

CNNs From Scratch

Advay Singh

CONTENTS

0.1.	Introduction	4
0.2.	Image Processing	4
0.3.	Fully-Connected Multi Layer Perception	4
0.4.	Efficiency Testing	10
0.5.	Application in Full-Stack Web Application	10

0.1. Introduction

0.2. Image Processing

0.2.1. Convolution.

0.2.2. Pooling.

0.3. Fully-Connected Multi Layer Perception

Now, with our inputs in place, will begin discussion regarding the a Fully-Connected Multi Layer Perception.

- (1) $x_i = ith$ input
- (2) $w_{ij}^k =$ weight from ith input to jth output in the kth layer
- (3) $b_{ji}^k =$ bias from ith input to jth output in the kth layer
- (4) $z_j^k = jth$ input for activation at layer k
- (5) $\alpha_j^k = jth$ activated val at layer k
- (6) $y_i = ith$ true label
- (7) $l =$ number of hidden layers
- (8) $n^k =$ number of activations for layer k
- (9) $Y_i = ith$ predicted output
- (10) $E_i = ith$ error
- (11) $\Delta w =$ change to weights
- (12) $\Delta b =$ change to biases
- (13) $a =$ learning rate.

0.3.1. Forward Propagation. The x and y are both lists of inputs and labels respectively. Other terms, require derivation.

0.3.1.1. *Weights and Biases.* Weights and Biases are both 2-dimensional arrays. Note that there are multiple of these, intact, $l + 1$ of each in a Fully-Connected Neural Network.

Initially, these arrays are filled with random values which are eventually re-defined through gradient-decent. The length of these arrays can be determined through matrix multiplication. The matrix w^k 's aXb where $a = n^{k-1}$ and $b = n^k$. For this, the initial layer $k - 1$ is reshaped to $n^{k-1} \cdot 1$ and k is reshaped $1 \cdot n^k$. It's a similar process for the biases.

0.3.1.2. *Zs (Activation Inputs)*. With the weights biases, and inputs in-place, we may derive our zs .

$$(14) \quad z_j^k = \sum_i x_i^{k-1} w_{ij}^k.$$

Note we need not derive zs for $k \geq l$.

0.3.1.3. *Activation Function*. There are many activation functions. ReLU, Sigmoid, and TanH to name a few. The following defines the **Sigmoid Activation Function**, $\sigma(z) : \mathbb{R} \rightarrow (0, 1)$. The reason for these function is to **introduce non-linearity** in the dataset. The more hidden layers there are, the more complex problems the CNN is capable of solving. The cost? Increased computation time.

$$(15) \quad \alpha_j^k = \sigma(z_j^k) = \frac{1}{1 + e^{-(z_j^k)}}.$$

Note that $\forall \alpha, \alpha \in (0, 1)$. Such normalization is often seen in activation functions. However, the purpose is the non-linear nature of the sigmoid graph.

We also have the **ReLU** activation function: $R'(x) =$

$$\begin{cases} 0 & x < 0 \\ x & x \geq 0. \end{cases}$$

More on this as we compute their derivatives.

0.3.1.4. *Softmax Activation Function*. Because $\forall \alpha, \alpha \in (0, 1)$ it's impossible to use them to predict values outside that range or non-numerical objects. In the latter case, we apply **Softmax** to the z^l list while a simple regression activation function often suffices for numerical classifications.

$$(16) \quad Y_i = S(x_i) = \frac{e^{x_i}}{\sum_j e^{x_j}}.$$

Due to the Softmax function, $\forall Y_i, Y_i \in (0, 1)$ meaning that $S(x) : \mathbb{R} \rightarrow [0, 1]$. To compare this to y_i , we convert $\forall y_i \in \{0, 1\}$. Where only the true classification $y_{true} = 1$ while every other $y_i = 0$.

With this, we have the forward propagation function can be generated. The inputs for the function are x , the inputs and true labels. The outputs: Y an array of the predicted outputs. Additionally, other functors and parameters can be used, including which activation function, etc.

0.3.2. Back Propagation & Gradient Decent. Firstly, we must convert the labels y_i to values $\in \{0, 1\}$

0.3.2.1. *Error Function.* The error function compares $Y_i \wedge y_i$ and is responsible for generating $\Delta w \wedge \Delta b$ for gradient decent. For this, we define a cost function.

$$(17) \quad E_i = C(i) = \frac{(Y_i - y_i)^2}{2}$$

In our computations, we use

$$(18) \quad E_{tot} = C(n^l) = \sum_i^{n^l} \frac{(Y_i - y_i)^2}{2}$$

Therefore, through each iteration, we derive $\Delta w_{ji}^l \wedge \Delta b_{ji}^l$ based on E_i . And update the w and b through back propagation. Iterating through the training set is considered one **epoch**. Notably, every epoch increases accuracy for regular gradient decent, however, **Stochastic Gradient Decent**'s accuracy plateaus over time. More on the later.

0.3.2.2. *Calculating Δw And Δb .* Firstly, we begin with our final layer. $\frac{\delta E_i}{\delta w_{ij}}$ and $\frac{\delta E_i}{\delta b_{ij}}$.
Breaking down the function:

$$(19) \quad \frac{\delta E_i}{\delta w_{ij}} = \frac{\delta E_i}{\delta S(z_i)} \cdot \frac{\delta S(z_i)}{\delta z_i^l} \cdot \frac{\delta z_i^l}{\delta w_{ij}}.$$

Hence we derive all the parts.

$$(20) \quad \frac{\delta E_i}{\delta S(x_i)} = \frac{\delta}{\delta S(x_i)} \frac{(Y_i - y_i)^2}{2}$$

$$(21) \quad = \frac{\delta}{\delta S(x_i)} \frac{(S(x_i) - y_i)^2}{2}$$

$$(22) \quad = (S(x_i) - y_i)$$

$$(23) \quad = (Y_i - y_i).$$

Note that for our purposes, x_i refers to z_i^l . For the following calculation, l is implicit for z_i because we are deriving for the final layer.

$$(24) \quad \frac{S(z_i)}{\delta z_i} = S(z_i)(1 - S(z_i)).$$

$J_{softmax}$ is known as the **Jacobian Matrix** and as modeled like such

$$\begin{bmatrix} \frac{\delta S(1)}{\delta z_1} & \frac{\delta S(1)}{\delta z_2} & \dots & \frac{\delta S(1)}{\delta z_i} \\ \frac{\delta S(2)}{\delta z_1} & \frac{\delta S(2)}{\delta z_2} & \dots & \frac{\delta S(2)}{\delta z_i} \\ \vdots & \vdots & \vdots & \vdots \\ \frac{\delta S(i)}{\delta z_1} & \frac{\delta S(i)}{\delta z_2} & \dots & \frac{\delta S(i)}{\delta z_i} \end{bmatrix}$$

where $i = n^l$. We will derive the $S(x_i)$ function later.

Now we derive Δw for the final layer (layer going from j th input to i th output).

$$(25) \quad \frac{\delta z_i}{\delta w_{ij}} = \frac{\delta}{\delta w_{ij}} \sum_p^{n^{i-1}} x_p w_{ip} + b_{ip}$$

Note that for $f(w) = \frac{\delta}{\delta w_{ij}} \sum_p^{n^{i-1}} x_p w_{ip} + b_{ip}$, $(p = j) \iff f(z) = x_j$ and $(j \neq i \iff f(z) = 0)$. Therefore,

$$(26) \quad \frac{\delta}{\delta w_{ij}} \sum_p^{n^{i-1}} x_p w_{ip} + b_{ip} = x_j$$

Note that the $\frac{\delta E_i}{\delta S(x_i)} \cdot \frac{\delta S(x_i)}{\delta z_i^l}$ remain the same for Δb , however, $\frac{\delta z_i}{\delta b_{ij}} = 1$. So we have

$$(27) \quad \Delta w_{ij} = \frac{\delta E_i}{\delta w_{ij}} = (Y_i - y_i) \cdot S(z_i)(1 - S(z_i)) \cdot x_j$$

$$(28) \quad \Delta b_{ij} = \frac{\delta E_i}{\delta b_{ij}} = (Y_i - y_i) \cdot S(z_i)(1 - S(z_i)).$$

0.3.2.3. Softmax Derivative and Jacobin Matrix. Our goal in this segment is to define $\frac{dS(z_i)}{dz_i}$. Firstly,

$$(29) \quad \frac{\delta \ln(S(z_i))}{\delta z_i} = \frac{1}{S(z_i)} \frac{\delta S(z_i)}{\delta z_i}.$$

Hence, we initially compute $\frac{\delta \ln(S(z_i))}{\delta z_i}$.

$$(30) \quad \frac{\delta \ln(S(z_i))}{\delta z_i} = \frac{\delta}{\delta z_i} \ln\left(\frac{e^{z_i}}{\sum_j^{n^j} e^{z_j}}\right)$$

$$(31) \quad = \frac{\delta}{\delta z_i} \ln(e^{z_i}) - \frac{\delta}{\delta z_i} \ln\left(\sum_j^{n^j} e^{z_j}\right)$$

$$(32) \quad = 1 - \frac{1}{\sum_j^{n^j} e^{z_j}} \frac{\delta}{\delta z_i} \left(\sum_j^{n^j} e^{z_j}\right).$$

Note that for $g(z) = \frac{\delta}{\delta z_i} \left(\sum_j^{n^j} e^{z_j}\right)$, $(j = i) \iff g(z) = 1$ and $(j \neq i) \iff g(z) = 0$. Therefore,

$$(33) \quad \frac{\delta \ln(S(z_i))}{\delta z_i} = 1 - S(z_i).$$

And so we have that

$$(34) \quad \frac{\delta S(z_i)}{\delta z_i} = S(z_i) \frac{\delta \ln(S(z_i))}{\delta z_i}$$

$$(35) \quad = S(z_i)(1 - S(z_i)).$$

There is also an idea of **Temperature** T which adds another layer of non-linearity to S . A smaller temperature favors greater z_i and vice-versa.

$$(36) \quad Y_i = S(x_i) = \frac{e^{\frac{x_i}{t}}}{\sum_j^{n^j} e^{\frac{x_j}{t}}}$$

0.3.2.4. *Activation Function Derivative.* The derivative of the activation function plays a similar role as the Softmax function derivative, just for

earlier, hidden layers. For now, we will be focusing on the Sigmoid function.

$$(37) \quad \frac{\delta}{\delta z_i} \alpha_i = \frac{\delta}{\delta z_i} \sigma(z_i)$$

$$(38) \quad = \frac{\delta}{\delta z_i} \frac{1}{1 + e^{-(z_i)}}$$

$$(39) \quad = \frac{\delta}{\delta z_i} (1 + e^{-(z_i)})^{-1}$$

$$(40) \quad = \frac{-1}{(1 + e^{-(z_i)})^2} \frac{\delta}{\delta z_i} e^{-(z_i)}$$

$$(41) \quad = \frac{-1}{(1 + e^{-(z_i)})^2} - e^{-(z_i)}$$

$$(42) \quad = \frac{e^{-(z_i)}}{(1 + e^{-(z_i)})^2}$$

$$(43) \quad = \frac{1}{1 + e^{-(z_i)}} \cdot \frac{e^{-(z_i)}}{1 + e^{-(z_i)}}$$

$$(44) \quad = \frac{1}{1 + e^{-(z_i)}} \cdot \left(1 - \frac{1}{1 + e^{-(z_i)}}\right)$$

$$(45) \quad = \sigma(z_i)(1 - \sigma(z_i)).$$

Look familiar? It's the same as S .

Notably, calculating the derivative of the ReLU function is much simpler.

$$R'(x) =$$

$$\begin{cases} 0 & x < 0 \\ 1 & x \geq 0. \end{cases}$$

Computationally, this often yields a much greater $\Delta w \wedge \Delta b$ allowing for a much faster back propagation process. Hence, we can now proceed to

calculate $\Delta w \wedge \Delta b$ for the hidden layers (layers $< l$). The following computes δw_{ji}^1 in a Neural Network with two layers ($l = 2$) with p inputs.

$$(46) \quad \frac{\delta E_i}{\delta w_{jp}^1} = \frac{\delta E_i}{\delta S(z_i^l)} \cdot \frac{\delta S(z_i^l)}{\delta z_i^l} \cdot \frac{\delta z_i^l}{\delta w_{jp}^1}$$

$$(47) \quad = \frac{\delta E_i}{\delta S(z_i^l)} \cdot \frac{\delta S(z_i^l)}{\delta z_i^l} \cdot w_{ij}^l \cdot \frac{\delta a_j^1}{\delta w_{jp}^1}$$

$$(48) \quad = \frac{\delta E_i}{\delta S(z_i^l)} \cdot \frac{\delta S(z_i^l)}{\delta z_i^l} \cdot w_{ij}^l \cdot \frac{\delta \sigma(z_j^1)}{\delta z_j^1} \cdot \frac{\delta z_j^1}{\delta w_{jp}^1}$$

$$(49) \quad = (Y_i - y_i) \cdot S(z_i^2)(1 - S(z_i^2)) \cdot w_{ij}^2 \cdot \sigma(z_j^1)(1 - \sigma(z_j^1)) \cdot x_p.$$

This way, we can also obtain Δb .

$$(50) \quad \frac{\delta E_i}{\delta b_{jp}^1} = \frac{\delta E_i}{\delta S(z_i^l)} \cdot \frac{\delta S(z_i^l)}{\delta z_i^l} \cdot \frac{\delta z_i^l}{\delta b_{jp}^1}$$

$$(51) \quad = \frac{\delta E_i}{\delta S(z_i^l)} \cdot \frac{\delta S(z_i^l)}{\delta z_i^l} \cdot w_{ij}^l \cdot \frac{\delta a_j^1}{\delta b_{jp}^1}$$

$$(52) \quad = \frac{\delta E_i}{\delta S(z_i^l)} \cdot \frac{\delta S(z_i^l)}{\delta z_i^l} \cdot w_{ij}^l \cdot \frac{\delta \sigma(z_j^1)}{\delta z_j^1} \cdot \frac{\delta z_j^1}{\delta b_{jp}^1}$$

$$(53) \quad = (Y_i - y_i) \cdot S(z_i^2)(1 - S(z_i^2)) \cdot w_{ij}^2 \cdot \sigma(z_j^1)(1 - \sigma(z_j^1)).$$

Now that we know how to compute $\Delta w \wedge \Delta b$ for all matrices, it's essential to understand how to automate the process for any computational application.

0.3.2.5. *Stochastic Gradient Decent.* The learning rate a has an impact but can cause for less overall accuracy.

0.4. Efficiency Testing

0.5. Application in Full-Stack Web Application