

Machine Learning Assignment 2

Neural Networks

Submission deadline: November 21, 2024

Please submit your solution in PDF format (preferably, but not necessarily, L^AT_EX—this .tex file can be found on iCorsi). Handwriting and scanned documents are not allowed. In case you need further help, please write on iCorsi or contact me at vincent.herrmann@usi.ch.

1 Problem 1. Calculating Gradients (20 points)

A two-layer neural network is defined by the following equation:

$$y_k := \mathbf{w}^{(2)\top} f(W^{(1)} \mathbf{x}_k + \mathbf{b}^{(1)}) + b^{(2)}$$

The values of the input vectors $\mathbf{x}_k \in \mathbb{R}^4$ (summarized in the matrix $X \in \mathbb{R}^{3 \times 4}$), targets t_k (summarized in the vector $\mathbf{t} \in \mathbb{R}^3$), weights $W^{(1)} \in \mathbb{R}^{2 \times 4}$ and $\mathbf{w}^{(2)} \in \mathbb{R}^2$ as well as biases $\mathbf{b}^{(1)} \in \mathbb{R}^2$ and $b^{(2)} \in \mathbb{R}$ are given below.

$$X = \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \mathbf{x}_3 \end{bmatrix} = [[0.2, -0.4, 0.1, 0.5], [0.8, 1.0, -0.1, 0.3], [-0.5, 0.9, 0.2, -0.8]]$$
$$\mathbf{t} = \begin{bmatrix} t_1 \\ t_2 \\ t_3 \end{bmatrix} = [-0.5, 0.1, 0.8]$$

$$W^{(1)} = [[-0.6, 0.7, 0.9, -0.3], [0.4, 0.5, -1.0, 0.1]]$$

$$\mathbf{b}^{(1)} = [-0.4, -0.5]$$

$$\mathbf{w}^{(2)} = [0.2, 0.6]$$

$$b^{(2)} = -0.7$$

The function f is the ReLU nonlinearity, applied element-wise. The loss of the model is

$$L := \frac{1}{2} \sum_{k=1}^3 (y_k - t_k)^2.$$

- Numerically calculate the loss and the gradients of L with respect to $W^{(1)}$, $\mathbf{w}^{(2)}$, $\mathbf{b}^{(1)}$ and $b^{(2)}$, and explain your process. You may use a calculator or math software (numpy, matlab, ...), but no auto-differentiation libraries like PyTorch or Jax. Tip: If you use matrix-matrix multiplications involving X , you don't need to do multiple explicit forward and backward passes.

- It is more computationally efficient to calculate gradients in a neural network when we start the chain of derivatives on the side of the loss and then going backward, as opposed to starting on the input side and going forward. Explain why this is the case using the given example.

Problem 2. Gradient Descent (15 points)

If we choose a learning rate that is too high, gradient descent can diverge (even in the case of convex functions). Assume we want to use gradient descent to find the value of $\mathbf{u} = \begin{bmatrix} x \\ y \end{bmatrix}$

that minimizes $f(\mathbf{u}) = \mathbf{u}^T \begin{bmatrix} 5 & 2 \\ 2 & 2 \end{bmatrix} \mathbf{u}$, or equivalently $f(x, y) = 5x^2 + 2y^2 + 4xy$. The update via gradient descent can be written as $\mathbf{u}_{n+1} = \mathbf{u}_n - \eta \nabla f(\mathbf{u}_n)$, with η being the learning rate. What is the range of values that η can take so that gradient descent eventually converges to the minimum, no matter the starting point? Explain your reasoning—the correct approach is more important than the numerical result.

Problem 3. Convolutional Networks (15 points)

We start with an input image of size 32×32 and three channels. Provide a sequence of alternating convolutional and pooling layers that transform this image into 64 feature maps of size 4×4 .

For each convolutional layer, specify the kernel size (e.g., 3×3 , 5×5 or 7×7), the number of input channels, and the number of filters. All convolutional layers should have a stride of one in either direction. No padding is used. All pooling layers should be max-pooling layers with a kernel size of 2×2 and a stride of two in each direction. Give the size and the number of the intermediate feature maps (i.e., the outputs of each layer). Calculate the total number of learnable parameters in your network (the number of weights in biases in all layers).

Your answer should look roughly like this:

Architecture 1:

- input feature map: 32×32 , 3 feature maps
- conv layer: kernel size $? \times ?$, ? input feature maps, ? output feature maps
- output conv layer: $? \times ?$, ? feature maps
- pooling layer: 2×2
- output pooling layer: $? \times ?$, ? channels/feature maps
- conv layer: ...
- ...
- output pooling/conv layer: 4×4 , 64 channels/feature maps

Total number of learnable parameters (weights and biases): ?

PROBLEM 1

$$Y_K = \begin{vmatrix} 0.2 & 0.6 \end{vmatrix} F \left(\begin{vmatrix} -0.6 & 0.7 & 0.9 & -0.3 \\ 0.4 & 0.5 & -1.0 & 0.1 \end{vmatrix} X_K + \begin{vmatrix} -0.4 \\ -0.5 \end{vmatrix} \right)$$

$$+ |-0.7| \quad F(x) = \text{ReLU}(z) = \max(0, x)$$

$$X_1 = \begin{vmatrix} 0.2 \\ -0.4 \\ 0.1 \\ 0.5 \end{vmatrix} \quad X_2 = \begin{vmatrix} 0.8 \\ 1.0 \\ -0.1 \\ 0.3 \end{vmatrix} \quad X_3 = \begin{vmatrix} -0.5 \\ 0.9 \\ 0.2 \\ -0.8 \end{vmatrix} \quad b = \begin{vmatrix} -0.5 \\ 0.1 \\ 0.8 \end{vmatrix}$$

$$Z_K = F(W^{(k)} X_K + b^{(k)})$$

QUESTION 1

$$K=1$$

$$\begin{vmatrix} -0.6 & 0.7 & 0.9 & -0.3 \\ 0.4 & 0.5 & -1.0 & 0.1 \end{vmatrix} \cdot \begin{vmatrix} 0.2 \\ -0.4 \\ 0.1 \\ 0.5 \end{vmatrix} + \begin{vmatrix} -0.4 \\ -0.5 \end{vmatrix} =$$

$$= \begin{vmatrix} -0.6 \cdot 0.2 - 0.7 \cdot 0.4 + 0.9 \cdot 0.1 - 0.3 \cdot 0.5 & -0.4 \\ 0.4 \cdot 0.2 - 0.5 \cdot 0.4 - 1.0 \cdot 0.1 + 0.1 \cdot 0.5 & -0.5 \end{vmatrix} =$$

$$= \begin{vmatrix} -0.86 \\ -0.67 \end{vmatrix} \Rightarrow Y_1 = \begin{vmatrix} 0.2 & 0.6 \end{vmatrix} \cdot F \left(\begin{vmatrix} -0.86 \\ -0.67 \end{vmatrix} \right) + |-0.7|$$

$$= \begin{vmatrix} 0.2 & 0.6 \end{vmatrix} \cdot \begin{vmatrix} 0 \\ 0 \end{vmatrix} + |-0.7| = |-0.7|$$

K=2

$$\begin{vmatrix} -0.6 & 0.7 & 0.9 & -0.3 \\ 0.4 & 0.5 & -1.0 & 0.1 \end{vmatrix} \cdot \begin{vmatrix} 0.8 \\ 1.0 \\ -0.1 \\ 0.3 \end{vmatrix} + \begin{vmatrix} -0.4 \\ -0.5 \end{vmatrix} =$$

$$= \begin{vmatrix} -0.6 \cdot 0.8 + 0.7 \cdot 1.0 - 0.9 \cdot 0.1 - 0.3 \cdot 0.3 & -0.4 \\ 0.4 \cdot 0.8 + 0.5 \cdot 1.0 + 1.0 \cdot 0.1 + 0.1 \cdot 0.3 & -0.5 \end{vmatrix} =$$

$$= \begin{vmatrix} -0.36 \\ 0.45 \end{vmatrix} \Rightarrow \gamma_2 = |0.2 \quad 0.6| \cdot F\left(\begin{vmatrix} -0.36 \\ 0.45 \end{vmatrix}\right) + |-0.7|$$

$$= |0.2 \quad 0.6| \cdot \begin{vmatrix} 0 \\ 0.45 \end{vmatrix} + |-0.7| = [-0.43]$$

K=3

$$\begin{vmatrix} -0.6 & 0.7 & 0.9 & -0.3 \\ 0.4 & 0.5 & -1.0 & 0.1 \end{vmatrix} \cdot \begin{vmatrix} -0.5 \\ 0.9 \\ 0.2 \\ -0.8 \end{vmatrix} + \begin{vmatrix} -0.4 \\ -0.5 \end{vmatrix} =$$

$$= \begin{vmatrix} +0.6 \cdot 0.5 + 0.7 \cdot 0.9 + 0.9 \cdot 0.2 + 0.3 \cdot -0.8 & -0.4 \\ -0.4 \cdot 0.5 + 0.5 \cdot 0.9 - 1.0 \cdot 0.2 - 0.1 \cdot -0.8 & -0.5 \end{vmatrix} =$$

$$= \begin{vmatrix} 0.95 \\ -0.53 \end{vmatrix} \Rightarrow \gamma_3 = |0.2 \quad 0.6| \cdot F\left(\begin{vmatrix} 0.95 \\ 0 \end{vmatrix}\right) + |-0.7|$$

$$= |0.2 \quad 0.6| \cdot \begin{vmatrix} 0.95 \\ 0 \end{vmatrix} + |-0.7| = [-0.51]$$

$$\Rightarrow \gamma = \begin{vmatrix} -0.7 \\ -0.43 \\ -0.51 \end{vmatrix} \quad \zeta = \begin{vmatrix} -0.5 \\ 0.1 \\ 0.8 \end{vmatrix}$$

$$L = \frac{1}{2} \left[(\gamma_1 - \zeta_1)^2 + (\gamma_2 - \zeta_2)^2 + (\gamma_3 - \zeta_3)^2 \right] = 1.0185$$

$$\frac{\partial L}{\partial w^{(2)}}$$

$$\frac{\partial L}{\partial w^{(2)}} = \frac{\partial}{\partial w^{(2)}} \sum_{k=1}^3 \frac{1}{2} (\gamma_k - \zeta_k)^2 = \sum_{k=1}^3 \frac{1}{2} \frac{\partial (\gamma_k - \zeta_k)^2}{\partial \gamma_k} \cdot \frac{\partial \gamma_k}{\partial z_k} \cdot \frac{\partial z_k^{(2)}}{\partial w^{(2)}}$$

$$= \sum_{k=1}^3 (\gamma_k - \zeta_k) \cdot 1 \cdot \frac{\partial}{\partial z_k^{(2)}} (w^{(2)\top} \cdot z_k^{(2)} + b_2^{(2)}) \cdot \frac{\partial}{\partial w^{(2)}} (w^{(2)\top} \cdot x_k + b^{(2)}) = \sum_{k=1}^3 (\gamma_k - \zeta_k) \cdot w^{(2)\top} \cdot F'(z_k) \cdot x_k^\top$$

$$F(z) = \max(0, z) \Rightarrow F'(z) = \begin{cases} 1 & z \geq 0 \\ 0 & z < 0 \end{cases}$$

$$\Rightarrow \frac{\partial L}{\partial w^{(2)}} = (\gamma_1 - \zeta_1) \cdot \begin{vmatrix} 0.2 \\ 0.6 \end{vmatrix} \cdot F'(z_1) \cdot x_1^\top +$$

$$(\gamma_2 - \zeta_2) \cdot \begin{vmatrix} 0.2 \\ 0.6 \end{vmatrix} \cdot F'(z_2) \cdot x_2^\top +$$

$$(\gamma_3 - \zeta_3) \cdot \begin{vmatrix} 0.2 \\ 0.6 \end{vmatrix} \cdot F'(z_3) \cdot x_3^\top =$$

$$= (-0.2) \cdot \begin{vmatrix} 0.2 \\ 0.6 \end{vmatrix} \cdot F' \left(\begin{vmatrix} -0.86 \\ -0.67 \end{vmatrix} \right) \cdot [0.2 \ 0.4 \ 0.1 \ 0.5] +$$

$$(-0.53) \cdot \begin{vmatrix} 0.2 \\ 0.6 \end{vmatrix} \cdot F' \left(\begin{vmatrix} -0.36 \\ 0.45 \end{vmatrix} \right) \cdot [0.8 \ 1.0 \ -0.1 \ 0.31] +$$

$$(-1.31) \cdot \begin{vmatrix} 0.2 \\ 0.6 \end{vmatrix} \cdot F' \left(\begin{vmatrix} 0.95 \\ -0.53 \end{vmatrix} \right) \cdot \begin{vmatrix} -0.5 & 0.9 & 0.2 & -0.8 \end{vmatrix}$$

$$(1) = 0 - 0.53 \cdot \begin{vmatrix} 0 \\ 0.6 \end{vmatrix} \cdot \begin{vmatrix} 0.8 & 1.0 & -0.1 & 0.31 \end{vmatrix}$$

$$-1.31 \cdot \begin{vmatrix} 0.2 \\ 0 \end{vmatrix} \cdot \begin{vmatrix} -0.5 & 0.9 & 0.2 & -0.8 \end{vmatrix} =$$

$$= \begin{vmatrix} 0 \\ -0.318 \end{vmatrix} \cdot \begin{vmatrix} 0.8 & 1.0 & -0.1 & 0.31 \end{vmatrix} +$$

$$\begin{vmatrix} -0.262 \\ 0 \end{vmatrix} \cdot \begin{vmatrix} -0.5 & 0.9 & 0.2 & -0.8 \end{vmatrix} =$$

$$= \begin{vmatrix} 0 & 0 & 0 & 0 \\ -0.2544 & -0.318 & 0.0318 & -0.0954 \end{vmatrix} +$$

$$+ \begin{vmatrix} 0.431 & -0.2358 & -0.0524 & 0.2096 \\ 0 & 0 & 0 & 0 \end{vmatrix}$$

$$= \begin{vmatrix} 0.431 & -0.2358 & -0.0524 & 0.2096 \\ -0.2544 & -0.318 & 0.0318 & -0.0954 \end{vmatrix}$$

$\frac{\partial L}{\partial w^{(2)}}$

$$\frac{\partial L}{\partial w^{(2)}} = \frac{1}{N} \sum_{k=1}^N \frac{1}{2} (\gamma_k - \hat{\gamma}_k)^2 = \sum_{k=1}^N \frac{(\gamma_k - \hat{\gamma}_k)^2}{N} \cdot \frac{\partial \gamma_k}{\partial w^{(2)}} =$$

$$= \sum_{k=1}^3 (\gamma_k - \bar{\gamma}_k) \cdot \frac{\partial}{\partial w^{(2)}} (w^{(2)\top} f(w^{(2)} x_k + b^{(2)})) + b^{(2)} =$$

$$= \sum_{k=1}^3 (\gamma_k - \bar{\gamma}_k) \cdot z_k \Rightarrow \frac{\partial L}{\partial w^{(2)}} = (\gamma_1 - \bar{\gamma}_1) z_1 + (\gamma_2 - \bar{\gamma}_2) z_2$$

$$+ (\gamma_3 - \bar{\gamma}_3) z_3 = (-0.2) \cdot \begin{vmatrix} 0 \\ 0 \end{vmatrix} + 0.53 \cdot \begin{vmatrix} 0 \\ 0.45 \end{vmatrix} + 1.31 \cdot \begin{vmatrix} 0.95 \\ 0 \end{vmatrix}$$

$$= \begin{vmatrix} -1.2445 \\ -0.2385 \end{vmatrix}$$

$$\frac{\partial L}{\partial b^{(2)}}$$

$$\frac{\partial L}{\partial b^{(2)}} = \frac{\partial}{\partial b^{(2)}} \sum_{k=1}^3 \frac{1}{2} (\gamma_k - \bar{\gamma}_k)^2 = \sum_{k=1}^3 \frac{1}{2} \cdot \frac{\partial (\gamma_k - \bar{\gamma}_k)^2}{\partial \gamma_k} \cdot \frac{\partial \gamma_k}{\partial z_k^{(2)}} \cdot \frac{\partial z_k^{(2)}}{\partial b^{(2)}} =$$

$$= \sum_{k=1}^3 (\gamma_k - \bar{\gamma}_k) \cdot w^{(2)\top} \cdot f'(z_k)$$

$$\Rightarrow \frac{\partial L}{\partial b^{(2)}} = (\gamma_1 - \bar{\gamma}_1) \cdot \begin{vmatrix} 0.2 \\ 0.6 \end{vmatrix} \cdot f'(z_1) + (\gamma_2 - \bar{\gamma}_2) \cdot \begin{vmatrix} 0.2 \\ 0.6 \end{vmatrix} \cdot f'(z_2) +$$

$$+ (\gamma_3 - \bar{\gamma}_3) \cdot \begin{vmatrix} 0.2 \\ 0.6 \end{vmatrix} \cdot f'(z_3) = (-0.2) \cdot \begin{vmatrix} 0.2 \\ 0.6 \end{vmatrix} \cdot P' \left(\begin{vmatrix} -0.86 \\ -0.67 \end{vmatrix} \right) +$$

$$+ (-0.53) \cdot \begin{vmatrix} 0.2 \\ 0.6 \end{vmatrix} \cdot P' \left(\begin{vmatrix} -0.36 \\ 0.45 \end{vmatrix} \right) + 1.31 \cdot \begin{vmatrix} 0.2 \\ 0.6 \end{vmatrix} \cdot P' \left(\begin{vmatrix} 0.95 \\ -0.53 \end{vmatrix} \right) =$$

$$= 0 - 0.53 \cdot \begin{vmatrix} 0 \\ 0.6 \end{vmatrix} - 1.31 \cdot \begin{vmatrix} 0.2 \\ 0 \end{vmatrix} = \begin{vmatrix} 0 \\ -0.318 \end{vmatrix} + \begin{vmatrix} -0.262 \\ 0 \end{vmatrix} =$$

$$= \begin{vmatrix} -0.262 \\ -0.318 \end{vmatrix}$$

$$\frac{\partial L}{\partial b^{(2)}}$$

$$\frac{\partial L}{\partial b^{(2)}} = \frac{1}{\partial b^{(2)}} \sum_{k=1}^3 \frac{1}{2} (\gamma_k - t_k)^2 = \sum_{k=1}^3 \frac{1}{2} \frac{(\gamma_k - t_k)^2}{\partial \gamma_k} \cdot \frac{\partial \gamma_k}{\partial b_k} =$$

$$= \sum_{k=1}^3 (\gamma_k - t_k) \cdot \frac{1}{\partial b^{(2)}} \left(\mathbf{w}^{(2)T} F(\mathbf{w}^{(2)} \mathbf{x}_k + b^{(2)}) + b^{(2)} \right) =$$

$$= \sum_{k=1}^3 (\gamma_k - t_k) \cdot 1 \Rightarrow \frac{\partial L}{\partial b^{(2)}} = (-0.7 + 0.5) + (-0.43 - 0.1) + (-0.51 - 0.8) = -2.04$$

QUESTION 2

THE STATEMENT REFERS TO THE FACT THAT BACKPROPAGATION IS MORE EFFICIENT THAN A FORWARD ALGORITHM. IN GENERAL, BACKPROPAGATION HAS SEVERAL ADVANTAGES:

- AVOIDS REPETITIVE COMPUTATIONS, STARTING FROM L AND PROPAGATING THE GRADIENTS ACROSS THE LAYERS
- WE MULTIPLY THE PARTIAL DERIVATIVES FROM LEFT-TO-RIGHT INSTEAD OF RIGHT-TO-LEFT BECAUSE, THEN, WE'LL HAVE VECTOR-MATRIX MULTIPLICATIONS INSTEAD OF MATRIX-MATRIX MULTIPLICATIONS. THIS IS ONLY THE CASE IF WE HAVE MULTIPLE INPUTS AND A SINGLE LOSS, WHICH IS USUALLY THE CASE IN NEURAL NETWORKS.

- AS A CONSEQUENCE, FOR EXAMPLE, ONCE WE GOT A $\beta'(z_k) = 0$
WE CAN IGNORE THE WHOLE NEURON JUST WITH THE FIRST COMPUTATION

PROBLEM 2

$f(w)$ IS A FUNCTION IN QUADRATIC FORM \Rightarrow MATRIX

$\begin{vmatrix} 5 & 2 \\ 2 & 2 \end{vmatrix}$ IS THE HESSIAN MATRIX. SINCE IT IS ALSO

SIMMETRIC, THE EIGENVALUES CAN DETERMINE THE DIRECTION'S

SHAPE. λ 'S CAN BE COMPUTED BY CALCULATING

$$\det\left(\begin{vmatrix} 5 & 2 \\ 2 & 2 \end{vmatrix} - \lambda I\right) = 0$$

$$\det\left(\begin{vmatrix} 5-\lambda & 2 \\ 2 & 2-\lambda \end{vmatrix}\right) = 10 - 7\lambda + \lambda^2 - 4 = 0$$

$$\lambda^2 - 7\lambda + 6 = 0 \quad \begin{array}{c|cc|c} & 1 & -7 & 6 \\ \hline 1 & & 1 & -6 \\ \hline & 1 & -6 & 0 \end{array}$$

$$(\lambda-6)(\lambda-1)=0 \Rightarrow \lambda_1=6; \lambda_2=1$$

IN ORDER FOR GRADIENT DESCENT TO CONVERGE, THE LEARNING

RATE γ SHOULD SATISFY THE CONDITION $0 < \gamma \leq \frac{1}{\lambda_{\max}}$

WHERE $\lambda_{\max} = \max(\lambda_1, \lambda_2) \Rightarrow 0 < \gamma \leq \frac{1}{6}$

PROBLEM 3

$$\text{OUTPUT SIZE} = \frac{\text{INPUT SIZE} - \text{KERNEL SIZE}}{\text{STRIDE}} + \text{NUMBER OF CHANNELS}$$

ARCHITECTURE 1

- INPUT FEATURE MAP: 32×32 , 3 FEATURE MAPS
- CONVOLUTIONAL LAYER 1: KERNEL SIZE 3×3 , 3 INPUT FEATURES MAPS, 16 OUTPUT FEATURES MAPS

$$\text{STRIDE} = 2$$

• NO PADDING

$$\text{OUTPUT SIZE} : \frac{32-3}{2} + 1 = 15 \times 15 \Rightarrow \text{OUTPUT} = 15 \times 15, 16 \text{ MAPS}$$

$$\text{LEARNABLE PARAMETERS: WEIGHTS} = 3 \cdot 3 \cdot 3 \cdot 16 = 432$$

$$\text{BIASES} = 16$$

$$\text{TOTAL} = 432 + 16 = 448$$

- POOLING LAYER 1: KERNEL SIZE 2×2

$$\text{STRIDE} = 2$$

• NO PADDING

$$\text{OUTPUT SIZE} : \frac{15}{2} = 8 \times 8 \Rightarrow \text{OUTPUT} = 8 \times 8, 16 \text{ FEATURE MAPS}$$

- CONVOLUTIONAL LAYER 2: KERNEL SIZE 3×3 , 16 INPUT FEATURES MAPS, 32 OUTPUT FEATURES MAPS

$$\text{STRIDE} = 2$$

- NO PADDING

- OUTPUT SIZE: $\frac{8-3}{2} + 1 = 4 \times 4 \Rightarrow$ OUTPUT = 4×4 , 32 FEATURE MAPS

LEARNABLE PARAMETERS: WEIGHTS = $3 \cdot 3 \cdot 16 \cdot 32 = 4608$

BIASES = 32

$$\text{TOTAL} = 4608 + 32 = 4640$$

TOTAL PARAMETERS = $448 + 4640 = 5088$

ARCHITECTURE 2

- INPUT FEATURE MAP: 32×32 , 3 FEATURE MAPS

- CONVOLUTIONAL LAYER 1: KERNEL SIZE 5×5 , 3 INPUT

FEATURE MAPS, 16 OUTPUT FEATURES MAPS

- STRIDE = 2

- NO PADDING

- OUTPUT SIZE: $\frac{32-5}{2} + 1 = 14 \times 14 \Rightarrow$ OUTPUT = 14×14 , 16 FEATURE MAPS

- LEARNABLE PARAMETERS: WEIGHTS = $5 \cdot 5 \cdot 3 \cdot 16 = 1200$

BIASES = 16

$$\text{TOTAL} = 1200 + 16 = 1216$$

- POOLING LAYER 1: KERNEL SIZE 2×2

- STRIDE = 2

- NO PADDING

- OUTPUT SIZE: $\frac{14}{2} = 7 \times 7 \Rightarrow$ OUTPUT = 7×7 , 16 FEATURE MAPS

- CONVOLUTIONAL LAYER 2: KERNEL SIZE 5×5 , 16 INPUT MAPS, 32 OUTPUT MAPS
 - STRIDE = 2
 - NO PADDING
 - OUTPUT SIZE: $\frac{7-5}{2} + 1 = 2 \times 2 \Rightarrow$ OUTPUT = 2×2 , 32 FEATURE MAPS
LEARNABLE PARAMETERS: WEIGHTS = $5 \cdot 5 \cdot 16 \cdot 32 = 12800$
BIASES = 32 TOTAL = $12800 + 32 = 12832$
 - CONVOLUTIONAL LAYER 3: KERNEL SIZE 5×5 , 32 INPUT MAPS, 64 OUTPUT MAPS
 - STRIDE = 2
 - NO PADDING
 - OUTPUT SIZE: $\frac{2-5}{2} + 1 = 4 \times 4 \Rightarrow$ OUTPUT = 4×4 , 64 FEATURE MAPS
LEARNABLE PARAMETERS: WEIGHTS = $5 \cdot 5 \cdot 32 \cdot 64 = 51200$
BIASES = 64 TOTAL = $12800 + 32 = 51264$
- TOTAL PARAMETERS = $1216 + 12832 + 51264 = 65312$

ARCHITECTURE 3

- INPUT FEATURE MAP: 32×32 , 3 FEATURE MAPS
- CONVOLUTIONAL LAYER 1: KERNEL SIZE 7×7 , 3 INPUT FEATURES MAPS, 16 OUTPUT FEATURES MAPS

• STRIDE = 2

• NO PADDING

• OUTPUT SIZE: $\frac{32-7}{2} + 1 = 13 \times 13 \Rightarrow$ OUTPUT = 13×13 , 16 FEATURE MAPS
LEARNABLE PARAMETERS: WEIGHTS = $7 \cdot 7 \cdot 3 \cdot 16 = 2352$

BIASES = 16

TOTAL = $2352 + 16 = 2368$

• POOLING LAYER 1: KERNEL SIZE 2×2

• STRIDE = 2

• NO PADDING

• OUTPUT SIZE: $\frac{13}{2} = 6 \times 6 \Rightarrow$ OUTPUT = 6×6 , 16 FEATURE MAPS
CONVOLUTIONAL LAYER 2: KERNEL SIZE 7×7 , 32 INPUT FEATURES MAPS, 64 OUTPUT FEATURES MAPS

• STRIDE = 2

• NO PADDING

- OUTPUT SIZE: $\frac{6-7}{2} + 1 = 3 \times 3 \Rightarrow$ OUTPUT $3 \times 3, 16$ FEATURE MAPS

- LEARNABLE PARAMETERS: WEIGHTS = $7 \cdot 7 \cdot 16 \cdot 32 = 25088$

BIASES = 32

TOTAL = $25088 + 32 = 25120$

- CONVOLUTIONAL LAYER 3: KERNEL SIZE 7×7 , 32 INPUT

FEATURE MAPS, 64 OUTPUT FEATURES MAPS

- STRIDE = 2

- NO PADDING

- OUTPUT SIZE: $\frac{3-7}{2} + 1 = 4 \times 4 \Rightarrow$ OUTPUT $4 \times 4, 64$ FEATURE MAPS

LEARNABLE PARAMETERS: WEIGHTS = $7 \cdot 7 \cdot 32 \cdot 64 = 100352$

BIASES = 64

TOTAL = $100352 + 64 = 100416$

TOTAL PARAMETERS = $2368 + 25120 + 100416 = 127904$

COMPARED WITH THE INITIAL 3×3 KERNEL SIZE:

- THE 5×5 KERNEL SLIGHTLY INCREASES THE CONVOLUTIONAL LAYERS

AND THE TOTAL PARAMETERS

- THE 7×7 KERNEL REACHES THE 4×4 TARGET DIMENSION WITH LESS

CONVOLUTIONAL LAYERS AND WITH AN ACCEPTABLE NUMBER OF

TOTAL PARAMETERS