Structural Priors in Deep Neural Networks

YANI IOANNOU, MAR. 12^{TH} 2018



About Me

• Yani loannou (yu-an-nu)

- Ph.D. Student, University of Cambridge
 Dept. of Engineering, Machine Intelligence Lab
 Prof. Roberto Cipolla, Dr. Antonio Criminisi (MSR)
- \circ Research scientist at Wayve
 - \odot Self-driving car start-up in Cambridge

• Have lived in 4 countries (Canada, UK, Cyprus and Japan)





Research Background

o M.Sc. Computing, Queen's University

- Prof. Michael Greenspan
- o 3D Computer Vision
- Segmentation and recognition in massive unorganized point clouds of urban environments
- "Difference of Normals" multi-scale operator

(Published at 3DIMPVT)



Research Background

Ph.D. Engineering, University of Cambridge (2014 - 2018)
 Prof. Roberto Cipolla, Dr. Antonio Criminisi (Microsoft Research)
 Microsoft PhD Scholarship, 9-month internship at Microsoft Research





Ph.D. – Collaborative Work

- \odot Segmentation of brain tumour tissues with CNNs
 - D. Zikic, Y. Ioannou, M. Brown, A. Criminisi (MICCAI-BRATS 2014) MICCAI-BRATS 2014
 - $\,\circ\,$ One of the first papers using deep learning for volumetric/medical imagery
- \circ Using CNNs for Malaria Diagnosis
 - Intellectual Ventures/Gates Foundation
 - $\,\circ\,$ Designed CNN for the classification of malaria parasites in blood smears
- Measuring Neural Net Robustness with Constraints
 - O. Bastani, Y. Ioannou, L. Lampropoulos, D. Vytiniotis, A. Nori, A. Criminisi NIPS 2016
 - $\,\circ\,$ Found that not all adversarial images can be used to improve network robustness
- \circ Refining Architectures of Deep Convolutional Neural Networks
 - S. Shankar, D. Robertson, Y. Ioannou, A. Criminisi, R. Cipolla
 - CVPR 2016
 - $\,\circ\,$ Proposed a method for adapting neural network architectures to new datasets





Ph.D. – First Author

- oThesis: "Structural Priors in Deep Neural Networks"
 - Training CNNs with Low-Rank Filters for Efficient Image Classification
 Yani Ioannou, Duncan Robertson, Jamie Shotton, Roberto Cipolla, Antonio Criminisi
 ICLR 2016
 - Deep Roots: Improving CNN Efficiency with Hierarchical Filter Groups Yani Ioannou, Duncan Robertson, Roberto Cipolla, Antonio Criminisi CVPR 2017
 - Decision Forests, Convolutional Networks and the Models In-Between
 Y. Ioannou, D. Robertson, D. Zikic, P. Kontschieder, J. Shotton, M. Brown, A. Criminisi Microsoft Research Tech. Report (2015)



- Deep Neural Networks are massive!
- AlexNet¹ (2012)
 - o 61 million parameters
 - 724 million FLOPS
 - Most compute in conv. layers



¹ Krizhevsky, Sutskever, and Hinton, "ImageNet Classification with Deep Convolutional Neural Networks" ² He, Zhang, Ren, and Sun, "Deep Residual Learning for Image Recognition"



- Deep Neural Networks are massive!
- AlexNet¹ (2012)
 - o 61 million parameters
 - 724 million FLOPS
 - 96% of param in F.C. layers!



¹ Krizhevsky, Sutskever, and Hinton, "ImageNet Classification with Deep Convolutional Neural Networks" ² He, Zhang, Ren, and Sun, "Deep Residual Learning for Image Recognition"



- Deep Neural Networks are massive!
- AlexNet¹ (2012)
 - o 61 million parameters
 - 7.24x10⁸ million FLOPS
- ResNet² 200 (2015)
 - o 62.5 million parameters
 - 5.65x10¹² FLOPS
 - 2-3 weeks of training on 8 GPUs



¹ Krizhevsky, Sutskever, and Hinton, "ImageNet Classification with Deep Convolutional Neural Networks" ² He, Zhang, Ren, and Sun, "Deep Residual Learning for Image Recognition"



- Until very recently, state-of-theart DNNs for Imagenet were only getting more computationally complex
- Each generation increased in depth and width
- Is it necessary to increase complexity to improve generalization?



Over-parameterization of DNNs

- There are many proposed methods for improving the test time efficiency of DNNs showing that trained DNNs are over-parameterized
- Compression
- Pruning
- Reduced Representation



Structural Prior

Incorporating our prior knowledge of the problem and its representation into the connective structure of a neural network

Optimization of neural networks needs to learn what weights not to use

- This is usually achieved with regularization
- Can we structure networks closer to the specialized components used for learning images with our prior knowledge of the problem/it's representation?
- Structural Priors ⊂ Network Architecture
 - architecture is a more general term, i.e., number of layers, activation functions, pooling, etc.



Regularization

• Regularization does help training, but is not a substitute for good structural priors

 MacKay (1991): regularization is not enough to make an over-parameterized network generalize as well as a network with a more appropriate parameterization

• We liken regularization to a weak structural prior

 \circ Used where our only prior knowledge is that our network is greatly over-parameterized



Rethinking Regularization

"Understanding deep learning requires rethinking generalization", Zhang et al., 2016
 "Deep neural networks easily fit random labels."

- Identifies types of "regularization":
 - "Explicit regularization" i.e. weight decay, dropout and data augmentation
 - "Implicit regularization" i.e. early stopping, batch normalization
 - o "Network architecture"
- Explicit regularization has little effect on fitting random labels, while implicit regularization and network architecture does
- Highlights the importance of network architecture, and by extension structural priors, for good generalization



Convolutional Neural Networks

Prior Knowledge for Natural Images:

Local correlations are very important

-> Convolutional filters

•We don't need to learn a different filter for each pixel

O -> Shared weights







Convolutional Neural Networks

Structural Prior for Natural Images





Input pixels

Convolutional Neural Networks

Structural Prior for Natural Images



Ph.D. Thesis Outline

My thesis is based on three novel contributions which have explored separate aspects of structural priors in DNN:

I. Spatial Connectivity

II. Inter-Filter Connectivity

III. Conditional Connectivity



Spatial Connectivity



Spatial Connectivity

Prior Knowledge:

- Many of the filters learned in CNNs appear to be representing vertical/horizontal edges/relationships
- Many others appear to be representable by combinations of low-rank filters
- Previous work had shown that full-rank filters could be replaced with low rank *approximations*, e.g. Jaderberg (2014)



Does every filter need to be square in a CNN?





Approximated Low-Rank Filters

Jaderberg, Max, Andrea Vedaldi, and Andrew Zisserman (2014) "Speeding up Convolutional Neural Networks with Low Rank Expansions".





CNN with Low-Dimensional Embedding

Typical sub-architecture found in Network-in-Network, ResNet/Inception





Proposed: Low-Rank Basis

Same total number of filters on each layer as original network, but 50% are 1x3, and 50% are 3x1





Proposed Structural Prior: Low-Rank + Full Basis

25% of total filters are full 3x3





Inception Learning a Filter-Size Basis – learning many small filters (1x1, 3x3), and fewer of the larger (5x5, 7x7)



ImageNet Results

 gmp: vgg-11 w/ global max pooling

○ gmp-lr-2x:

o 60% less computation

o gmp-lr-join-wfull:

16% less computation

o 1% pt. lower error









Structural Prior for CNNs





1. Introduction

Since the 2012 ImageNet competition [16] winning en-

has made it leasible to utilize inception networks in big-data scenarios[17], [13], where huge amount of data needed to

29

Inception v.3

Google's Inception architecture (v.3 and higher) uses our low-rank filters!



Inter-Filter Connectivity



Inter-filter Connectivity

Prior Knowledge:

• CNNs learn sparse, distributed representations

Most filters on adjacent layers have low correlation

Does every filter need to be connected to every other filter on a previous layer in a CNN?



AlexNet Filter Groups



• AlexNet¹ used model parallelization to fit in the GPU memory constraints of the time

o "filter groups" used to split the network into two on all conv layers (except conv3)

¹Krizhevsky, Sutskever, and Hinton, "ImageNet Classification with Deep Convolutional Neural Networks"

應應 UNIVERSITY OF 國國 CAMBRIDGE



AlexNet Filter Groups

Convolutional filters filters in layers with g groups only operate on 1/g of the # input channels



AlexNet Filter Groups



• Filter groups reduce connectivity between filters, allowing easier model parallelization

• Filter groups drastically reduce the number of parameters, and computation

o ... and they don't seem to affect the generalization of AlexNet?!





CNN with Low-Dimensional Embedding

Typical sub-architecture found in Network-in-Network, ResNet/Inception





Root-2 Module

Structural Prior for CNNs with Sparse Inter-Filter Relationships





Root-4 Module

Structural Prior for CNNs with Sparse Inter-Filter Relationships



Network in Network Filter Groups



Replace non-spatial convolutional layers with root modules



Filter Group Topologies

• But how many groups to use? Should this change with depth?

• We explored 3 basic topologies on CIFAR10:





CIFAR-10 Results



NiN: mean and standard deviation (error bars) are shown over 5 different random initializations.



CIFAR-10 Results



NiN: mean and standard deviation (error bars) are shown over 5 different random initializations.



Covariance



•Block-diagonal sparsity effected by a root-module is visible in the inter-layer correlation



ILSVRC12 Results – ResNet 50

o root-16

- 27% fewer parameters
- 37% less computation
- CPU 23% faster
- GPU 13% faster
 - (not optimized!)
- o 0.2% pt. lower error

o root-64

- 40% fewer parameters
- 45% less computation
- CPU 31% faster
- GPU 12% faster
- 0.1% pt. higher error



ILSVRC12 Results – ResNet 200

- o root-64
 - 27% fewer parameters
 - 48% less computation
 - o 0.2% pt. lower error
 - \circ 0.14% lower error





Root Module

Structural Prior for CNNs with Sparse Inter-Filter Relationships



[cs.NE] 489v1

Deep Roots: Improving CNN Efficiency with Hierarchical Filter Groups Yani Ioannou¹, Duncan Robertson², Roberto Cipolla¹, and Antonio Criminisi² ¹University of Cambridge, ²Microsoft Research

¹University of Cambridge, ²Microsoft Research

Abstract. We propose a new method for training computationally efficient and compact convolutional neural networks (CNNs) using a novel sparse connection structure that resembles a tree root. Our sparse connection structure facilitates a significant reduction in computational cost and number of parameters of state-of-the-art deep CNNs without compromising accuracy. We validate our approach by using it to train more efficient variants of state-of-the-art CNN architectures, evaluated on the CIFAR10 and ILSVRC datasets. Our results show similar or higher accuracy than the baseline architectures with much less compute, as measured

Deep Roots





Xception

Google's Xception architecture uses a form of root modules (#channels = #filter groups) - "Depthwise Separable Convolution"





ResNeXt

Facebook's ResNet architecture uses root modules, denoted "Aggregated Residual Transforms"



Conclusion

- Structural priors are important for both generalization and efficiency
- They are not simply replaced by strong regularization
- Simplify the optimization of deep neural networks by constraining the search space/dimensionality
- There is still a lot we don't understand about the optimization of deep neural networks!



I. Automatically Discovering Structural Priors

II. Learning with "Natural" Datasets

III. Jointly Exploiting Random Exploration and Imitation



I. Automatically Discovering Structural Priors

 Can we find methods of automatically discovering good structural priors from data?

- Pruning does not improving generalization
- Greedily growing networks leads to poor generalization
- Results by Han et. al¹ show some promise: pruning/growing cycle
- Infer connectivity by analyzing inter-channel correlations in training data?

Han, Song, Jeff Pool, John Tran, and William J. Dally (2015). "Learning both weights and connections for efficient neural networks



II. Learning with "Natural" Datasets

• Both ML and Computer Vision are dataset driven fields

• ImageNet, CIFAR and MNIST are class-balanced

• Current solutions involve either throwing away data or fiddling with loss weighting



III. Exploiting both random exploration and imitation

- RL is appealing agents learn entirely from experience in an environment
- For many problems this isn't data efficient enough or feasible:
 - e.g. learning to drive a car randomly exploring in a real environment is dangerous and time consuming
 - \odot But we can easily collect data from a human driver for real-world driving
- Use supervised learning to bootstrap RL



Questions

http://yani.io/annou

yai20@cam.ac.uk

