

Neural net applications for images

Victor Kitov

v.v.kitov@yandex.ru



Convolutional neural networks

- Convolutional neural network:
 - Used for image analysis
 - Consists of a set of convolutional layer / sub-sampling (pooling) layer pairs and several fully connected layers at the end.

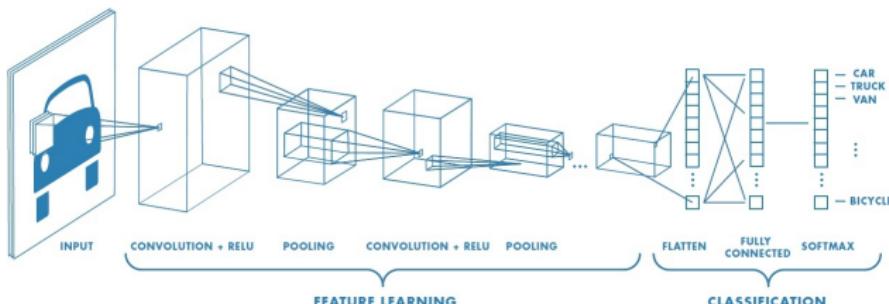
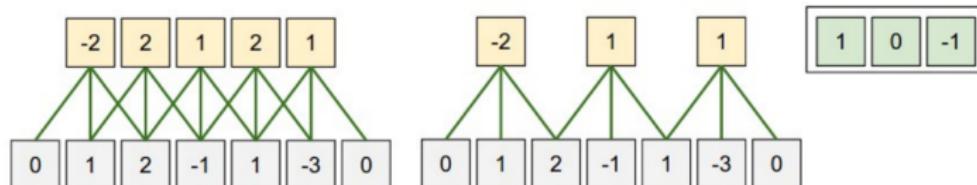


Table of Contents

- 1 ConvNets building blocks
- 2 Case study: ZIP codes recognition
- 3 Major classification architectures
- 4 Data augmentation methods
- 5 Image segmentation

1-D Convolution operation

1-D convolution [$W = 5$, $K = 3$, zero-padded with $P=1$, $S = 1$ (left) and $S = 2$ (right)]



Parameters¹:

- W - length of input
- K - kernel size
- P - amount of padding
- Type of padding (zero, extension, mirror)
- S - stride (offset of kernel)
- D - dilation (offset inside kernel)

¹Depending on these parameters, what would be the size of output layer?

Convolution²

Single layer convolution:

3 ₀	3 ₁	2 ₂	1	0
0 ₀	0 ₁	1 ₂	3	1
3 ₀	1 ₁	2 ₂	2	3
2 ₀	0 ₁	0 ₂	2	2
2 ₀	0 ₁	0 ₂	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

3	3 ₀	2 ₁	1 ₂	0
0	0 ₂	1 ₂	3 ₀	1
3	1 ₀	2 ₁	2 ₂	3
2	0	0	2	2
2	0	0	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

3	3 ₀	2 ₁	1 ₂	0
0	0	1 ₀	3 ₁	1 ₀
3	1	2 ₀	2 ₁	3 ₂
2	0	0	2	2
2	0	0	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

3	3 ₀	2 ₁	1 ₂	0
0 ₀	0 ₁	1 ₂	3	1
3 ₂	1 ₂	2 ₀	2	3
2 ₀	0 ₁	0 ₂	2	2
2 ₀	0 ₁	0 ₂	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

3	3 ₀	2 ₁	1 ₂	0
0	0 ₀	1 ₁	3 ₂	1
3	1 ₂	2 ₂	2 ₀	3
2	0 ₀	0 ₁	2 ₂	2
2	0	0	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

3	3 ₀	2 ₁	1 ₂	0
0	0	1 ₀	3 ₁	1 ₂
3	1	2 ₂	2 ₂	3 ₀
2	0	0 ₀	2 ₁	2 ₂
2	0	0	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

3	3 ₀	2 ₁	1 ₂	0
0 ₀	0 ₁	1 ₃	1	1
3 ₀	1 ₁	2 ₂	2	3
2 ₂	0 ₂	0 ₀	2	2
2 ₀	0 ₁	0 ₂	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

3	3 ₀	2 ₁	1 ₂	0
0	0	1 ₃	1	1
3	1 ₀	2 ₁	2 ₂	3
2	0 ₂	0 ₀	2 ₀	2
2	0 ₀	0 ₁	0 ₂	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

3	3 ₀	2 ₁	1 ₂	0
0	0	1	3	1
3	1	2	2	3
2	0	0	2	0
2	0	0	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

²Illustrations from Dumoulin et al. 2018.

Convolution

Convolution with stride and zero-padding:

0_0	0_1	0_2	0	0	0	0	0
0_2	3_2	3_0	2	1	0	0	
0_0	0_1	0_2	1	3	1	0	
0	3	1	2	2	3	0	
0	2	0	0	2	2	0	
0	2	0	0	0	1	0	
0	0	0	1	0	0	0	0



0	0	0_0	0_1	0_2	0	0
0	3	3_2	2_1	0	0	
0	0	0_0	1_1	3_2	1	0
0	3	1	2	2	3	0
0	2	0	0	2	2	0
0	2	0	0	0	1	0
0	0	0	0	0	0	0



0	0	0	0	0	0_0	0_1	0_2
0	3	3	2	1	0	0	
0	0	0	1	3	1	0_2	
0	3	1	2	2	3	0	
0	2	0	0	2	2	0	
0	2	0	0	0	1	0	
0	0	0	0	0	0	0	0



0_0	0_1	0_2	0	0	0	0	0
0_2	3_3	2_1	0	0			
0_0	0_1	0_2	1	3	1	0	
0	3	1	2	2	3	0	
0	2	0	0	2	2	0	
0	2	0	0	0	1	0	
0	0	0	1	0	0	0	0



0	0	0	0	0	0	0	0
0	3	3	2	1	0	0	
0	0	0	1	3	1	0	
0	3	1	2	2	3	0	
0	2	0	0	2	2	0	
0	2	0	0	0	1	0	
0	0	0	0	0	0	0	0



0	0	0	0	0	0	0	0
0	3	3	2	1	0	0	
0	0	0	1	3	1	0_2	
0	3	1	2	2	3	0	
0	2	0	0	2	2	0	
0	2	0	0	0	1	0	
0	0	0	0	0	0	0	0



0_0	0_1	0_2	0	0	0	0	0
0_2	3_3	2_0	0	2	2	0	
0_0	2_2	0_0	0	2	2	0	
0	2	0	0	0	0	1	0
0	0	1	0	0	0	0	0



0	0	0	0	0	0	0	0
0	3	3	2	1	0	0	
0	0	0	1	3	1	0	
0	3	1	2	2	3	0	
0	2	0	0	2	2	0	
0	2	0	0	0	1	0	
0	0	0	0	0	0	0	0



0	0	0	0	0	0	0	0
0	3	3	2	1	0	0	
0	0	0	1	3	1	0_2	
0	3	1	2	2	3	0	
0	2	0	0	2	2	0	
0	2	0	0	0	1	0	
0	0	0	0	0	0	0	0

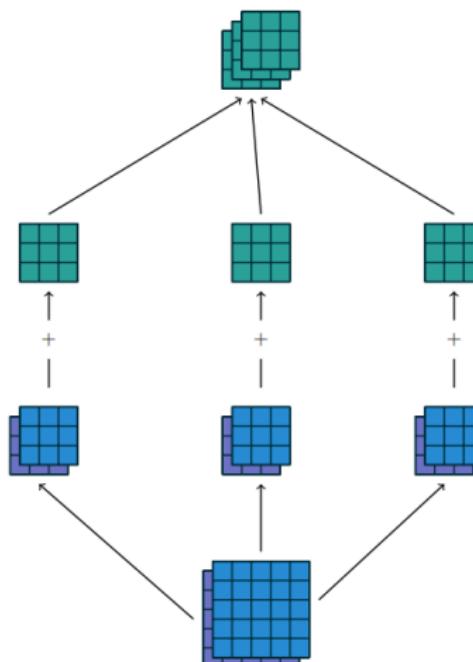


Padding

- Stride: to decrease dimensionality.
- Padding: to increase dimensionality.
- Padding types:
 - zero padding
 - same padding
 - mirror padding

Convolution

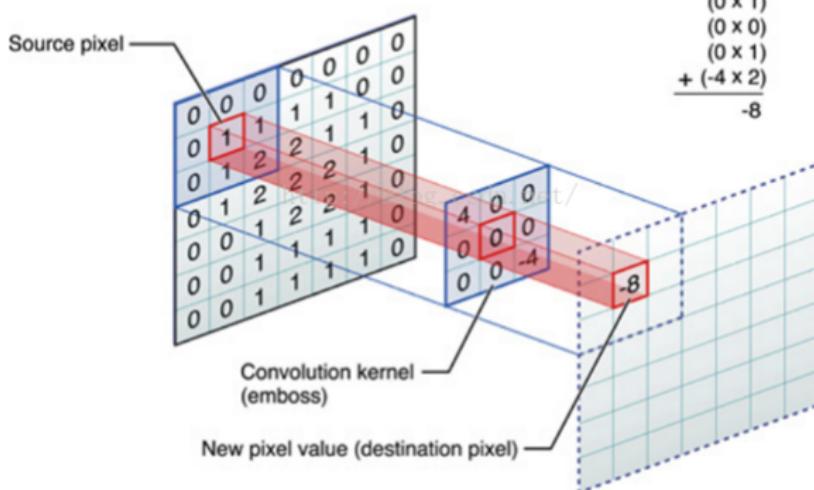
2 layer input, 3 layer output convolution:



Convolution operation

2-D convolution

Center element of the kernel is placed over the source pixel. The source pixel is then replaced with a weighted sum of itself and nearby pixels.



Comments

- Comments on convolution:
 - Locality: each neuron in the feature map takes output from small neighborhood of input layer neurons
 - Equivalence: the same transformation is applied by each neuron in the feature map
 - obtained by constraining sets of weights to each feature map layer neuron to be equal
 - This is feature extraction from a patch

Average pooling

Average 3x3 pooling:

3	3	2	1	0
0	0	1	3	1
3	1	2	2	3
2	0	0	2	2
2	0	0	0	1

1.7	1.7	1.7
1.0	1.2	1.8
1.1	0.8	1.3

3	3	2	1	0
0	0	1	3	1
3	1	2	2	3
2	0	0	2	2
2	0	0	0	1

1.7	1.7	1.7
1.0	1.2	1.8
1.1	0.8	1.3

3	3	2	1	0
0	0	1	3	1
3	1	2	2	3
2	0	0	2	2
2	0	0	0	1

1.7	1.7	1.7
1.0	1.2	1.8
1.1	0.8	1.3

3	3	2	1	0
0	0	1	3	1
3	1	2	2	3
2	0	0	2	2
2	0	0	0	1

1.7	1.7	1.7
1.0	1.2	1.8
1.1	0.8	1.3

3	3	2	1	0
0	0	1	3	1
3	1	2	2	3
2	0	0	2	2
2	0	0	0	1

1.7	1.7	1.7
1.0	1.2	1.8
1.1	0.8	1.3

3	3	2	1	0
0	0	1	3	1
3	1	2	2	3
2	0	0	2	2
2	0	0	0	1

1.7	1.7	1.7
1.0	1.2	1.8
1.1	0.8	1.3

3	3	2	1	0
0	0	1	3	1
3	1	2	2	3
2	0	0	2	2
2	0	0	0	1

1.7	1.7	1.7
1.0	1.2	1.8
1.1	0.8	1.3

3	3	2	1	0
0	0	1	3	1
3	1	2	2	3
2	0	0	2	2
2	0	0	0	1

1.7	1.7	1.7
1.0	1.2	1.8
1.1	0.8	1.3

3	3	2	1	0
0	0	1	3	1
3	1	2	2	3
2	0	0	2	2
2	0	0	0	1

1.7	1.7	1.7
1.0	1.2	1.8
1.1	0.8	1.3

Max pooling

Max 3x3 pooling:

3	3	2	1	0
0	0	1	3	1
3	1	2	2	3
2	0	0	2	2
2	0	0	0	1

3.0	3.0	3.0
3.0	3.0	3.0
3.0	2.0	3.0

3	3	2	1	0
0	0	1	3	1
3	1	2	2	3
2	0	0	2	2
2	0	0	0	1

3.0	3.0	3.0
3.0	3.0	3.0
3.0	2.0	3.0

3	3	2	1	0
0	0	1	3	1
3	1	2	2	3
2	0	0	2	2
2	0	0	0	1

3.0	3.0	3.0
3.0	3.0	3.0
3.0	2.0	3.0

3	3	2	1	0
0	0	1	3	1
3	1	2	2	3
2	0	0	2	2
2	0	0	0	1

3.0	3.0	3.0
3.0	3.0	3.0
3.0	2.0	3.0

3	3	2	1	0
0	0	1	3	1
3	1	2	2	3
2	0	0	2	2
2	0	0	0	1

3.0	3.0	3.0
3.0	3.0	3.0
3.0	2.0	3.0

3	3	2	1	0
0	0	1	3	1
3	1	2	2	3
2	0	0	2	2
2	0	0	0	1

3.0	3.0	3.0
3.0	3.0	3.0
3.0	2.0	3.0

3	3	2	1	0
0	0	1	3	1
3	1	2	2	3
2	0	0	2	2
2	0	0	0	1

3.0	3.0	3.0
3.0	3.0	3.0
3.0	2.0	3.0

3	3	2	1	0
0	0	1	3	1
3	1	2	2	3
2	0	0	2	2
2	0	0	0	1

3.0	3.0	3.0
3.0	3.0	3.0
3.0	2.0	3.0

3	3	2	1	0
0	0	1	3	1
3	1	2	2	3
2	0	0	2	2
2	0	0	0	1

3.0	3.0	3.0
3.0	3.0	3.0
3.0	2.0	3.0

Table of Contents

- 1 ConvNets building blocks
- 2 Case study: ZIP codes recognition
- 3 Major classification architectures
- 4 Data augmentation methods
- 5 Image segmentation

Case study (due to Hastie et al. The Elements of Statistical Learning)

ZIP code recognition task



Neural network structures

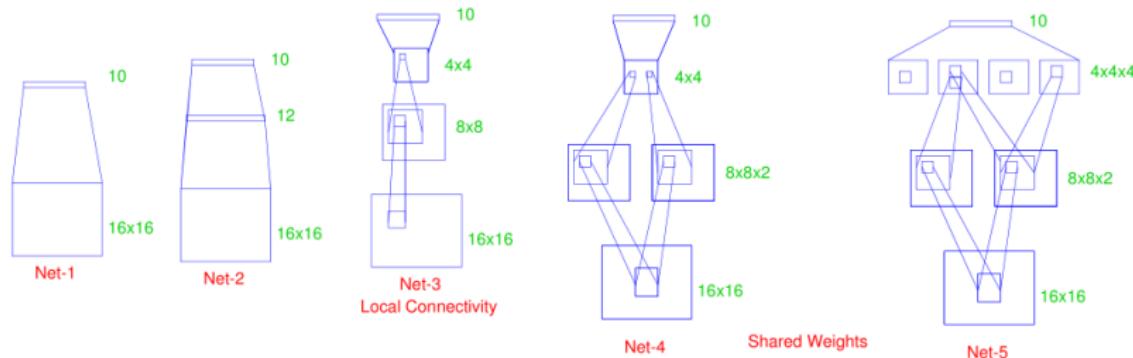
Net1: no hidden layer

Net2: 1 hidden layer, 12 hidden units fully connected

Net3: 2 hidden layers, locally connected

Net4: 2 hidden layers, locally connected with weight sharing

Net5: 2 hidden layers, locally connected, 2 levels of weight sharing



Results

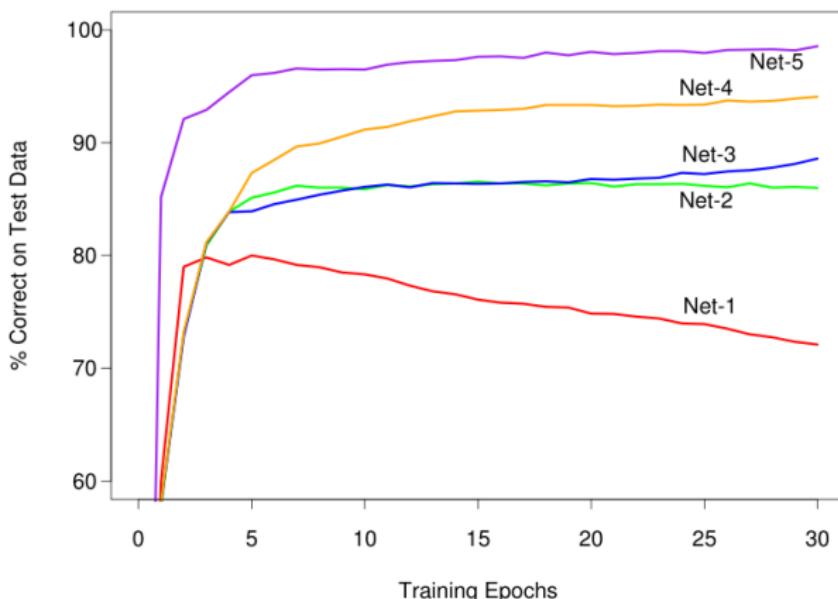
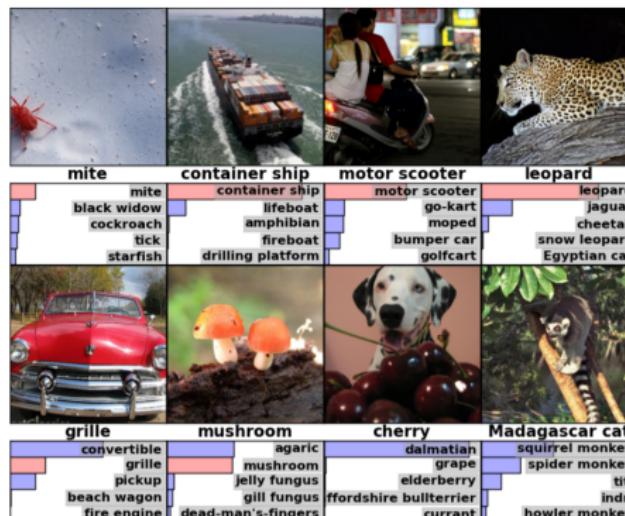


Table of Contents

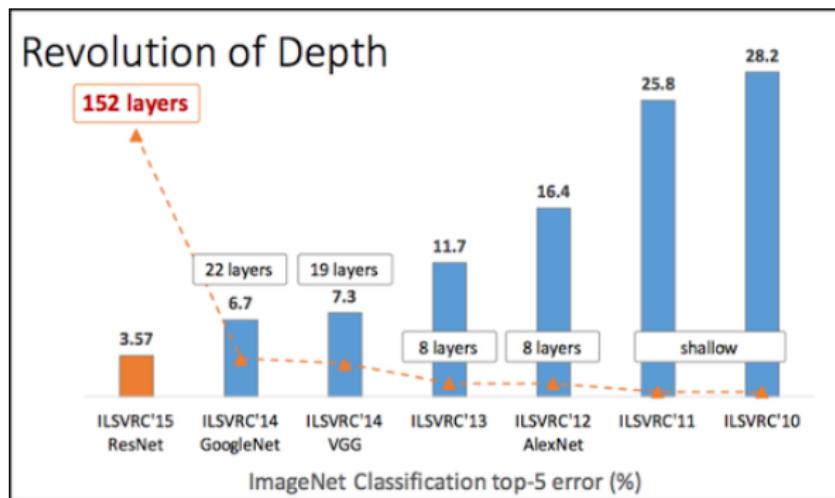
- 1 ConvNets building blocks
- 2 Case study: ZIP codes recognition
- 3 Major classification architectures
- 4 Data augmentation methods
- 5 Image segmentation

ImageNet classification challenge

- 1000 unambiguous classes (including 120 dog breeds!).
- >1 million hand annotated images.
- Classifiers evaluated by top-5 accuracy
 - is the true class present among top-5 predictions?



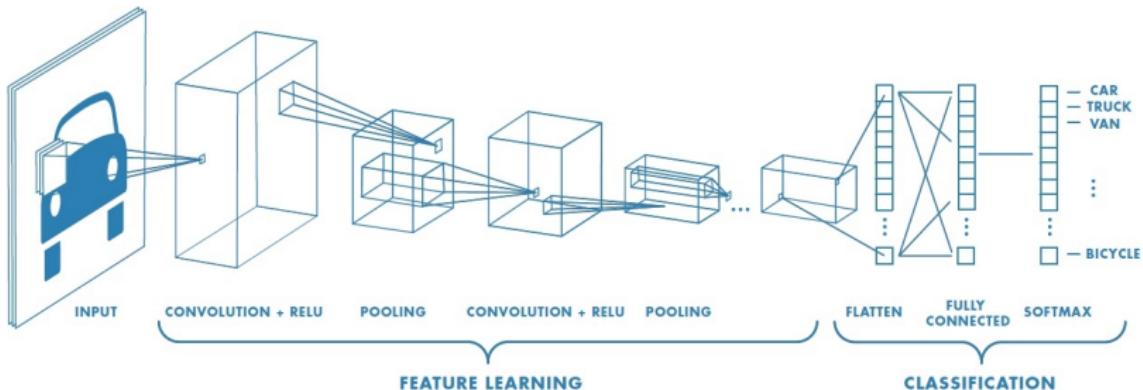
ImageNet challenge progress



- Starting from 2012 - triumph of deep convolutional networks.
- Human performance 3-15% (depending on acquaintance with the classes).³

³ Andrew Karpathy human test.

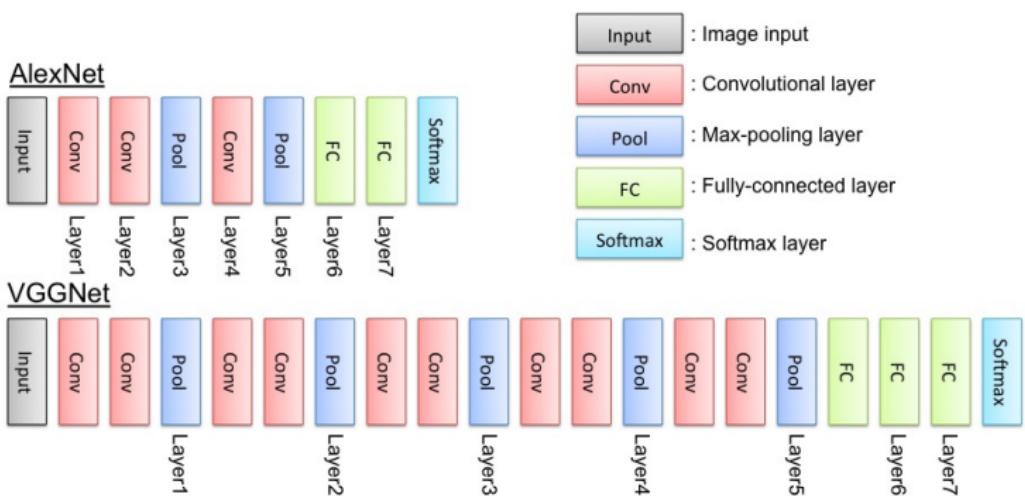
Convolution network



- Later layers learn more and more abstract features.
- Receptive field (in terms of original image) of neurons from deeper layers is wider.

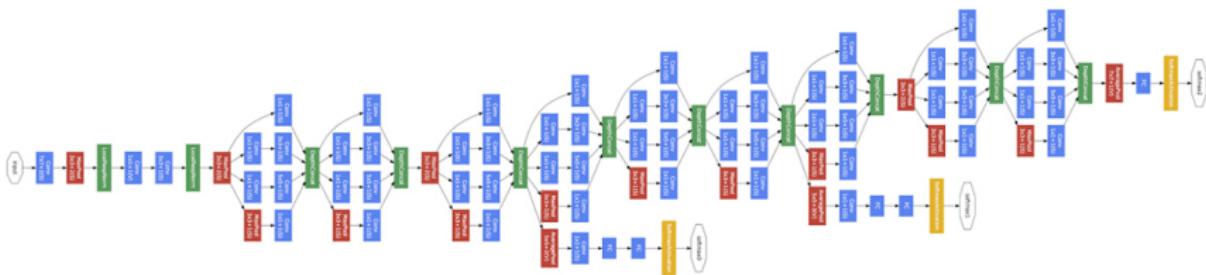
Major CNN architectures

AlexNet vs VGG



Each layer is followed by ReLu non-linearity.

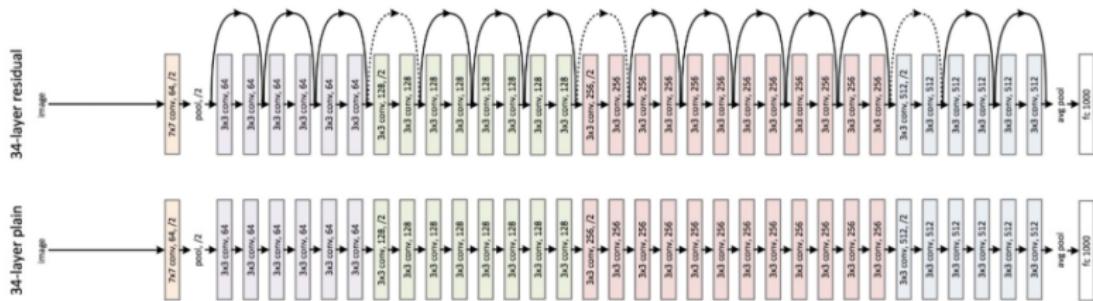
GoogleNet



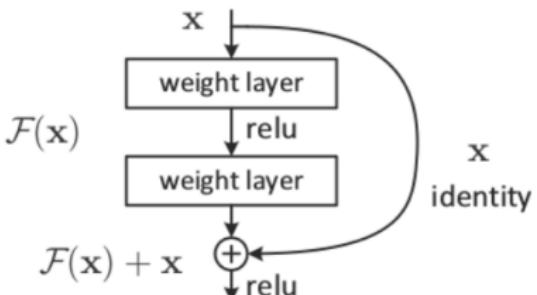
- add intermediate outputs during training
- reduce computation and # parameters by 1×1 convolutions

ResNet

ResNet vs. plain network



ResNet building block:



Skip identity connections allow:

- better propagate gradient backwards
- allow more natural initialization of weight layers
 - so that they are almost constant

VGG network⁴

Winners of ImageNet Challenge 2014!

- Data preprocessing - extract mean for R,G,B channels from all pixel intensities.
- Key idea: gradually reduce size and increase receptive field:
 - Filters with a very small receptive field: 3×3 , 1×1 .
 - Padding to keep original size (1 pixel for 3×3 conv)
 - Stride 1 for conv
 - Max-pooling is performed over a 2×2 pixel window, with stride 2.
- All hidden layers are followed by ReLU



⁴2015 - Very deep convolutional networks for large-scale image recognition - Simonyan et al.

VGG details

- Optimization: SGD with momentum
 - learning rate decreased 3 times.
- Parallelization over minibatches
 - gradients are then averaged
- First 2 fully connected layers:
 - weight decay regularization
 - dropout regularization
- Train more shallow net, then with learned weights initialize deeper network.
- Dataset augmentation: random scaling, cropping.

Table of Contents

- 1 ConvNets building blocks
- 2 Case study: ZIP codes recognition
- 3 Major classification architectures
- 4 Data augmentation methods
 - Padding unknown data
- 5 Image segmentation

Definition of data augmentation

Data augmentation - extend training set $\{(x_n, y_n)\}_{n=1}^N$ with $\{(f_\theta(x_m), y_m)\}_{m=1}^M$ where

- $f_\theta(x)$ is label preserving transformation
 - in case of image: horizontal flipping, cropping, scaling, small rotation, small change in brightness, contrast, saturation, hue.
 - it may combine several transformations.
- θ are randomly sampled parameters of the transformation
- M is amount of new samples.

Applications of data augmentation

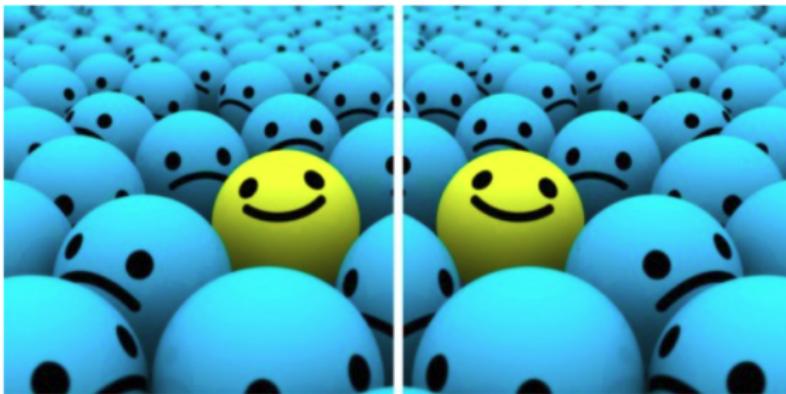
- Enlarge training set - estimate network parameters better.
 - feasible to estimate many parameters
 - less overfitting
- Enforce prediction invariance to transformation.
 - transfer learning: $f_\theta(\cdot)$ maps training domain to target domain.
 - e.g. have labelled *day-time* images, but need *night-time* classification.

#objects vs. #parameters⁵

	VGGNet	DeepVideo	GNMT
Used For	Identifying Image Category	Identifying Video Category	Translation
Input	Image 	Video 	English Text 
Output	1000 Categories	47 Categories	French Text
Parameters	140M	~100M	380M
Data Size	1.2M Images with assigned Category	1.1M Videos with assigned Category	6M Sentence Pairs, 340M Words
Dataset	ILSVRC-2012	Sports-1M	WMT'14

⁵Image source.

Horizontal flipping⁶



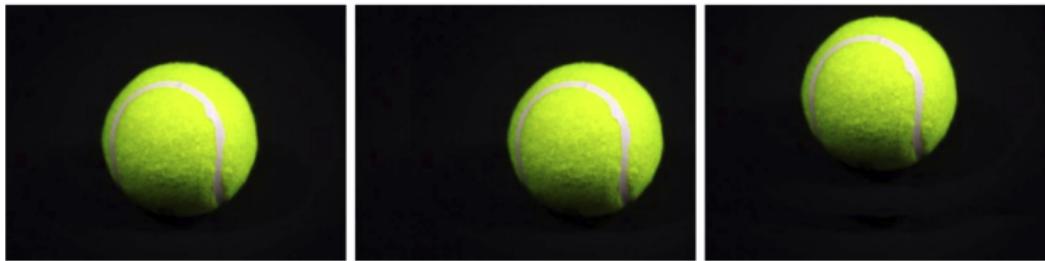
⁶ Image source.

Rotations⁷



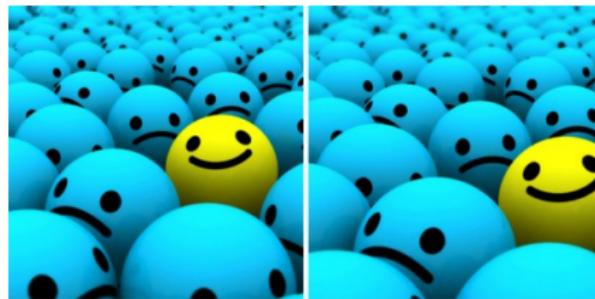
⁷ Image source.

Translations⁸



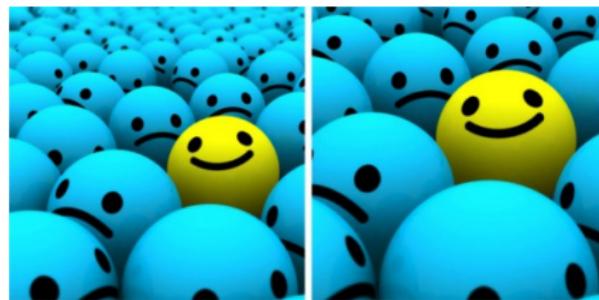
⁸Image source.

Scaling⁹



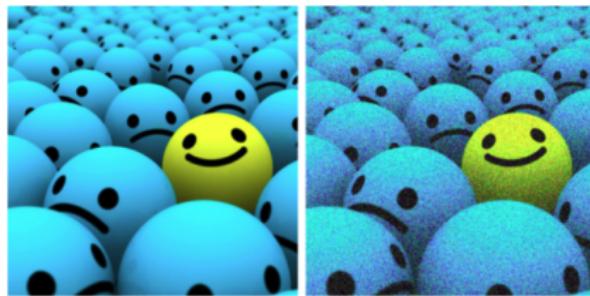
⁹ Image source.

Cropping¹⁰



¹⁰ Image source.

Adding noise (Gaussian, salt&pepper, etc.)¹¹



¹¹ Image source.

Other augmentation methods

- Randomly adjust:
 - brightness (average level)
 - contrast (std. deviation of levels)
 - hue (add random offsets to RGB channels)
- Fill random patches with constant or noise
 - like dropout on input

Data augmentation methods

Padding unknown data

④ Data augmentation methods

- Padding unknown data

Data augmentation methods

Padding unknown data

Transformation uncovers unknown data¹²

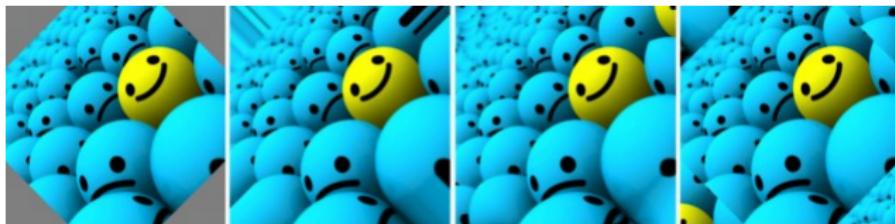


¹² Image source.

Data augmentation methods

Padding unknown data

Methods to pad unknown data¹³



Unknown data padding methods (left to right):

- constant
- extend edge values
- mirror
- wrap around

¹³ [Image source](#).

Table of Contents

- 1 ConvNets building blocks
- 2 Case study: ZIP codes recognition
- 3 Major classification architectures
- 4 Data augmentation methods
- 5 Image segmentation

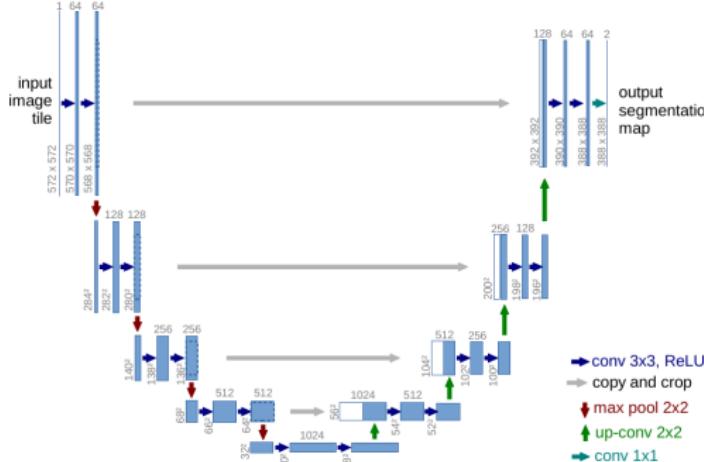
Image segmentation¹⁴



- Segmentation - classification of every pixel of the image.
- Applications:
 - surveillance systems, autonomous driving, image classification, activity recognition on videos, etc.
- Model needs:
 - high level features to reconstruct object type
 - low level features to reconstruct boundaries

¹⁴ Picture source.

U-net architecture¹⁵



Horizontal numbers = # [channels]; vertical numbers = spatial size.
 White blocks - copied output of earlier layers; up-conv - rescaling & convolution.

¹⁵Ronneberger et al [2015].

Discussion

Key ideas of U-net:

- preserve spatial info at each layer
 - use only convolution, pooling, scaling.
 - don't use vectorization & fully connected layers
- 1st half - encoder; 2nd half - decoder.
- Encoder aggregates wider and wider local information
 - creating more abstract features
- Decoder reconstructs local information from
 - more abstract features (green input on figure)
 - lower level features (gray input on figure)