

Mathematics for Neural Networks

styxofsyntax

July 9, 2024

1 Introduction

This document explains the mathematical concepts used in the implementation of a neural network from scratch.

2 Equations

2.1 Calculating net input

m = number of examples

n = number of inputs

k = number of neurons

$$\mathbf{Z} = \mathbf{X} \cdot \mathbf{W} + \mathbf{j} \otimes \mathbf{b}$$

- \mathbf{X} is the input matrix with dimensions $m \times n$.
- \mathbf{W} is the weight matrix with dimensions $n \times k$
- \mathbf{j} is all ones with dimensions $m \times 1$
- \mathbf{b} is the bias matrix with dimensions $1 \times k$

The result \mathbf{Z} will have dimensions $m \times k$.

$$\mathbf{X} = \begin{bmatrix} x_1^1 & x_2^1 & x_3^1 & \cdots & x_n^1 \\ x_1^2 & x_2^2 & x_3^2 & \cdots & x_n^2 \\ x_1^3 & x_2^3 & x_3^3 & \cdots & x_n^3 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_1^m & x_2^m & x_3^m & \cdots & x_n^m \end{bmatrix}_{m \times n}$$
$$\mathbf{W} = \begin{bmatrix} w_1^1 & w_1^2 & w_1^3 & \cdots & w_1^k \\ w_2^1 & w_2^2 & w_2^3 & \cdots & w_2^k \\ w_3^1 & w_3^2 & w_3^3 & \cdots & w_3^k \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ w_n^1 & w_n^2 & w_n^3 & \cdots & w_n^k \end{bmatrix}_{n \times k}$$
$$\mathbf{J} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}_{m \times 1}$$
$$\mathbf{b} = [b_1 \quad b_2 \quad b_3 \quad \cdots \quad b_k]_{1 \times k}$$

Note: The kronecker product is not required to be explicitly implemented in code, as \mathbf{B} will be added to all rows of \mathbf{W} regardless.

2.2 Gradients

- **Gradient of the Error with respect to the activation of the last layer**

The error functions is given by:

$$E = \frac{1}{m} \sum_{i=0}^m \left(\mathbf{a}_i^{[L]} - \mathbf{y}_i \right)^2$$

Where $\mathbf{a}_i^{[L]}$ is the activation of the last layer and \mathbf{y}_i is the true output, over one example.

The gradient of the error with respect to the activations of the last layer is:

$$\begin{aligned} \frac{\partial E}{\partial \mathbf{a}_j^{[L]}} &= \frac{1}{m} \cdot 2 \sum_{i=0}^m \left(\mathbf{a}_i^{[L]} - \mathbf{y}_i \right) \cdot 1 \\ &= \frac{2}{m} \left(\mathbf{a}_j^{[L]} - \mathbf{y}_j \right) \end{aligned}$$

Where $\mathbf{a}_j^{[L]}$ is one particular activation of the last layer over all examples.

Matrix notation of gradient of the error with respect to activations of the last layer is:

$$\frac{\partial E}{\partial \mathbf{A}^{[L]}} = \frac{2}{m} \left(\mathbf{A}^{[L]} - \mathbf{Y} \right) \quad (\text{dimensions: } m \times k)$$

Where $\mathbf{A}^{[L]}$ is the activation of the last layer and \mathbf{Y} are true outputs, over all examples.

- **Gradient of the activations of the last layer with respect to the net input of the last layer**

$$\frac{\partial \mathbf{A}^{[L]}}{\partial \mathbf{Z}^{[L]}} = h'(\mathbf{Z}^{[L]}) \quad (\text{dimensions: } m \times k)$$

where $\mathbf{Z}^{[L]}$ are the net inputs over all examples.

- **Gradient of the net input of the last layer with respect to the activations of the previous layer**

Given that:

$$\mathbf{Z}^{[L]} = \mathbf{A}^{[L-1]} \cdot \mathbf{W}^{[L]} + \mathbf{B}^{[L]}$$

The partial derivative of $\mathbf{Z}^{[L]}$ with respect to $\mathbf{A}^{[L-1]}$ is:

$$\frac{\partial \mathbf{Z}^{[L]}}{\partial \mathbf{A}^{[L-1]}} = \mathbf{W}^{[L]} \quad (\text{dimensions: } n \times k)$$

- **Gradient of the net input of the last layer with respect to the weights**

The partial derivative of $\mathbf{Z}^{[L]}$ with respect to $\mathbf{W}^{[L]}$ is:

$$\frac{\partial \mathbf{Z}^{[L]}}{\partial \mathbf{W}^{[L]}} = \mathbf{A}^{[L-1]} \quad (\text{dimensions: } m \times n)$$

where $\mathbf{A}^{[L-1]}$ are the activations of previous layer over all examples.

- **Gradient of the net input of the last layer with respect to the biases**

The partial derivative of $\mathbf{Z}^{[L]}$ with respect to $\mathbf{b}^{[L]}$ is:

$$\frac{\partial \mathbf{Z}^{[L]}}{\partial \mathbf{b}^{[L]}} = \mathbf{1}$$

2.3 Combining Gradients

Combining these gradients using the chain rule:

To generalize the equations over all layers, $\boldsymbol{\partial}(l)$ can be used in place of $\frac{\partial E}{\partial \mathbf{A}^{[l]}}$ where l indicates l -th Layer.

The gradient of the last layer $\boldsymbol{\partial}(L)$ is given as,

$$\boldsymbol{\partial}(L) = \frac{2}{m} (\mathbf{A}^{[L]} - \mathbf{Y})$$

- **Gradient with respect to the previous layer's activations**

$$\begin{aligned} \boldsymbol{\partial}(l-1) &= \left(\frac{\partial \mathbf{A}^{[l]}}{\partial \mathbf{Z}^{[l]}} \times \boldsymbol{\partial}(l) \right) \cdot \left(\frac{\partial \mathbf{Z}^{[l]}}{\partial \mathbf{A}^{[l-1]}} \right)^T \\ &= \left[h'(\mathbf{Z}^{[l]}) \times \boldsymbol{\partial}(l) \right] \cdot (\mathbf{W}^{[l]})^T \quad (\text{dimensions: } m \times n) \end{aligned}$$

- **Gradients with respect to the weights**

$$\begin{aligned} \frac{\partial E}{\partial w^{[l]}} &= \left(\frac{\partial \mathbf{Z}^{[l]}}{\partial w^{[l]}} \right)^T \cdot \left(\frac{\partial \mathbf{A}^{[l]}}{\partial \mathbf{Z}^{[l]}} \times \boldsymbol{\partial}(l) \right) \\ &= (\mathbf{A}^{[l-1]})^T \cdot \left[h'(\mathbf{Z}^{[l]}) \times \boldsymbol{\partial}(l) \right] \quad (\text{dimensions: } n \times k) \end{aligned}$$

- **Gradients with respect to the biases**

$$\begin{aligned} \frac{\partial E}{\partial \mathbf{B}^{[l]}} &= \frac{\partial \mathbf{Z}^{[l]}}{\partial \mathbf{b}^{[l]}} \cdot \left(\frac{\partial \mathbf{A}^{[l]}}{\partial \mathbf{Z}^{[l]}} \times \boldsymbol{\partial}(l) \right) \\ \text{since } \frac{\partial \mathbf{Z}^{[l]}}{\partial \mathbf{b}^{[l]}} &= 1, \text{ we have} \\ &= \frac{\partial \mathbf{A}^{[l]}}{\partial \mathbf{Z}^{[l]}} \times \boldsymbol{\partial}(l) \\ &= h'(\mathbf{Z}^{[l]}) \times \boldsymbol{\partial}(l) \quad (\text{dimensions: } m \times k) \end{aligned}$$

2.4 Updating parameters

- Updating Weights

$$\mathbf{W}'^{[l]} = \mathbf{W}^{[l]} - \alpha \frac{\partial E}{\partial \mathbf{W}^{[l]}}$$

- Updating Biases

The partial derivatives over all examples are added,

$$\mathbf{b}'^{[l]} = \mathbf{b}^{[l]} - \alpha \sum_{i=0}^m \frac{\partial E}{\partial \mathbf{B}_i^{[l]}}$$

Here α is the learning rate.