

Neural Network Hand Calculations

Hanhee Lee

May 2024

1 Introduction

This paper explores simple regression models to demonstrate the underlying mechanics by implementing these models both through MATLAB and using a paper-and-pen approach for a shallow and deep neural network.

2 Notation

The following table summarizes the notation used in the neural network models:

Table 1: Detailed Explanation of Variables

Header	Dimension	Explanation
Superscripts		
$[l]$	1	Current layer
(i)	1	ith training example
Subscripts		
j	1	jth node of the current layer
k	1	jth node of the previous layer
Sizes		
m	1	Number of training examples in the dataset
n_x	1	Number of nodes in the input layer
n_y	1	Number of nodes in the output layer
$n^{[l]}$	1	Number of nodes in the current layer
$n^{[l-1]}$	1	Number of nodes in the previous layer
L	1	Number of layers in the network
Objects		
\mathbf{X}	$n_x \times m$	Input matrix
$\mathbf{x}^{(i)}$	$n_x \times 1$	ith example represented as a column vector
$\mathbf{W}^{[l]}$	$n^{[l]} \times n^{[l-1]}$	Weight matrix of the current layer
$z_j^{[l](i)}$	1	A weighted sum of the activations of the previous layer, shifted by a bias
$w_{j,k}^{[l]}$	1	A weight that scales the k th activation of the previous layer
$b^{[l]}$	$n^{[l]} \times 1$	Bias vector in the current layer
$b_j^{[l]}$	$n^{[l]} \times 1$	Bias in the current layer
$a_j^{[l](i)}$	1	An activation in the current layer
$a_k^{[l-1](i)}$	1	An activation in the previous layer
$g_j^{[l]}$	1	An activation function used in the current layer

3 Neural Network Formulas

This section describes the foundational equations used in neural networks, detailing the computation involved in forward propagation.

Layer-Wise Computation

$$z_j^{[l](i)} = \sum_k w_{j,k}^{[l]} a_k^{[l-1](i)} + b_j^{[l]} \quad (1)$$

- This equation calculates the weighted sum of the activations from the previous layer, combined with a bias term, to produce the pre-activation value for neuron j in layer l for a given input i .

$$a_k^{[l-1](i)} = g_j^{[l]} \left(z_1^{[l](i)}, \dots, z_j^{[l](i)}, \dots, z_{n^{[l]}}^{[l](i)} \right) \quad (2)$$

- Here, the activation of the k th neuron in layer $l - 1$ for input i is calculated using the activation function $g_j^{[l]}$, which is applied to the vector of all pre-activation values from layer l .

$$L = \sqrt{\frac{1}{m} \sum_{i=1}^m (\hat{y}_i - y_i)^2} \quad (3)$$

- This equation represents the root mean squared error, a common cost function used to measure the difference between the predicted outputs (\hat{y}_i) and the actual targets (y_i) over all m training examples.

4 Simple Regression Models

To demonstrate the neural network's learning mechanism, we will use simple regression models. These models illustrate how a neural network can be trained to fit data according to specific mathematical relationships.

Model Descriptions

The first model we will consider is described by the following linear relationship:

$$y = 4x_1 + x_2^2 + x_3 \quad (4)$$

- This model expresses y as a linear combination of x_1 and x_3 , with x_2 contributing quadratically. This demonstrates how different types of input features can be integrated into the prediction.

Additional models can be added here following the same format, explaining the mathematical relation and its potential learning implications for a neural network.

5 Neural Network Diagram

This section details the architecture of a simple neural network, which is diagrammed below:

- **Input Layer:** Consists of 3 nodes, representing the input features (x_1, x_2, x_3) .
- **Hidden Layer:** Comprises 4 nodes, which facilitate the learning of non-linear relationships.
- **Output Layer:** Contains a single node, which outputs the prediction of the network.

Please refer to the accompanying diagram for a visual representation of the network's structure:

6 Hyperparameters

Hyperparameters are configurations external to the model that influence how the network is structured and trained. Hyperparameters play a crucial role in determining the model's performance by affecting how quickly and effectively it learns from the training data.

Table 2: Detailed Explanation of Variables

Hyperparameter	Description
Number of hidden layers	1
Optimizer	Stochastic gradient descent
Number of nodes in the hidden layer	4
Activation function of the hidden layer	ReLU function
Activation function of the output layer	Sigmoid function
Loss function	Root mean squared error
Learning rate	1
Number of epochs	1

7 Example: Mapping Function in Neural Network

This example illustrates how a specific input maps through a neural network to produce an output. The neural network is structured with three input nodes and one output node. The inputs for this example are chosen with feature dimensions of 2, 1, and 3, respectively.

Input Layer: The network receives a single training point with features:

- $x_1 = 2$
- $x_2 = 1$
- $x_3 = 3$

These inputs are fed into the network to process through the neural architecture.

8 Forward Propagation

Forward Propagation involves processing input data through the network from the input to the output layer using current weights and biases, generating predictions that are used to calculate the error against actual targets.

1. Input to Hidden Layer:

Let $\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$, where each x_i is an input feature.

- Here, \mathbf{x} represents the input vector to the network, consisting of features x_1 , x_2 , and x_3 . These features are the data points that you want the network to learn from.

Let $\mathbf{W}^{[1]} = \begin{bmatrix} w_{1,1}^{[1]} & w_{1,2}^{[1]} & w_{1,3}^{[1]} \\ w_{2,1}^{[1]} & w_{2,2}^{[1]} & w_{2,3}^{[1]} \\ w_{3,1}^{[1]} & w_{3,2}^{[1]} & w_{3,3}^{[1]} \\ w_{4,1}^{[1]} & w_{4,2}^{[1]} & w_{4,3}^{[1]} \end{bmatrix}$, be the weight matrix connecting the input to the hidden layer.

- $\mathbf{W}^{[1]}$ is the weight matrix associated with the first layer. Each element, such as $w_{1,1}^{[1]}$, represents the weight connecting the input node x_1 to the first neuron in the hidden layer.

Let $\mathbf{b}^{[1]} = \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix}$, for the hidden layer.

- $\mathbf{b}^{[1]}$ represents the bias vector for the hidden layer, with each entry like $b_1^{[1]}$ adding a bias term to the corresponding neuron's output. This helps to adjust the threshold at which the neuron activates.

$$\mathbf{z}^{[1]} = \mathbf{W}^{[1]}\mathbf{x} + \mathbf{b}^{[1]}$$

- This equation computes the linear combination of inputs and weights, adjusted by the bias. The result, $\mathbf{z}^{[1]}$, is the pre-activation output of the hidden layer.

$$\begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]} \end{bmatrix} = \begin{bmatrix} w_{1,1}^{[1]} & w_{1,2}^{[1]} & w_{1,3}^{[1]} \\ w_{2,1}^{[1]} & w_{2,2}^{[1]} & w_{2,3}^{[1]} \\ w_{3,1}^{[1]} & w_{3,2}^{[1]} & w_{3,3}^{[1]} \\ w_{4,1}^{[1]} & w_{4,2}^{[1]} & w_{4,3}^{[1]} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix}$$

$$\begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]} \end{bmatrix} = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 2 \\ 1 \\ 3 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

- This step shows the explicit matrix multiplication and addition for the given example. It is a practical computation where each neuron's input is the sum of products of each input feature and the corresponding weight plus a bias term.

$$\mathbf{z}^{[1]} = \begin{bmatrix} 6 \\ 6 \\ 6 \\ 6 \end{bmatrix}$$

- The result is a vector of pre-activation values for each neuron in the hidden layer.

2. Activation in Hidden Layer:

$$\mathbf{a}^{[1]} = \text{ReLU}(\mathbf{z}^{[1]})$$

- The ReLU function is applied to each pre-activation value. This non-linear function outputs the input directly if it is positive; otherwise, it outputs zero.

$$\begin{bmatrix} a_1^{[1]} \\ a_2^{[1]} \\ a_3^{[1]} \\ a_4^{[1]} \end{bmatrix} = \text{ReLU} \left(\begin{bmatrix} 6 \\ 6 \\ 6 \\ 6 \end{bmatrix} \right)$$

- Since all inputs are positive (6), the ReLU output matches the input. This vector represents the activated output of the hidden layer.

$$\mathbf{a}^{[1]} = \begin{bmatrix} 6 \\ 6 \\ 6 \\ 6 \end{bmatrix}$$

- These are the activation values from the hidden layer that will be fed to the next layer.

3. Hidden Layer to Output Layer:

$$\text{Let } \mathbf{W}^{[2]} = \begin{bmatrix} w_{1,1}^{[2]} & w_{1,2}^{[2]} & w_{1,3}^{[2]} & w_{1,4}^{[2]} \end{bmatrix} \text{ and } \mathbf{b}^{[2]} = \begin{bmatrix} b_1^{[2]} \end{bmatrix}$$

- $\mathbf{W}^{[2]}$ and $\mathbf{b}^{[2]}$ are the weight vector and bias for the output layer. Here, we're transitioning from a hidden layer with multiple neurons to an output layer with potentially one neuron.

$$\mathbf{z}^{[2]} = \mathbf{W}^{[2]} \mathbf{a}^{[1]} + \mathbf{b}^{[2]}$$

- This equation calculates the linear combination of the activated outputs from the hidden layer, weighed by $\mathbf{W}^{[2]}$, and adjusted by the bias $\mathbf{b}^{[2]}$. The result is the input to the output layer's activation function.

$$\mathbf{z}^{[2]} = [24]$$

- The computation is simplified to show the input to the output layer's activation function as 24.

4. Activation in Output Layer:

$$a^{[2]} = \hat{y} = \sigma(\mathbf{z}^{[2]})$$

- The sigmoid function σ is used at the output layer to map the input value into a (0,1) range, which is typical for binary classification tasks or probability estimation.

$$\hat{y} = \sigma([24]) = 0.99$$

- Given the high input value of 24, the sigmoid function outputs a value close to 1, indicating a high confidence level in the positive class, assuming a binary classification context.

9 Backward Propagation

Backward Propagation adjusts the network's parameters by calculating the loss function's gradient and updating weights and biases to minimize prediction errors, optimizing network performance over training iterations.

1. Calculate Loss:

$$L = \sqrt{\frac{1}{m} \sum_{i=1}^m (\hat{y}_i - y_i)^2}$$

$$L = \sqrt{(0.99 - 12)^2} = \sqrt{121.2201}$$

$$L = 11.01$$

- This formula calculates the Root Mean Squared Error (RMSE) between the predicted values (\hat{y}_i) and the actual values (y_i). Here, $\hat{y}_i = 0.99$ and $y_i = 12$.
- The RMSE provides a measure of how well the model is predicting the output, quantifying the difference in terms of the model's accuracy. In this case, an RMSE of 11.01 indicates a significant error, showing that the model's prediction is far from the actual value.

2. Output to Hidden Layer:

Using the chain rule, we first calculate the gradient of the loss function with respect to the output predictions:

$$\frac{\partial L}{\partial \hat{y}} = \frac{2}{m}(\hat{y} - y) \frac{1}{2\sqrt{\text{mean squared error}}}$$

For the next layer's pre-activation output, we derive the gradient with respect to the sigmoid function:

$$\frac{\partial L}{\partial z^{[2]}} = \frac{\partial L}{\partial \hat{y}} \cdot \sigma'(z^{[2]})$$

- $\sigma'(z^{[2]})$ is the derivative of the sigmoid activation function, applied to the pre-activation outputs at the output layer.

3. Hidden Layer to Input Layer:

The gradients of the weights and biases are calculated as follows:

$$\begin{aligned}\frac{\partial L}{\partial W^{[2]}} &= \frac{\partial L}{\partial z^{[2]}} \cdot a^{[1]} \\ \frac{\partial L}{\partial b^{[2]}} &= \frac{\partial L}{\partial z^{[2]}}\end{aligned}$$

For the activations of the previous layer, we use the transpose of the weights:

$$\frac{\partial L}{\partial a^{[1]}} = W^{[2]T} \cdot \frac{\partial L}{\partial z^{[2]}}$$

The derivative through the ReLU function is computed next:

$$\frac{\partial L}{\partial z^{[1]}} = \frac{\partial L}{\partial a^{[1]}} \cdot \text{ReLU}'(z^{[1]})$$

- $\text{ReLU}'(z^{[1]})$ is the derivative of the ReLU function, which is 1 for positive inputs and 0 otherwise.

4. Update Parameters:

Finally, we update the weights and biases using the calculated gradients and a learning rate α :

$$\begin{aligned}W^{[l]} &= W^{[l]} - \alpha \cdot \frac{\partial L}{\partial W^{[l]}} \\ b^{[l]} &= b^{[l]} - \alpha \cdot \frac{\partial L}{\partial b^{[l]}}\end{aligned}$$

- These updates adjust the weights and biases to minimize the loss, thereby improving the model with each iteration.

10 References