

# Understanding Matrix Multiplication in Neural Networks

## A Comprehensive Guide with a Real-World Example

Debasish Maji  
LinkedIn  
GitHub

### Abstract

Matrix multiplication is at the heart of neural network computations. This article provides an in-depth exploration of how matrix operations facilitate the functioning of neural networks, using a real-world example of text classification. Through step-by-step examples and intuitive explanations, we delve into the mathematical foundations and practical implementations of matrix multiplication in neural networks. Fundamental concepts are revisited where necessary to ensure a cohesive and comprehensive learning experience.

## Contents

|          |   |          |
|----------|---|----------|
| <b>1</b> | <b>Introduction</b>   | <b>2</b> |
| <b>2</b> | <b>Fundamentals of Matrices</b>                             | <b>2</b> |
| 2.1      | What is a Matrix? . . . . .                                 | 2        |
| 2.2      | Matrix Operations . . . . .                                 | 3        |
| 2.2.1    | Addition and Subtraction . . . . .                          | 3        |
| 2.2.2    | Scalar Multiplication . . . . .                             | 3        |
| 2.2.3    | Matrix Multiplication . . . . .                             | 3        |
| 2.3      | Vectors . . . . .   | 3        |
| <b>3</b> | <b>Vectors and Matrices in Neural Networks</b>              | <b>3</b> |
| 3.1      | Representation of Data and Parameters . . . . .             | 3        |
| 3.2      | Why Use Matrices? . . . . .                                 | 3        |
| <b>4</b> | <b>Real-World Example: Text Classification</b>              | <b>4</b> |
| 4.1      | Problem Setup . . . . .                                     | 4        |
| 4.2      | Preprocessing: Bag-of-Words Model . . . . .                 | 4        |
| 4.2.1    | Vocabulary . . . . .  | 4        |
| 4.2.2    | Encoding the Document . . . . .                             | 4        |
| 4.3      | Defining the Neural Network . . . . .                       | 4        |
| 4.4      | Initializing Parameters . . . . .                           | 5        |
| 4.5      | Forward Propagation Step-by-Step . . . . .                  | 5        |
| 4.5.1    | Input Vector ( $\mathbf{x}$ ) . . . . .                     | 5        |
| 4.5.2    | Computing Hidden Layer Activations . . . . .                | 5        |
| 4.5.3    | Computing Output Layer Activation . . . . .                 | 6        |
| 4.6      | Summary of Forward Propagation . . . . .                    | 6        |
| <b>5</b> | <b>Understanding Each Step in Depth</b>                     | <b>6</b> |
| 5.1      | Why Use the Bag-of-Words Model? . . . . .                   | 6        |
| 5.2      | Matrix Multiplication in Hidden Layer Computation . . . . . | 6        |
| 5.2.1    | Breaking Down the Computation . . . . .                     | 7        |
| 5.3      | Activation Functions . . . . .                              | 7        |

|          |   |           |
|----------|---|-----------|
| 5.4      | Interpretation of Network Output . . . . .                  | 7         |
| <b>6</b> | <b>Backpropagation and Matrix Multiplication</b>            | <b>7</b>  |
| 6.1      | Computing the Loss . . . . .                                | 7         |
| 6.2      | Backpropagation Step-by-Step . . . . .                      | 7         |
| 6.2.1    | Error at Output Layer . . . . .                             | 7         |
| 6.2.2    | Gradient with Respect to $\mathbf{W}^{[2]}$ . . . . .       | 7         |
| 6.2.3    | Error at Hidden Layer . . . . .                             | 8         |
| 6.2.4    | Gradient with Respect to $\mathbf{W}^{[1]}$ . . . . .       | 8         |
| 6.3      | Updating Parameters . . . . .                               | 8         |
| 6.3.1    | Update $\mathbf{W}^{[2]}$ . . . . .                         | 9         |
| 6.3.2    | Update $\mathbf{b}^{[2]}$ . . . . .                         | 9         |
| 6.3.3    | Update $\mathbf{W}^{[1]}$ . . . . .                         | 9         |
| 6.3.4    | Update $\mathbf{b}^{[1]}$ . . . . .                         | 9         |
| 6.4      | Summary of Backpropagation . . . . .                        | 9         |
| <b>7</b> | <b>Connecting the Concepts</b>                              | <b>9</b>  |
| <b>8</b> | <b>Fundamental Concepts Revisited</b>                       | <b>10</b> |
| 8.1      | Chain Rule in Calculus . . . . .                            | 10        |
| 8.2      | Dot Product and Its Role in Matrix Multiplication . . . . . | 10        |
| 8.3      | Importance of Activation Functions . . . . .                | 10        |
| <b>9</b> | <b>Conclusion</b>   | <b>10</b> |

# 1 Introduction

Neural networks have become a cornerstone of modern artificial intelligence, powering applications ranging from image recognition to natural language processing. At the core of these networks lies the mathematical operation of **matrix multiplication**. Understanding how and why matrices are used in neural networks is crucial for both developing new models and improving existing ones.

This article aims to demystify the role of matrix multiplication in neural networks by providing a logical, step-by-step exploration of the topic. We'll cover everything from the basics of matrices to their application in forward and backward propagation within neural networks. We'll use a real-world example—text classification—to illustrate key concepts and ensure that you gain a thorough understanding of the material. Along the way, we'll revisit fundamental concepts to reinforce your learning.

# 2 Fundamentals of Matrices

Before we delve into neural networks, it's essential to have a solid grasp of matrix mathematics. This section covers the foundational concepts necessary for understanding how matrices function within neural networks.

## 2.1 What is a Matrix?

A **matrix** is a two-dimensional array of numbers arranged in rows and columns. Mathematically, an  $m \times n$  matrix has  $m$  rows and  $n$  columns:

$$\mathbf{A} = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{bmatrix}$$

Matrices are powerful tools for representing and manipulating data in a structured format. They are widely used in various fields such as physics, engineering, computer science, and, of course, machine learning.

## 2.2 Matrix Operations

### 2.2.1 Addition and Subtraction

Matrices of the same dimensions can be added or subtracted element-wise:

$$\mathbf{C} = \mathbf{A} + \mathbf{B}, \quad c_{ij} = a_{ij} + b_{ij}$$

### 2.2.2 Scalar Multiplication

A matrix can be multiplied by a scalar (a single number), multiplying each element:

$$\mathbf{C} = k\mathbf{A}, \quad c_{ij} = k \times a_{ij}$$

### 2.2.3 Matrix Multiplication

Matrix multiplication involves the dot product of rows and columns:

$$\mathbf{C} = \mathbf{AB}, \quad c_{ij} = \sum_{k=1}^n a_{ik}b_{kj}$$

For matrix multiplication to be valid, the number of columns in  $\mathbf{A}$  must equal the number of rows in  $\mathbf{B}$ .

## 2.3 Vectors

A **vector** is a special case of a matrix with either one row or one column:

- **Column Vector:**  $n \times 1$  matrix.
- **Row Vector:**  $1 \times n$  matrix.

Vectors are used to represent quantities that have both magnitude and direction and are essential in representing data inputs and weights in neural networks.

## 3 Vectors and Matrices in Neural Networks

In neural networks, data is often represented as vectors and matrices to facilitate batch processing and efficient computation.

### 3.1 Representation of Data and Parameters

- **Input Features ( $\mathbf{x}$ ):** Represented as a column vector.
- **Weights ( $\mathbf{W}$ ):** Represented as a matrix where each row corresponds to a neuron.
- **Biases ( $\mathbf{b}$ ):** Represented as a column vector.
- **Activations ( $\mathbf{a}$ ):** Output of neurons after applying the activation function.

### 3.2 Why Use Matrices?

Matrices allow us to compute the outputs of all neurons in a layer simultaneously using matrix multiplication. This is both computationally efficient and aligns well with hardware acceleration on GPUs and other specialized processors.

## 4 Real-World Example: Text Classification

To make the concepts more tangible, we'll walk through a real-world example of text classification using a simple neural network. We'll classify a document as either **positive** or **negative** sentiment based on the words it contains.

### 4.1 Problem Setup

Given a document (text), we want to classify its sentiment. We'll use a neural network to perform this classification.

### 4.2 Preprocessing: Bag-of-Words Model

First, we need to convert the text into a numerical format that the neural network can process.

#### 4.2.1 Vocabulary

Assume we have a vocabulary of 5 words:

1. happy
2. sad
3. great
4. terrible
5. okay

#### 4.2.2 Encoding the Document

We use the **Bag-of-Words** model to encode the document into a vector. Each element in the vector represents the count of a word in the document.

**Example Document:**

"I feel happy because today is a great day."

**Encoded Vector (x):**

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$$

The document contains "happy" and "great" once each.

### 4.3 Defining the Neural Network

Our neural network will have:

- **Input Layer:** 5 neurons (one for each word in the vocabulary).
- **Hidden Layer:** 3 neurons.
- **Output Layer:** 1 neuron (outputting a value between 0 and 1, where values closer to 1 indicate positive sentiment).

## 4.4 Initializing Parameters

**Weights for Hidden Layer ( $\mathbf{W}^{[1]}$ ):**

$$\mathbf{W}^{[1]} = \begin{bmatrix} w_{11}^{[1]} & w_{12}^{[1]} & w_{13}^{[1]} & w_{14}^{[1]} & w_{15}^{[1]} \\ w_{21}^{[1]} & w_{22}^{[1]} & w_{23}^{[1]} & w_{24}^{[1]} & w_{25}^{[1]} \\ w_{31}^{[1]} & w_{32}^{[1]} & w_{33}^{[1]} & w_{34}^{[1]} & w_{35}^{[1]} \end{bmatrix}$$

Assume the weights are initialized as:

$$\mathbf{W}^{[1]} = \begin{bmatrix} 0.5 & -0.2 & 0.3 & 0.0 & -0.1 \\ -0.4 & 0.6 & -0.1 & -0.5 & 0.2 \\ 0.1 & -0.3 & 0.7 & -0.2 & 0.4 \end{bmatrix}$$

**Biases for Hidden Layer ( $\mathbf{b}^{[1]}$ ):**

$$\mathbf{b}^{[1]} = \begin{bmatrix} 0.1 \\ -0.1 \\ 0.05 \end{bmatrix}$$

**Weights for Output Layer ( $\mathbf{W}^{[2]}$ ):**

$$\mathbf{W}^{[2]} = \begin{bmatrix} w_{11}^{[2]} & w_{12}^{[2]} & w_{13}^{[2]} \end{bmatrix} = \begin{bmatrix} 0.2 & -0.3 & 0.5 \end{bmatrix}$$

**Bias for Output Layer ( $b^{[2]}$ ):**

$$b^{[2]} = 0.0$$

## 4.5 Forward Propagation Step-by-Step

Let's compute the output of the neural network for our example document.

### 4.5.1 Input Vector ( $\mathbf{x}$ )

$$\mathbf{x} = \begin{bmatrix} 1 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$$

### 4.5.2 Computing Hidden Layer Activations

**Linear Combination ( $\mathbf{z}^{[1]}$ ):**

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

Compute  $\mathbf{W}^{[1]}\mathbf{x}$ :

$$\mathbf{W}^{[1]}\mathbf{x} = \begin{bmatrix} (0.5)(1) + (-0.2)(0) + (0.3)(1) + (0.0)(0) + (-0.1)(0) \\ (-0.4)(1) + (0.6)(0) + (-0.1)(1) + (-0.5)(0) + (0.2)(0) \\ (0.1)(1) + (-0.3)(0) + (0.7)(1) + (-0.2)(0) + (0.4)(0) \end{bmatrix}$$

Simplify:

$$\mathbf{W}^{[1]}\mathbf{x} = \begin{bmatrix} 0.5 + 0 + 0.3 + 0 + 0 \\ -0.4 + 0 - 0.1 + 0 + 0 \\ 0.1 + 0 + 0.7 + 0 + 0 \end{bmatrix} = \begin{bmatrix} 0.8 \\ -0.5 \\ 0.8 \end{bmatrix}$$

Add biases:

$$\mathbf{z}^{[1]} = \begin{bmatrix} 0.8 + 0.1 \\ -0.5 - 0.1 \\ 0.8 + 0.05 \end{bmatrix} = \begin{bmatrix} 0.9 \\ -0.6 \\ 0.85 \end{bmatrix}$$

#### Applying Activation Function

Using the ReLU activation function:

$$f(z) = \max(0, z)$$

Compute  $\mathbf{a}^{[1]}$ :

$$\mathbf{a}^{[1]} = \begin{bmatrix} \max(0, 0.9) \\ \max(0, -0.6) \\ \max(0, 0.85) \end{bmatrix} = \begin{bmatrix} 0.9 \\ 0 \\ 0.85 \end{bmatrix}$$

### 4.5.3 Computing Output Layer Activation

**Linear Combination** ( $z^{[2]}$ ):

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

Compute:

$$z^{[2]} = (0.2)(0.9) + (-0.3)(0) + (0.5)(0.85) + 0 = 0.18 + 0 + 0.425 = 0.605$$

#### Applying Activation Function

Using the sigmoid activation function:

$$\hat{y} = f(z^{[2]}) = \frac{1}{1 + e^{-z^{[2]}}}$$

Compute:

$$\hat{y} = \frac{1}{1 + e^{-0.605}} \approx 0.646$$

#### Interpretation

A predicted output of approximately 0.646 suggests that the document has a positive sentiment.

## 4.6 Summary of Forward Propagation

We have processed the input document through the neural network using matrix multiplication at each layer to compute the activations, resulting in a sentiment prediction.

## 5 Understanding Each Step in Depth

To ensure a comprehensive understanding, let's delve deeper into each step, revising relevant concepts along the way.

### 5.1 Why Use the Bag-of-Words Model?

The Bag-of-Words model simplifies text data by representing it as numerical vectors, making it suitable for neural network input. It disregards word order but captures word frequency.

### 5.2 Matrix Multiplication in Hidden Layer Computation

The operation  $\mathbf{z}^{[1]} = \mathbf{W}^{[1]}\mathbf{x} + \mathbf{b}^{[1]}$  is a matrix-vector multiplication followed by vector addition. Each neuron's pre-activation value is a weighted sum of the input features plus a bias.

### 5.2.1 Breaking Down the Computation

For the first neuron in the hidden layer:

$$z_1^{[1]} = \sum_{i=1}^5 w_{1i}^{[1]} x_i + b_1^{[1]} = w_{11}^{[1]} x_1 + w_{12}^{[1]} x_2 + \dots + w_{15}^{[1]} x_5 + b_1^{[1]}$$

## 5.3 Activation Functions

Activation functions introduce non-linearity, allowing the network to learn complex patterns.

- **ReLU:**  $f(z) = \max(0, z)$
- **Sigmoid:**  $f(z) = \frac{1}{1+e^{-z}}$

## 5.4 Interpretation of Network Output

The output of the network is a probability between 0 and 1 due to the sigmoid activation, making it suitable for binary classification tasks like sentiment analysis.

# 6 Backpropagation and Matrix Multiplication

To train the network, we need to adjust the weights and biases based on the error between the predicted output and the true label.

## 6.1 Computing the Loss

Assuming the true label  $y = 1$  (positive sentiment), we use the binary cross-entropy loss function:

$$L = -[y \log(\hat{y}) + (1 - y) \log(1 - \hat{y})]$$

Compute:

$$L = -[1 \times \log(0.646) + 0 \times \log(1 - 0.646)] \approx 0.437$$

## 6.2 Backpropagation Step-by-Step

### 6.2.1 Error at Output Layer

$$\delta^{[2]} = \hat{y} - y = 0.646 - 1 = -0.354$$

### 6.2.2 Gradient with Respect to $\mathbf{W}^{[2]}$

$$\frac{\partial L}{\partial \mathbf{W}^{[2]}} = \delta^{[2]} (\mathbf{a}^{[1]})^\top$$

Compute:

$$\frac{\partial L}{\partial \mathbf{W}^{[2]}} = -0.354 \begin{bmatrix} 0.9 & 0 & 0.85 \end{bmatrix}$$

$$\frac{\partial L}{\partial \mathbf{W}^{[2]}} = \begin{bmatrix} -0.3186 & 0 & -0.3009 \end{bmatrix}$$

### 6.2.3 Error at Hidden Layer

Compute the derivative of the ReLU activation function:

$$f'(z) = \begin{cases} 1 & \text{if } z > 0 \\ 0 & \text{if } z \leq 0 \end{cases}$$

Since  $\mathbf{z}^{[1]} = \begin{bmatrix} 0.9 \\ -0.6 \\ 0.85 \end{bmatrix}$ , the derivative vector is:

$$f'(\mathbf{z}^{[1]}) = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}$$

Compute:

$$\delta^{[1]} = (\mathbf{W}^{[2]})^\top \delta^{[2]} \odot f'(\mathbf{z}^{[1]})$$

Transpose  $\mathbf{W}^{[2]}$ :

$$(\mathbf{W}^{[2]})^\top = \begin{bmatrix} 0.2 \\ -0.3 \\ 0.5 \end{bmatrix}$$

Compute:

$$(\mathbf{W}^{[2]})^\top \delta^{[2]} = -0.354 \begin{bmatrix} 0.2 \\ -0.3 \\ 0.5 \end{bmatrix} = \begin{bmatrix} -0.0708 \\ 0.1062 \\ -0.177 \end{bmatrix}$$

Apply element-wise multiplication with  $f'(\mathbf{z}^{[1]})$ :

$$\delta^{[1]} = \begin{bmatrix} -0.0708 \times 1 \\ 0.1062 \times 0 \\ -0.177 \times 1 \end{bmatrix} = \begin{bmatrix} -0.0708 \\ 0 \\ -0.177 \end{bmatrix}$$

### 6.2.4 Gradient with Respect to $\mathbf{W}^{[1]}$

$$\frac{\partial L}{\partial \mathbf{W}^{[1]}} = \delta^{[1]}(\mathbf{x})^\top$$

Compute:

$$\delta^{[1]} = \begin{bmatrix} -0.0708 \\ 0 \\ -0.177 \end{bmatrix}$$

$$(\mathbf{x})^\top = [1 \quad 0 \quad 1 \quad 0 \quad 0]$$

Compute the gradient matrix:

$$\frac{\partial L}{\partial \mathbf{W}^{[1]}} = \begin{bmatrix} -0.0708 \\ 0 \\ -0.177 \end{bmatrix} [1 \quad 0 \quad 1 \quad 0 \quad 0] = \begin{bmatrix} -0.0708 & 0 & -0.0708 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ -0.177 & 0 & -0.177 & 0 & 0 \end{bmatrix}$$

## 6.3 Updating Parameters

Using a learning rate  $\alpha = 0.1$ :



### 6.3.1 Update $\mathbf{W}^{[2]}$

$$\mathbf{W}_{\text{new}}^{[2]} = \mathbf{W}^{[2]} - \alpha \frac{\partial L}{\partial \mathbf{W}^{[2]}}$$

Compute:

$$\mathbf{W}_{\text{new}}^{[2]} = \begin{bmatrix} 0.2 \\ -0.3 \\ 0.5 \end{bmatrix} - 0.1 \begin{bmatrix} -0.3186 \\ 0 \\ -0.3009 \end{bmatrix} = \begin{bmatrix} 0.2 + 0.03186 \\ -0.3 + 0 \\ 0.5 + 0.03009 \end{bmatrix} = \begin{bmatrix} 0.23186 \\ -0.3 \\ 0.53009 \end{bmatrix}$$

### 6.3.2 Update $\mathbf{b}^{[2]}$

$$b_{\text{new}}^{[2]} = b^{[2]} - \alpha \delta^{[2]} = 0 - 0.1(-0.354) = 0.0354$$

### 6.3.3 Update $\mathbf{W}^{[1]}$

$$\mathbf{W}_{\text{new}}^{[1]} = \mathbf{W}^{[1]} - \alpha \frac{\partial L}{\partial \mathbf{W}^{[1]}}$$

Compute updated weights (only non-zero elements change):

$$\mathbf{W}_{\text{new}}^{[1]} = \mathbf{W}^{[1]} - 0.1 \times \begin{bmatrix} -0.0708 & 0 & -0.0708 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ -0.177 & 0 & -0.177 & 0 & 0 \end{bmatrix}$$

Update first row:

$$\mathbf{W}_{\text{new},1}^{[1]} = [0.5 - (-0.00708) \quad -0.2 \quad 0.3 - (-0.00708) \quad 0.0 \quad -0.1] = [0.50708 \quad -0.2 \quad 0.30708 \quad 0.0 \quad -0.1]$$

Update third row:

$$\mathbf{W}_{\text{new},3}^{[1]} = [0.1 - (-0.0177) \quad -0.3 \quad 0.7 - (-0.0177) \quad -0.2 \quad 0.4] = [0.1177 \quad -0.3 \quad 0.7177 \quad -0.2 \quad 0.4]$$

### 6.3.4 Update $\mathbf{b}^{[1]}$

$$\mathbf{b}_{\text{new}}^{[1]} = \mathbf{b}^{[1]} - \alpha \delta^{[1]}$$

Compute:

$$\mathbf{b}_{\text{new}}^{[1]} = \begin{bmatrix} 0.1 - 0.1(-0.0708) \\ -0.1 - 0.1(0) \\ 0.05 - 0.1(-0.177) \end{bmatrix} = \begin{bmatrix} 0.1 + 0.00708 \\ -0.1 \\ 0.05 + 0.0177 \end{bmatrix} = \begin{bmatrix} 0.10708 \\ -0.1 \\ 0.0677 \end{bmatrix}$$

## 6.4 Summary of Backpropagation

Through matrix multiplication and application of the chain rule, we've computed the gradients and updated the network's parameters to minimize the loss.

## 7 Connecting the Concepts

By walking through a real-world example, we've connected the abstract mathematical concepts to practical neural network operations. Each step, from input encoding to parameter updates, involves matrix multiplication, highlighting its central role in neural networks.

## 8 Fundamental Concepts Revisited

### 8.1 Chain Rule in Calculus

The chain rule is essential in backpropagation for computing derivatives of composite functions. In matrix form, it allows us to efficiently compute gradients across layers.

### 8.2 Dot Product and Its Role in Matrix Multiplication

The dot product calculates the sum of the products of corresponding entries in two sequences of numbers. In matrix multiplication, it computes the weighted sum of inputs for each neuron.

### 8.3 Importance of Activation Functions

Activation functions introduce non-linearity, enabling neural networks to learn complex patterns. Understanding their properties and derivatives is crucial for both forward and backward propagation.

## 9 Conclusion

Matrix multiplication is integral to the functioning of neural networks, enabling efficient computation of forward and backward passes. By representing inputs, weights, and activations as matrices, we can leverage optimized mathematical operations to build scalable and high-performing neural networks.

Through a detailed, real-world example, we've explored how these mathematical concepts are applied in practice, reinforcing the learning experience. Understanding the role of matrices not only enhances your theoretical knowledge but also improves practical implementation skills.

## References

- Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep Learning*. MIT Press. Available Online
- Nielsen, M. (2015). *Neural Networks and Deep Learning*. Available Online
- NumPy Documentation
- PyTorch Documentation
- TensorFlow Documentation

## About the Author

**Debasish Maji** is a deep learning researcher and enthusiast with a passion for unraveling the mathematical foundations of neural networks. With expertise in machine learning and data science, Debasish is dedicated to sharing knowledge and advancing the field.

- **LinkedIn:** <https://www.linkedin.com/in/debasish-maji-88170a96/>
- **GitHub:** <https://github.com/DebasishMaji>