 3$:
stride 1 => $(7 - 3)/1 + 1 = 5$
stride 2 => $(7 - 3)/2 + 1 = 3$
stride 3 => $(7 - 3)/3 + 1 = 2.33 :)$

Slide credit: Fei-Fei Li, Justin Johnson, Serena Yeung 52

In practice: Common to zero pad the border

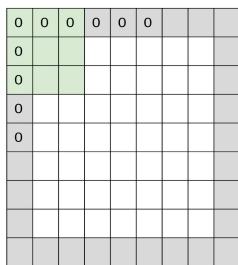


e.g. input 7x7
3x3 filter, applied with **stride 1**
pad with 1 pixel border => what is the output?

(recall):
 $(N - F) / \text{stride} + 1$

Slide credit: Fei-Fei Li, Justin Johnson, Serena Yeung 53

In practice: Common to zero pad the border

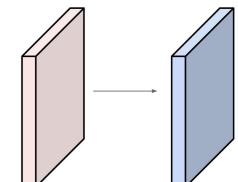


e.g. input 7x7
3x3 filter, applied with **stride 1**
pad with 1 pixel border => what is the output?

7x7 output!
in general, common to see CONV layers with
stride 1, filters of size $F \times F$, and zero-padding with
 $(F-1)/2$. (will preserve size spatially)
e.g. $F = 3 \Rightarrow$ zero pad with 1
 $F = 5 \Rightarrow$ zero pad with 2
 $F = 7 \Rightarrow$ zero pad with 3

Slide credit: Fei-Fei Li, Justin Johnson, Serena Yeung 54

Example



Convolution layer requires 4 hyper-parameters:
number of filters, **their spatial dimension**, **stride**, **pad**

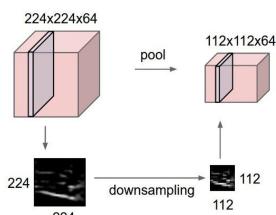
Input volume: 32x32x3
10 5x5 filters with stride 1, pad 2

Output volume size:
 $(32+2*2-5)/1+1 = 32$ spatially, so
32x32x10

Slide credit: Fei-Fei Li, Justin Johnson, Serena Yeung 55

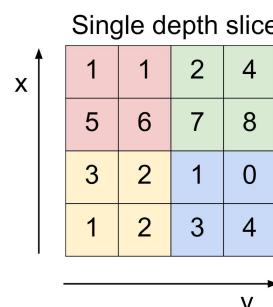
Pooling

- Makes the representations smaller and more manageable
- Operates over each activation map independently

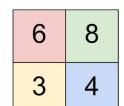


Slide credit: Fei-Fei Li, Justin Johnson, Serena Yeung 57

Max Pooling

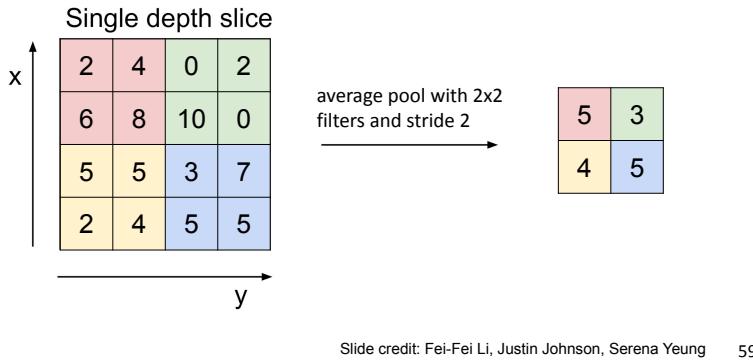


max pool with 2x2 filters
and stride 2



Slide credit: Fei-Fei Li, Justin Johnson, Serena Yeung 58

Average Pooling



Classic ConvNet Architecture

Input

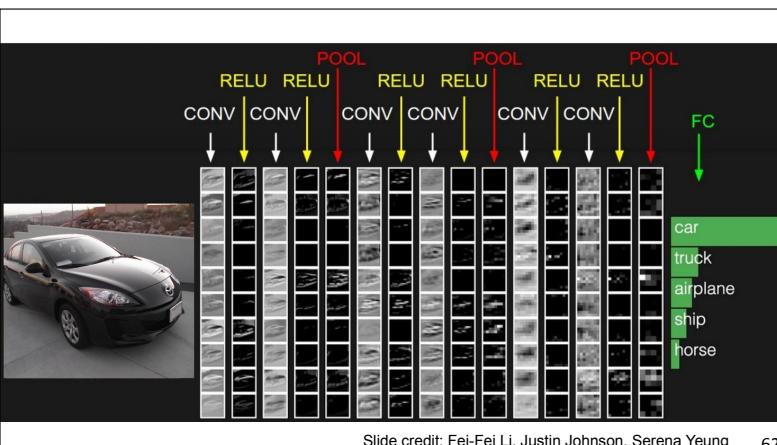
Conv blocks

- Convolution + non-linear activation (relu)
- Convolution + non-linear activation (relu)
- ...
- Maxpooling 2x2

Output layer (e.g. for classification)

- Fully connected layers
- Softmax

61



ConvNetJS demo: training on CIFAR-10

ConvNetJS CIFAR-10 demo

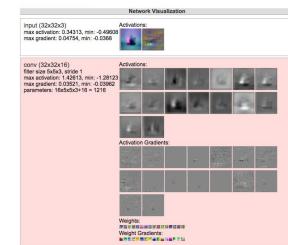
Description

This demo trains a Convolutional Neural Network on the CIFAR-10 dataset in your browser, with nothing but Javascript. The state of the art on this dataset is about 90% accuracy and human performance is at about 94% (not including the 6% of images that are labeled as 'unknowable'). You can upload your own original files (python version) into batches of images that can be easily loaded into page DOM with img tags.

This dataset is more difficult and it takes longer to train a network. Data augmentation includes random flipping and random image shifts by up to 2px horizontally and vertically.

By default, in this demo we're using Adadelta which is one of per-parameter adaptive step size methods, so we don't have to worry about changing learning rates or momentum over time. However, I still included the text fields for changing these if you'd like to play around with SGD+Momentum trainer.

Report questions/bugs/suggestions to @karpathy.

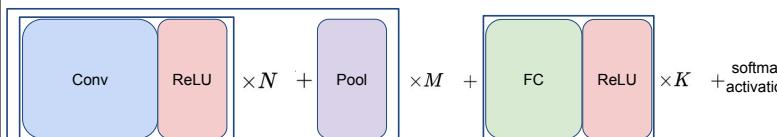


<http://cs.stanford.edu/people/karpathy/convnetjs/demo/cifar10.html>

Slide credit: Fei-Fei Li, Justin Johnson, Serena Yeung 63

Summary

- ConvNets stack CONV, POOL, FC layers
- Trend towards smaller filters and deeper architectures
- Trend towards getting rid of POOL/FC layers (just CONV)
- Typical architectures look like



where N is usually up to ~5, M is large, 0 <= K <= 2.

- but recent advances such as ResNet/GoogLeNet challenge this paradigm

Slide credit: Fei-Fei Li, Justin Johnson, Serena Yeung 64

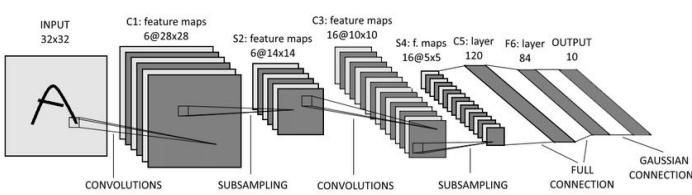
Outline

- Optimization
- CNN basics
- Examples of CNN Architectures
- Applications of CNN

65

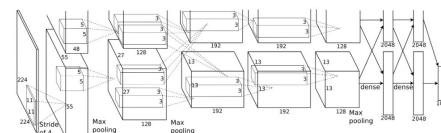
LeCun et al. 1989

Gradient-based learning applied to document recognition [LeCun, Bottou, Bengio, Haffner, 1989]



66

AlexNet

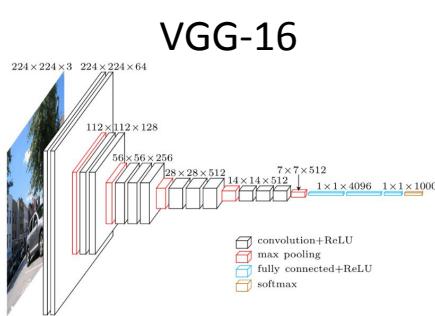


First conv layer: kernel 11x11x3x96 stride 4

- Kernel shape: (11, 11, 3, 96)
- Output shape: (55, 55, 96)
- Number of parameters: 34,944
- Equivalent MLP parameters: 43.7×10^6

Simplified version of Krizhevsky, Alex, Sutskever, and Hinton. "Imagenet classification with deep convolutional neural networks." NIPS 2012

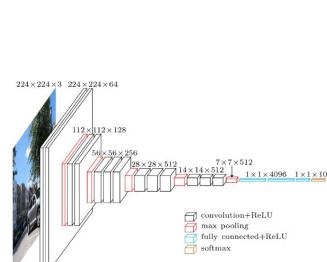
67



Simonyan, Karen, and Zisserman. "Very deep convolutional networks for large-scale image recognition." (2014)

68

VGG-16



Lots of 3×3 convolution layers: more non-linearity than a single 7×7 conv layer

ConvNet Configuration				
A	B - LRN	C	D	E
11 weight layers	11 weight layers	13 weight layers	16 weight layers	16 weight layers
input (224×224 RGB image)				
conv3-64	conv3-64	conv3-64	conv3-64	conv3-64
LRN	conv3-64	conv3-64	conv3-64	conv3-64
conv3-128	conv3-128	conv3-128	conv3-128	conv3-128
	conv3-128	conv3-128	conv3-128	conv3-128
conv3-256	conv3-256	conv3-256	conv3-256	conv3-256
conv3-256	conv3-256	conv3-256	conv3-256	conv3-256
	conv1-256	conv1-256	conv1-256	conv1-256
conv3-512	conv3-512	conv3-512	conv3-512	conv3-512
conv3-512	conv3-512	conv3-512	conv3-512	conv3-512
	conv1-512	conv1-512	conv1-512	conv1-512
conv3-512	conv3-512	conv3-512	conv3-512	conv3-512
conv3-512	conv3-512	conv3-512	conv3-512	conv3-512
	conv1-512	conv1-512	conv1-512	conv1-512
maxpool				
FC-4096				
FC-1000				
softmax				

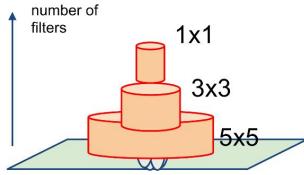
69

Very Deep Models

Going deep with convolutions, Szegedy et al.

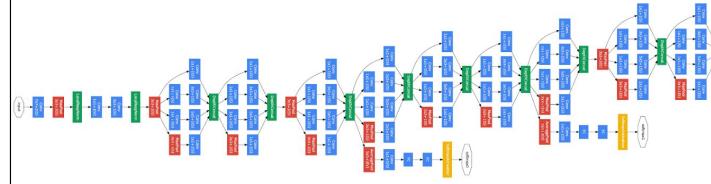
GoogLeNet inception module:

1. Multiple filter scales at each layer
2. Dimensionality reduction to keep computational requirements down



Very Deep Models

Going deep with convolutions, Szegedy et al.

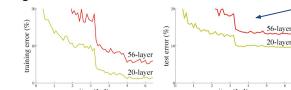


71

Residual Networks

[He, Zhang, Ren, Sun, CVPR 2016]

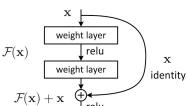
Really, really deep convnets don't train well, E.g. CIFAR10:



deeper networks perform worse due to vanishing gradient

Key idea: introduce "pass through" into each layer

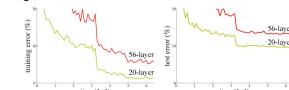
Thus only residual now needs to be learned



Residual Networks

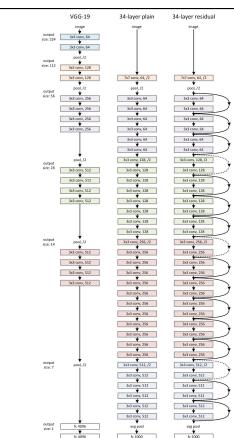
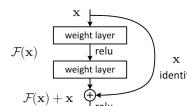
[He, Zhang, Ren, Sun, CVPR 2016]

Really, really deep convnets don't train well, E.g. CIFAR10:



Key idea: introduce "pass through" into each layer

Thus only residual now needs to be learned

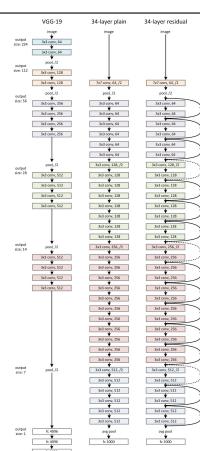


73

Residual Networks

ResNet50 Compared to VGG:

- Superior accuracy in all vision tasks
5.25% top-5 error vs 7.1%
- Less parameters
25M vs 138M
- Computational complexity
3.8B Flops vs 15.3B Flops
- Fully convolutional until the last layer



Deeper is better

ImageNet Large Scale Visual Recognition Challenge (ILSVRC)

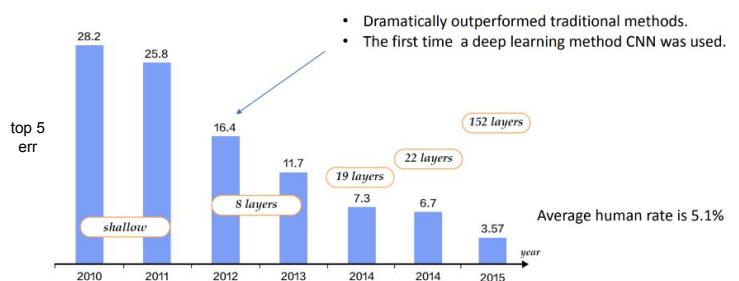
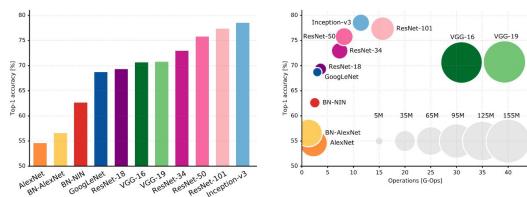


Figure source: http://www.rt-rk.uns.ac.rs/sites/default/files/materijali/predavanja/DL_5.pdf

75

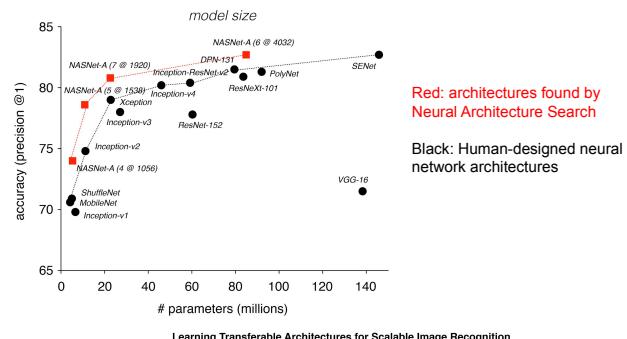
ImageNet classification

Top 1-accuracy, performance and size on ImageNet



Canziani, Paszke, and Culurciello. "An Analysis of Deep Neural Network Models for Practical Applications." (May 2016), 76

Learned architectures surpass human-invented



Red: architectures found by Neural Architecture Search

Black: Human-designed neural network architectures

Learning Transferable Architectures for Scalable Image Recognition
B. Zoph, V. Vasudevan, J. Shlens, Q. Le (2017)

77

Outline

- Optimization
- CNN basics
- Examples of CNN Architectures
- Applications of CNN

76

Pre-trained models

- Training a model on ImageNet from scratch takes **days or weeks**.
- Many models trained on ImageNet and their weights are publicly available!

Transfer learning

- Use pre-trained weights, remove last layers to compute representations of images
- Train a classification model from these features on a new classification task
- The network is used as a generic feature extractor
- Better than handcrafted feature extraction on natural images

78

79

Pre-trained models

- Training a model on ImageNet from scratch takes **days or weeks**.
- Many models trained on ImageNet and their weights are publicly available!

Fine-tuning

Retraining the (some) parameters of the network (given enough data)

- Truncate the last layer(s) of the pre-trained network
- Freeze the remaining layers weights
- Add a (linear) classifier on top and train it for a few epochs
- Then fine-tune the whole network or the few deepest layers
- Use a smaller learning rate when fine tuning

Applications of CNN

- Case study 1: image classification**
- Case study 2: object detection
- Case study 3: segmentation

80

81

Object classification

- Given an input image, classify the object within it
- Input?
- Output? "Cat" (class probability): $\hat{y} = (\hat{p}_1, \hat{p}_2, \dots, \hat{p}_C)$
- Loss? $\mathcal{L}(y, \hat{y}) = -y \log(\hat{y}) = -\sum_{c=1}^C p_c \log(\hat{p}_c)$

Object classification

- Imagenet dataset
- Very large dataset of many object classes
- Enable learning models capable to identify many object features



83

84

Object Classification: Deeper is better

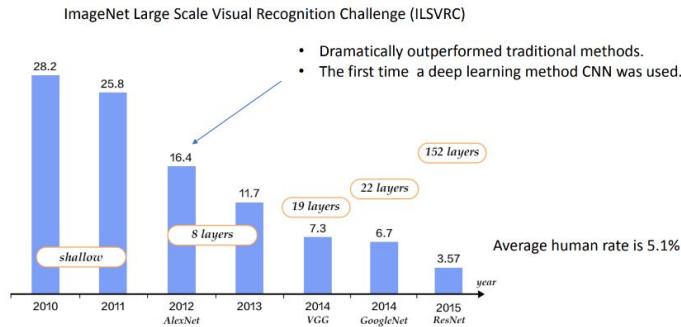


Figure source: http://www.rt-rk.uns.ac.rs/sites/default/files/materijali/predavanja/DL_5.pdf

85

Emotion Classification

- Given an input human face, label the emotion
- Input?
- Output? "fear"
- Loss? Multi-class cross entropy



Examples of emotion categories

Figure from: <https://www.youtube.com/watch?v=wlaR5F30hiU> 86

Emotion Classification



Figure credit: https://github.com/carriga/face_classification

87

Emotion Classification

- Pipeline:
 - Input preprocessing (alignment, normalization, etc)
 - Given the preprocessed data, pass it to a CNN and get output
 - Output label with highest probability is the assigned emotion

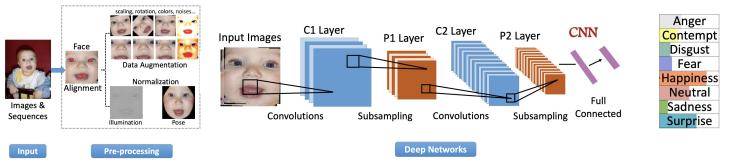


Figure credit: Shan Li and Weihong Deng. Deep Facial Expression Recognition: A Survey.

88

Applications of CNN

- Case study 1: Image Classification
- Case study 2: Object Detection**
- Case study 3: Segmentation

Object Detection

- Given an image, detect all objects and classify
- Input?
- Output?
- Loss? Error in bounding box position? Bounding box scale? Class label?

89

90

Object Detection



Classes = [cat, dog, duck]

Cat ? YES

Dog ? NO

Duck? NO

Object Detection



Classes = [cat, dog, duck]

Cat ? NO

Dog ? NO

Duck? NO

Object Detection



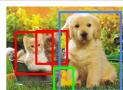
Problem:
Too many positions & scales to test

Solution: If your classifier is fast enough, go for it

Slide credit: Amaia Salvador

93

Object Detection

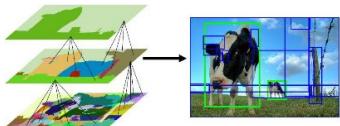


- Two sub-problems:
 - Object locations/size: localize each object with a tight bounding box in the image
 - Object classification: classify the object in each bounding box
- Typically solved by “**proposing candidate bounding boxes**” and “**classifying them**”
- Two main families:
 - A grid in the image where each cell is a proposal (YOLO)
 - Region proposal (Fast RCNN, Faster RCNN)

94

Bounding Box Proposals

- Instead of having a predefined set of box proposals, find them on the image:
 - Selective Search (from pixels, not learnt)
 - Region Proposal Network (RPN, learnt)

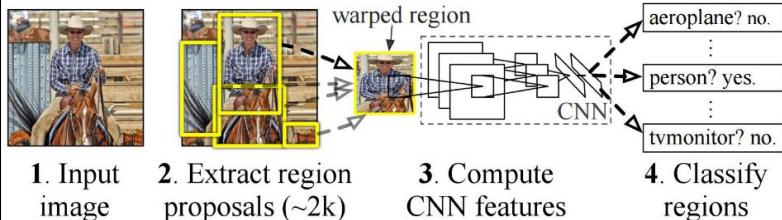


Uijlings, J. R., Van De Sande, K. E., Gevers, T., & Smeulders, A. W. (2013). Selective search for object recognition. *IJCV*, 104(2), 154-171.

95

Region-CNN (RCNN)

For each region proposal, use CNN to predict the categories (with confidence)

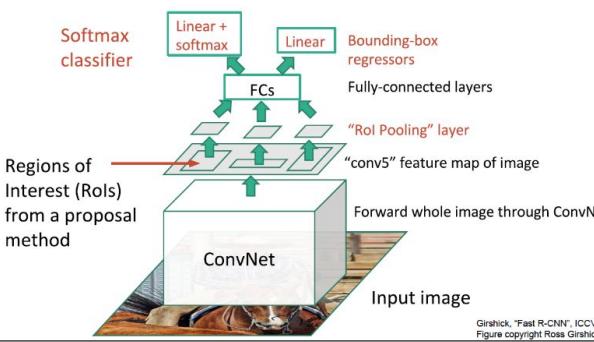


Girshick, R., Donahue, J., Darrell, T., & Malik, J. (2014). Rich feature hierarchies for accurate object detection and semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 580-587).

97

Fast RCNN

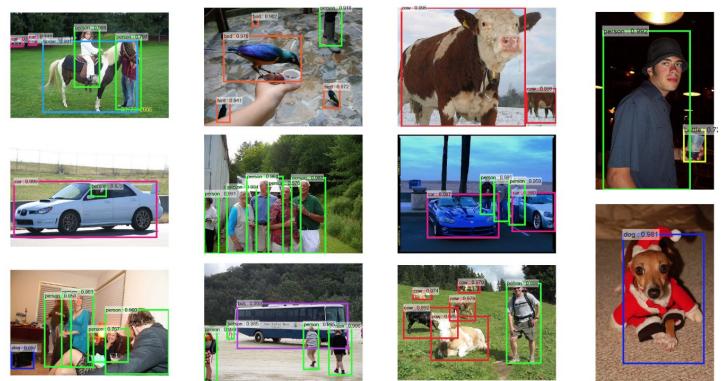
Main idea: calculate convolutional feature map only once and reuse it!



Girshick, "Fast R-CNN", ICCV 2015, Figure copyright Ross Girshick, 2015.

100

Faster RCNN



101

Applications of CNN

- Case study 1: Image Classification
- Case study 2: Object Detection
- Case study 3: Segmentation**

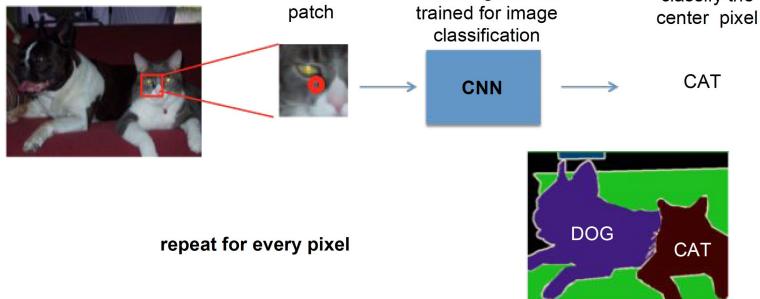
Semantic Segmentation

- Given an image, partition the image into coherent parts
- Input?
- Output? Pixel-level class map
- Loss? Pixel level classification loss

102

103

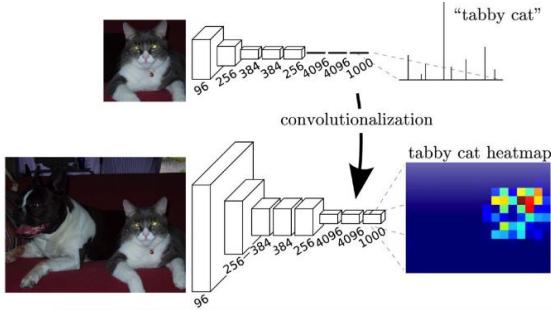
From Classification to Semantic Segmentation



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

From Classification to Semantic Segmentation

- A classification network becoming fully convolutional

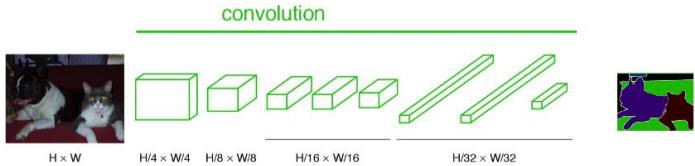


104

105

From Classification to Semantic Segmentation

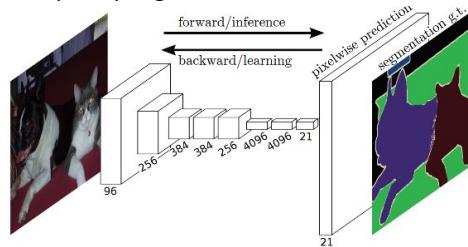
- Dense prediction: fully convolutional, end to end, pixel-to-pixel network
- Problems
 - Output is smaller than input → Add upsampling layers



106

From Classification to Semantic Segmentation

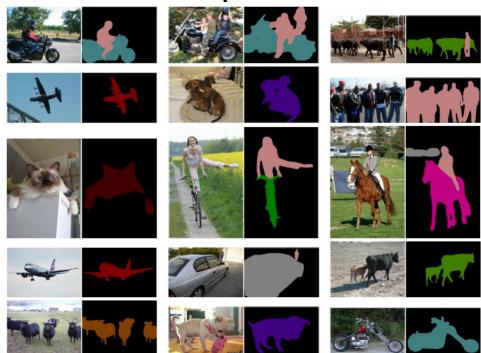
- Dense prediction: fully convolutional, end to end, pixel-to-pixel network
- Learnable upsampling



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

107

Example results



Chen, L. C., Papandreou, G., Kokkinos, I., Murphy, K., & Yuille, A. L. (2018). DeepLab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE transactions on pattern analysis and machine intelligence*, 40(4), 834-848.

108