

argmax_x

Learning the Number of Neurons in Deep Networks

Jose M. Alvarez¹ Mathieu Salzmann²

¹Data61 @ CSIRO, Canberra, ACT 2601, Australia

²CVLab, EPFL, CH-1015 Lausanne, Switzerland

NIPS, 2016

Presenter: Arshdeep Sekhon

Motivation

- 1 Designing a deep NN architecture
- 2 Configure:
 - number of layers
 - number of units
- 3 mostly hand-designed
- 4 have redundant parameters

Automatic Model Selection: Constructive Approaches

- ① Incrementally add layers/parameters
- ② But Shallow networks are less expressive
- ③ bad initialization when incrementally adding layers

Automatic Model Selection: Destructive Approaches

- 1 Very Deep networks more expressive
- 2 Start from very deep networks, eliminate redundant parameters
- 3 check influence of every parameter
- 4 for example, check network Hessian wrt every parameter in an over complete network
- 5 not scalable to large networks

Automatic Model Selection

- ① Automatically get number of neurons for each layer
- ② cancel effects of individual neurons
- ③ jointly as we learn
- ④ no preprocessing

Model Selection

- 1 Start with an overcomplete network
- 2 Find neurons for each layer
- 3 A general deep network:
 L layers in network architecture
 N_l neurons in each layer
- 4 weights $\Theta = [\theta_l, b_l]$ for layer l $\theta_l = [\theta_l^n]$ $1 \leq l \leq L$ and $1 \leq n \leq N_l$
- 5 The optimization problem:

$$\min_{\Theta} \frac{1}{N} \sum_{i=1}^N \ell(y_i, f(x_i, \Theta)) + r(\Theta) \quad (1)$$

Model Selection: Regularizer

- 1 The optimization problem:

$$\min_{\Theta} \frac{1}{N} \sum_{i=1}^N \ell(y_i, f(x_i, \Theta)) + r(\Theta) \quad (2)$$

- 2 Goal: Cancel entire neurons

Model Selection: Regularizer

- 1 The optimization problem:

$$\min_{\Theta} \frac{1}{N} \sum_{i=1}^N \ell(y_i, f(x_i, \Theta)) + r(\Theta) \quad (2)$$

- 2 Goal: Cancel entire neurons
- 3 traditional regularizers: ℓ_1 or ℓ_2
- 4 cannot cancel entire neurons because they control weights individually.

Model Selection: Regularizer

- 1 The optimization problem:

$$\min_{\Theta} \frac{1}{N} \sum_{i=1}^N \ell(y_i, f(x_i, \Theta)) + r(\Theta) \quad (2)$$

- 2 Goal: Cancel entire neurons
- 3 traditional regularizers: ℓ_1 or ℓ_2
- 4 cannot cancel entire neurons because they control weights individually.
- 5 Neurons are *groups* of parameters
- 6 weights $\Theta = [\theta_l, b_l]$ for layer l $\theta_l = [\theta_l^n]$ $1 \leq l \leq L$ and $1 \leq n \leq N_l$

Model Selection: Regularizer

- 1 The optimization problem:

$$\min_{\Theta} \frac{1}{N} \sum_{i=1}^N \ell(y_i, f(x_i, \Theta)) + r(\Theta) \quad (2)$$

- 2 Goal: Cancel entire neurons
- 3 traditional regularizers: ℓ_1 or ℓ_2
- 4 cannot cancel entire neurons because they control weights individually.
- 5 Neurons are *groups* of parameters
- 6 weights $\Theta = [\theta_l, b_l]$ for layer l $\theta_l = [\theta_l^n]$ $1 \leq l \leq L$ and $1 \leq n \leq N_l$
- 7 Use new regularizer: *group sparsity*

Model Selection: Group Sparsity

- 1 Parameters associated with a neuron are grouped together
- 2 Penalty on groups of weights instead of individual weights
- 3 parameters of each neuron in layer l are grouped in a vector of size P_l
- 4 New regularizer:

$$r(\Theta) = \sum_{l=1}^L \lambda_l \sqrt{P_l} \sum_{n=1}^{N_l} \|\theta_n^l\|_2 \quad (3)$$

- 5 θ_n^l are the parameters for neuron n in layer l
- 6 ℓ_2 norm followed by ℓ_1 norm
- 7 λ_l sets the influence of the penalty.

Model Selection: Group Sparsity

- 1 But does not lead to sparsity within a group

$$r(\Theta) = \sum_{l=1}^L (1 - \alpha) \lambda_l \sqrt{P_l} \sum_{n=1}^{N_l} \|\theta'_n\|_2 + \alpha \lambda_l \|\theta_l\|_1 \quad (4)$$

- 2 more general penalty that leads to sparsity both at and within group level.

Training: Proximal Gradient Descent

$$\text{minimize } f(x) = g(x) + h(x) \quad (5)$$

proximal gradient algorithm:

$$x^{k+1} = \mathbf{prox}_{t_k h} \left(x^{k-1} - t_k (\nabla g(x^{k-1})) \right) \quad (6)$$

proximal operator:

$$\mathbf{prox}_h(x) = \arg \min_u h(u) + \frac{1}{2} \|x - u\|_2^2 \quad (7)$$

$$x^{k+1} = \arg \min_u \left(h(u) + \frac{1}{2t} \|u - x^{k-1} + t_k (\nabla g(x^{k-1}))\|_2^2 \right) \quad (8)$$

Training: Proximal Gradient Descent

The objective:

$$\min_{\Theta} \frac{1}{N} \sum_{i=1}^N \ell(y_i, f(x_i, \Theta)) + r(\Theta) \quad (9)$$

$$r(\Theta) = \sum_{l=1}^L (1 - \alpha) \lambda_l \sqrt{P_l} \sum_{n=1}^{N_l} \|\theta_n^l\|_2 + \alpha \lambda_\ell \|\theta_\ell\|_1 \quad (10)$$

loss function is $g(x)$ and regularizer $h(x)$ in proximal gradient algorithm

Training: Proximal Gradient Descent

- 1 Update: Take gradient of loss and apply proximal operator of the regularizer

2

$$\tilde{\theta}_l^n = \arg \min_{\tilde{\theta}_l^n} \frac{1}{2t} \|\tilde{\theta}_l^n - \hat{\theta}_l^n\|_2^2 + r(\Theta) \quad (11)$$

where $\hat{\theta}_l^n$ is update by gradient of loss function

- 3 This has a closed form solution:

$$\tilde{\theta}_l^n = \left(1 - \frac{t(1-\alpha)\lambda_l\sqrt{P_l}}{\|S(\hat{\theta}_l^n, t\alpha\lambda_l)\|_2} \right)_+ S(\hat{\theta}_l^n, t\alpha\lambda_l)$$

$$(S(\mathbf{z}, \tau))_j = \text{sign}(z_j)(\|\mathbf{z}\| - \tau)_+$$

Experiments and Model Architectures

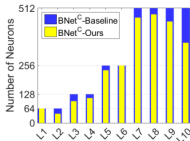
- ① Dataset: ImageNet , Places2-401
- ② Models:
 - ① VGG-B Net: 10 convolutional layers followed by three fully-connected layers
 - ② DecomposeMe₈ (Dec₈): 16 Conv layers with 1D kernels

Experiments and Model Architectures

Table 1: Top-1 accuracy results for several state-of-the-art architectures and our method on ImageNet.

Model	Top-1 acc. (%)	Model	Top-1 acc. (%)
BNet	62.5	Ours-Bnet ^C _{GS}	62.7
BNet ^C	61.1	Ours-Dec ₈ - _{GS}	64.8
ResNet50 ^a [He et al., 2015]	67.3	Ours-Dec ₈ -640 _{SGL}	67.5
Dec ₈	64.8	Ours-Dec ₈ -640 _{GS}	68.6
Dec ₈ -640	66.9	Ours-Dec ₈ -768 _{GS}	68.0
Dec ₈ -768	68.1		

^a Trained over 55 epochs using a batch size of 128 on two TitanX with code publicly available.



BNet ^C on ImageNet (in %)	
	GS
neurons	12.70
group param	13.59
total param	13.59
total induced	27.38
accuracy gap	1.6

Figure 1: Parameter reduction on ImageNet using BNet^C. (Left) Comparison of the number of neurons per layer of the original network with that obtained using our approach. (Right) Percentage of zeroed-out neurons and parameters, and accuracy gap between our network and the original one. Note that we outperform the original network while requiring much fewer parameters.