CS60010: Deep Learning

Sudeshna Sarkar

Spring 2018

Residual Net

• Deep Residual Learning for Image Recognition, Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun (CVPR 2016)

The deeper, the better?

- The deeper network can cover more complex problems
 - Receptive field size \uparrow
 - Non-linearity 个
- However, training the deeper network is more difficult because of vanishing/exploding gradients problem



ImageNet Classification top-5 error (%)

Deep NN

- Escape from few layers
 - ReLU for solving gradient vanishing problem
 - Dropout ...
- Escape from 10 layers
 - Normalized initialization
 - Intermediate normalization layers
- Escape from 100 layers
 - Residual network

Plain Network

- Plain nets: stacking 3x3 conv layers
- 56-layer net has higher training error and test error than 20layers net



Plain Network

- "Overly deep" plain nets have higher training error
- A general phenomenon, observed in many datasets



The residual module

- Introduce *skip* or *shortcut* connections
- Make it easy for network layers to represent the identity mapping



Residual block

- If identity were optimal, easy to set weights as 0
- If optimal mapping is closer to identity, easier to find small fluctuations
 - -> Appropriate for treating perturbation as keeping a base information

• Difference between an original image and a changed image



can treat perturbation

Deeper ResNets have lower training error



- Residual block
 - Very simple
 - Parameter-free







Performances increase absolutely

task	2nd-place winner	MSRA	margin (relative)
ImageNet Localization (top-5 error)	12.0	9.0	27%
ImageNet Detection (mAP@.5)	53.6 abso 8.5%	plute 62.1	16%
COCO Detection (mAP@.5:.95)	33.5	37.3	11%
COCO Segmentation (mAP@.5:.95)	25.1	28.2	12%

- Based on ResNet-101
- Existing techniques can use residual networks or features from it

Summary: ILSVRC 2012-2015

Team	Year	Place	Error (top-5)	External data
SuperVision – Toronto (AlexNet, 7 layers)	2012	-	16.4%	no
SuperVision	2012	1st	15.3%	ImageNet 22k
Clarifai – NYU (7 layers)	2013	-	11.7%	no
Clarifai	2013	1st	11.2%	ImageNet 22k
VGG – Oxford (16 layers)	2014	2nd	7.32%	no
GoogLeNet (19 layers)	2014	1st	6.67%	no
ResNet (152 layers)	2015	1st	3.57%	
Human expert*			5.1%	

http://karpathy.github.io/2014/09/02/what-i-learned-from-competing-against-a-convnet-on-imagenet/

Design principles

- Reduce filter sizes (except possibly at the lowest layer), factorize filters aggressively
- Use 1x1 convolutions to reduce and expand the number of feature maps judiciously
- Use skip connections and/or create multiple paths through the network

Reading list

- <u>https://culurciello.github.io/tech/2016/06/04/nets.html</u>
- Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, <u>Gradient-based learning applied to document</u> recognition, Proc. IEEE 86(11): 2278–2324, 1998.
- A. Krizhevsky, I. Sutskever, and G. Hinton, <u>ImageNet Classification with Deep Convolutional</u> <u>Neural Networks</u>, NIPS 2012
- M. Zeiler and R. Fergus, <u>Visualizing and Understanding Convolutional Networks</u>, ECCV 2014
- K. Simonyan and A. Zisserman, <u>Very Deep Convolutional Networks for Large-Scale Image</u> <u>Recognition</u>, ICLR 2015
- M. Lin, Q. Chen, and S. Yan, <u>Network in network</u>, ICLR 2014
- C. Szegedy et al., <u>Going deeper with convolutions</u>, CVPR 2015
- C. Szegedy et al., <u>Rethinking the inception architecture for computer vision</u>, CVPR 2016
- K. He, X. Zhang, S. Ren, and J. Sun, <u>Deep Residual Learning for Image Recognition</u>, CVPR 2016

"You need a lot of a data if you want to train/use CNNs"

Fei-Fei Li & Andrej Karpathy

Lecture 7 -

"You need a lot of a data if you want to train/use CNNs"



The Unreasonable Effectiveness of Deep Features





Low-level: Pool1

High-level: FC6

Classes separate in the deep representations and transfer to many tasks. [DeCAF] [Zeiler-Fergus]

Can be used as a generic feature

("CNN code" = 4096-D vector before classifier)



nearest neighbors in the "code" space

query image

Fei-Fei Li & Andrej Karpathy

Lecture 7 -

Can be used as a generic feature ("CNN code" = 4096-D vector before classifier)



query image

nearest neighbors in the "code" space

Fei-Fei Li & Andrej Karpathy

Lecture 7 -

Transfer Learning with CNNs

image	1. Train on
conv-64	Imagenet
conv-64	
maxpool	
conv-128	
conv-128	
maxpool	
conv-256	
conv-256	
maxpool	
conv-512	
conv-512	
maxpool	
conv-512	
conv-512	
maxpool	
FC-4096	
FC-4096	
FC-1000	
softmax	

Fei-Fei Li & Andrej Karpathy

Lecture 7 -

Transfer Learning with CNNs

image

conv-64

conv-64

maxpool

conv-128

conv-128

maxpool

conv-256

conv-256

conv-512

maxpool

maxpool

FC-4096

FC-4096

FC-1000

softmax



1. Train on Imagenet maxpool conv-512 conv-512 conv-512

2. If small dataset: fix all weights (treat CNN as fixed feature extractor), retrain only the classifier

i.e. swap the Softmax layer at the end

Lecture 7 -

21 Jan 2015

Fei-Fei Li & Andrej Karpathy

Transfer Learning with CNNs

conv-64 conv-64 maxpool conv-128 conv-128 maxpool conv-256 conv-256 maxpool conv-512 conv-512 maxpool conv-512 conv-512 maxpool FC-4096 FC-4096 FC-1000 softmax

image

image 1. Train on Imagenet conv-64 conv-64 maxpool conv-128 conv-128 maxpool conv-256 conv-256 maxpool conv-512 conv-512 maxpool conv-512 conv-512 maxpool FC-4096 FC-4096 FC-1000 softmax

2. If small dataset: fix all weights (treat CNN as fixed feature extractor), retrain only the classifier i.e. swap the Softmax layer at the end

image conv-64 conv-64 maxpool conv-128 conv-128 maxpool conv-256 conv-256 maxpool conv-512 conv-512 maxpool conv-512 conv-512 maxpool FC-4096 FC-4096 FC-1000 softmax

3. If you have medium sized dataset, "finetune" instead: use the old weights as initialization, train the full network or only some of the higher layers

retrain bigger portion of the network, or even all of it.

Fei-Fei Li & Andrej Karpathy

Lecture 7 -

Object Detection: PASCAL VOC mean Average Precision (mAP)



Figure source: Ross Girshick