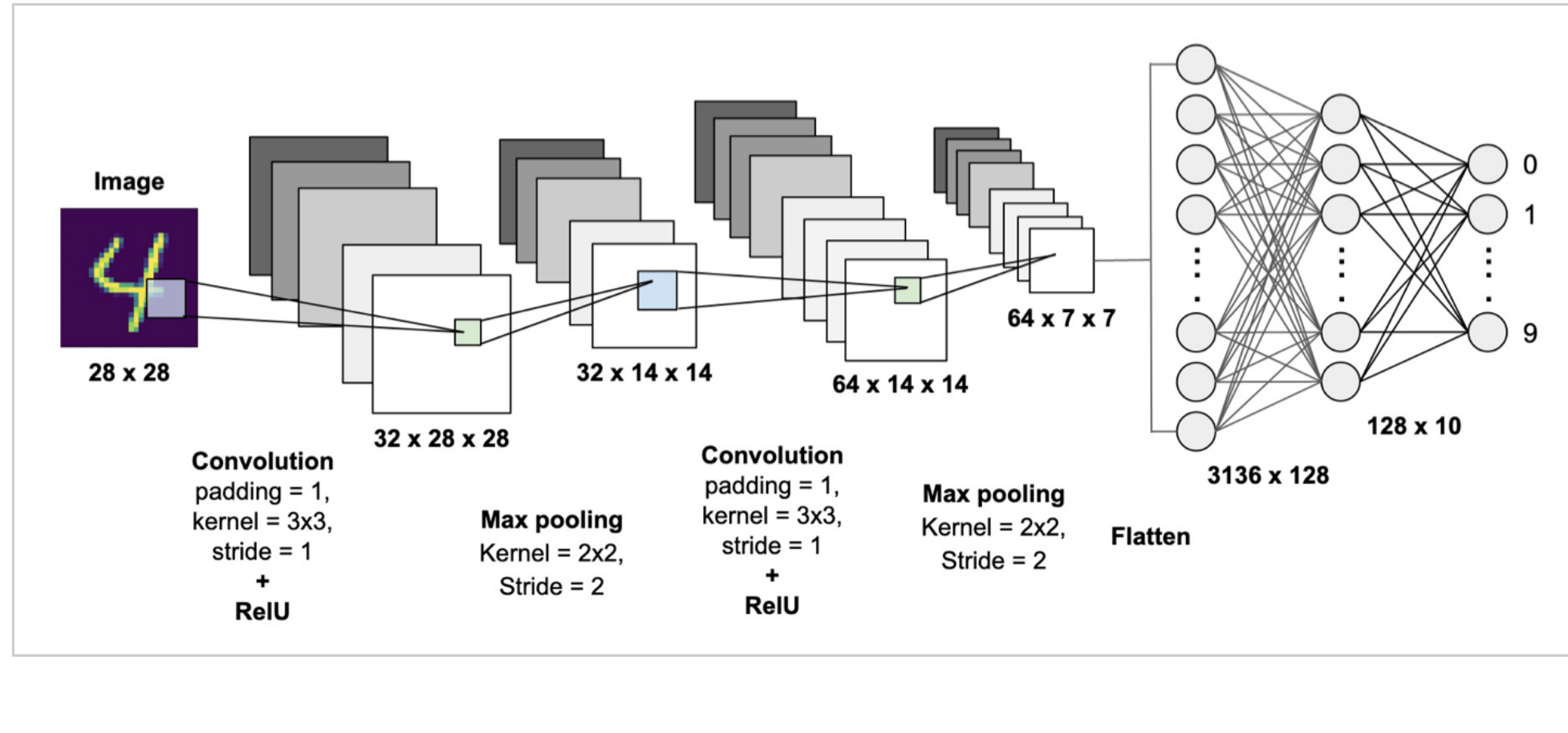


Derivation

Monday, 3 August 2020

4:23 PM



For FC layers:

Forward pass -

$$\text{Input layer} = \hat{x}$$

$$\begin{aligned} z_1 &:= \hat{x} \\ \text{for } i &= 1, 2, 3, \dots, N-1: \\ z_{i+1} &:= f_i(W_i z_i + b_i) \end{aligned}$$

$$\text{(also, let } a_i = W_i z_i + b_i \text{)}$$

Backward pass -

$$\text{Output vector} = z_N$$

$$\text{Loss} = L(\hat{y}, \theta)$$

$$\frac{\partial L(\hat{y}, \theta)}{\partial z_N^n} = \delta_N^n$$

$$\begin{aligned} \text{eg. } \frac{\partial L(\hat{y}, \theta)}{\partial z_N^1} &= \delta_N^1, \\ \text{and } \frac{\partial L(\hat{y}, \theta)}{\partial z_N} &= [\delta_N^1, \delta_N^2, \dots, \delta_N^N] \end{aligned}$$

To prove

$$\epsilon_i = \frac{\partial L(\hat{y}, \theta)}{\partial a_i} = \delta_i \circ f'(a_i) \quad - (1)$$

$$b_{i-1} = W_{i-1}^T \cdot \epsilon_i \quad - (2)$$

$$db_i = \epsilon_i \quad - (3)$$

$$dW_i = \epsilon_i \otimes z_{i+1} \quad - (4)$$

where,

$$\nabla L(\hat{y}, \theta) = [db_1, db_2, db_3, \dots, db_{N-1}, dW_1, dW_2, dW_3, \dots, dW_{N-1}]$$

and the backprop algorithm,

$$\begin{aligned} \text{Compute } \delta_{N-1} &:= \frac{\partial L(\hat{y}, \theta)}{\partial z_N} \\ \epsilon_{N-1} &:= \delta_{N-1} \circ f'(a_{N-1}) \\ \text{for } i &= N-1, N-2, \dots, 1: \\ db_i &:= \epsilon_i \\ dW_i &:= \epsilon_i \otimes z_{i+1} \\ b_{i-1} &:= W_{i-1}^T \cdot \epsilon_i \\ \epsilon_{i-1} &:= \delta_{i-1} \circ f'(a_{i-1}) \end{aligned}$$

$$\begin{aligned} (1) \quad \frac{\partial L(\hat{y}, \theta)}{\partial a_i} &= \frac{\partial L(\hat{y}, \theta)}{\partial z_i} \cdot \frac{\partial z_i}{\partial a_i} \\ &= \delta_i \cdot f'(a_i) \quad \blacksquare \end{aligned}$$

(2) ϵ_i is the derivative of loss function w.r.t input of neurons in i^{th} layers.

We want δ_{i-1} which is derivative of loss function with respect to output of neurons in $i-1^{\text{th}}$ layers.

Consider z_{i-1}^1 (first neuron in layer i)

and a_i .

$$a_i = W_{i-1} z_{i-1} + b_{i-1}$$

$$a_i^1 = W_{i-1}^1 \cdot z_{i-1} + b_{i-1}^1$$

$$\frac{\partial a_i^1}{\partial z_{i-1}^1} = \frac{\partial W_{i-1}^1 z_{i-1} + b_{i-1}^1}{\partial z_{i-1}^1}$$

$$\frac{\partial a_i^1}{\partial z_{i-1}^1} = W_{i-1}^{11} \quad \text{and in general,}$$

$$\frac{\partial a_i^k}{\partial z_{i-1}^k} = W_{i-1}^{kk}$$

$$\frac{\partial L(\hat{y}, \theta)}{\partial z_{i-1}^k} = \sum_{j=1}^n \frac{\partial L(\hat{y}, \theta)}{\partial a_i^j} \cdot \frac{\partial a_i^j}{\partial z_{i-1}^k}$$

$$= \sum_{j=1}^n \epsilon_i^j W_{i-1}^{jk}$$

$$\frac{\partial L(\hat{y}, \theta)}{\partial z_{i-1}} = \begin{bmatrix} \sum_{j=1}^n \epsilon_i^j W_{i-1}^{j1} \\ \sum_{j=1}^n \epsilon_i^j W_{i-1}^{j2} \\ \sum_{j=1}^n \epsilon_i^j W_{i-1}^{j3} \\ \vdots \\ \sum_{j=1}^n \epsilon_i^j W_{i-1}^{jk} \end{bmatrix}$$

where $\sum_{j=1}^n \epsilon_i^j W_{i-1}^{jk}$ is the dot product of ϵ_i and k^{th} column of $W_{i-1} \Rightarrow \sum_{j=1}^n \epsilon_i^j W_{i-1}^{jk} = \epsilon_i \cdot W_{i-1}^{T1}$

$$\therefore \frac{\partial L(\hat{y}, \theta)}{\partial z_{i-1}} = W_{i-1}^T \epsilon_i \quad \blacksquare$$

$$(3) \quad \frac{\partial L(\hat{y}, \theta)}{\partial b_i} = \frac{\partial L(\hat{y}, \theta)}{\partial a_i} \cdot \frac{\partial a_i}{\partial b_i}$$

$$= \epsilon_i \cdot 1$$

$$= \epsilon_i \quad \blacksquare$$

$$(4) \quad \frac{\partial L(\hat{y}, \theta)}{\partial W_{i-1}^{ab}} = \epsilon_i^b \cdot z_{i-1}^a$$

Expanding gives the following:

$$\frac{\partial L(\hat{y}, \theta)}{\partial W_{i-1}^{ab}} = \begin{bmatrix} \epsilon_i^1 z_{i-1}^1 & \epsilon_i^2 z_{i-1}^1 & \dots & \epsilon_i^b z_{i-1}^1 \\ \epsilon_i^1 z_{i-1}^2 & \epsilon_i^2 z_{i-1}^2 & \dots & \epsilon_i^b z_{i-1}^2 \\ \vdots & \vdots & \ddots & \vdots \\ \epsilon_i^1 z_{i-1}^a & \epsilon_i^2 z_{i-1}^a & \dots & \epsilon_i^b z_{i-1}^a \end{bmatrix}$$

$$= \epsilon_i \otimes z_{i-1} \quad \blacksquare$$