

# COMP9444

## Neural Networks and Deep Learning

### Term 2, 2024



## Week 5 Tutorial: Convolutional Neural Networks (Sample Solution)

### 1. Convolutional Network Architecture

One of the early papers on Deep Q-Learning for Atari games (Mnih et al., 2013) contains this description of its Convolutional Neural Network:

“The input to the neural network consists of an  $84 \times 84 \times 4$  image. The first hidden layer convolves 16  $8 \times 8$  filters with stride 4 with the input image and applies a rectifier nonlinearity. The second hidden layer convolves 32  $4 \times 4$  filters with stride 2, again followed by a rectifier nonlinearity. The final hidden layer is fully-connected and consists of 256 rectifier units. The output layer is full-connected linear layer with a single output for each valid action. The number of valid actions varied between 4 and 18 on the games we considered.”

You should assume the input images are gray-scale, there is no padding, and there are 18 valid actions (outputs).

For each layer in this network, compute the number of:

- (a) weight per filter in this layer (including bias)
- (b) width and height of layer (convolutional layers)
- (c) neurons in this layer
- (d) connections into the neurons in this layer
- (e) independent parameters in this layer

---

First Convolutional Layer:

$$J = K = 84, L = 4, M = N = 8, P = 0, s = 4$$

$$\text{Weights per filter: } 1 + M \times N \times L = 1 + 8 \times 8 \times 4 = 257$$

$$\text{Width and Height of layer: } 1 + (J - M)/s = 1 + (84 - 8)/4 = 20$$

$$\text{Neurons in layer: } 20 \times 20 \times 16 = 6400$$

$$\text{Connections: } 20 \times 20 \times 16 \times 257 = 1644800$$

$$\text{Independent parameters: } 16 \times 257 = 4112$$

Second Convolutional Layer:

$$J = K = 20, L = 16, M = N = 4, P = 0, s = 2$$

$$\text{Weights per filter: } 1 + M \times N \times L = 1 + 4 \times 4 \times 16 = 257$$

$$\text{Width and Height of layer: } 1 + (J - M)/s = 1 + (20 - 4)/2 = 9$$

Neurons in layer:  $9 \times 9 \times 32 = 2592$   
Connections:  $9 \times 9 \times 32 \times 257 = 666144$   
Independent parameters:  $32 \times 257 = 8224$

Fully Connected Layer:

Weights per neuron:  $1 + 2592 = 2593$   
Neurons in layer: 256  
Connections:  $256 \times 2593 = 663808$   
Independent parameters: 663808

Output Layer:

Weights per neuron:  $1 + 256 = 257$   
Neurons in layer: 18  
Connections:  $18 \times 257 = 4626$   
Independent parameters: 4626

## 2. Weight Initialization

Briefly describe the problem of vanishing or exploding gradients, and how Weight Initialization can be used to prevent it.

---

The differentials in a deep neural network tend to grow according to this equation:

$$Var[\partial/\partial x] \approx (\prod_{1 \leq i \leq D} G_1 n_{i+1} Var[w^{(i)}]) Var[\partial/\partial z]$$

where  $w^{(i)}$  are the weights at layer  $i$ ,  $n(i + 1)$  is the number of weights fanning out from each node in layer  $i$ , and  $G_1$  estimates the average value of the derivative of the transfer function. If the weights are initialised so that the factor in parentheses corresponding to each layer is approximately 1, then the differentials will remain in a healthy range. Otherwise, they may either grow or vanish exponentially.

## 3. Batch Normalization

Briefly describe the Batch Normalization algorithm.

---

The mean and variance of the activations  $x_k^{(i)}$  at layer  $i$  over a batch of training items are estimated or pre-computed, and normalized activations are calculated for each node:

$$\hat{x}_k^{(i)} = \frac{(x_k^{(i)} - Mean[x_k^{(i)}])}{\sqrt{Var[x_k^{(i)}]}}$$

These activations are then shifted and rescaled by  $y_k^{(i)} = \beta_k^{(i)} + \gamma_k^{(i)} \hat{x}_k^{(i)}$ , where  $\beta_k^{(i)}, \gamma_k^{(i)}$  are additional parameters to be learned by backpropagation.

#### 4. Residual and Dense Networks

Explain the difference between a Residual Network and a Dense Network.

---

A Residual Network includes “skip” connections which bypass each pair of consecutive layers. These intermediate layers therefore compute a residual component, which is added to the output from previous layers and corrects their errors, or provides additional details which they were not powerful enough to compute.

A Dense Network is built from densely connected blocks, separated by convolution and pooling layers. Within a dense block, each layer is connected by shortcut connections to all preceding layers.

#### 5. Other Questions

Any questions or discussion about image processing, societal impact, or any other aspect of the course.