

# The Role of Linear Algebra in Deep Learning

Thienkim Nguyen

April 2025

## 1 Introduction

Deep learning has revolutionized fields such as computer vision, natural language processing, and autonomous systems by enabling models to learn intricate patterns from vast datasets. At the heart of these deep learning models lies linear algebra—a mathematical framework that provides the tools to represent and manipulate data efficiently. Understanding the role of linear algebra in deep learning is crucial for comprehending how these models operate and for developing more effective algorithms. This paper explores the foundational linear algebra concepts that underpin deep learning architectures, highlighting their applications and significance in the design and training of neural networks.

## 2 Vectors and Matrices in Neural Networks

Neural networks fundamentally rely on linear transformations of data, represented through vectors and matrices. Each layer of a neural network transforms its inputs via matrix multiplication followed by an activation function. Formally, the transformation at a given layer can be expressed as:

where  $X$  is the input vector,  $W$  is the weight matrix,  $b$  is the bias vector, and  $Z$  is the resulting output. These simple matrix operations enable deep networks to perform complex tasks such as classification and regression by chaining multiple such transformations together.

In a typical feedforward neural network, the input data is a vector of features, which is transformed as it passes through successive layers of the network. At each layer, the dot product between the input vector and the weight matrix is computed, and a bias is added. This operation can be interpreted as a projection or rotation in a high-dimensional space, depending on the values of the weights. The non-linear activation function that follows introduces the ability to model complex, non-linear relationships.

To illustrate, consider a layer with input vector  $X \in \mathbb{R}^3$ , weight matrix  $W \in \mathbb{R}^{4 \times 3}$ , and bias vector  $b \in \mathbb{R}^4$ . The output  $Z \in \mathbb{R}^4$  is computed as:

$$\begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{23} & w_{24} \\ w_{31} & w_{33} & w_{34} \\ w_{41} & w_{43} & w_{44} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \\ b_3 \\ b_4 \end{bmatrix}$$

This results in a 4-dimensional vector  $Z$  that becomes the input to the next layer.

Moreover, in batch processing (where multiple data points are processed simultaneously), input data is represented as a matrix  $X \in \mathbb{R}^{n \times d}$ , where  $n$  is the number of samples and  $d$  is the number of features. The weight matrix  $W$  is transposed and applied to the batch as:

This matrix multiplication allows for parallelization, which is crucial for computational efficiency in large-scale deep learning models.

**Definition 1.** A **vector** is an ordered collection of numbers arranged either in a row or a column. In  $\mathbb{R}^n$ , a vector has  $n$  real-valued components and is commonly represented as a column matrix.

good intro, that algebra seems off, maybe that should be algebra, a?

this is an advanced level course, but perhaps defs of vector, matrix, and linear transformation would help here rather than later. Also proving something is a linear transformation might be a good step

love these explanations. good work.

**Definition 2.** A **matrix** is a rectangular array of numbers arranged in rows and columns. An  $m \times n$  matrix has  $m$  rows and  $n$  columns and is used to represent linear transformations between vector spaces.

### 3 Automatic Differentiation and Linear Algebra

Modern deep learning frameworks heavily depend on automatic differentiation to compute gradients efficiently. According to Seeger et al. (2017), automatic differentiation leverages linear algebra operations like the Jacobian and Hessian matrices to propagate gradients through complex networks [?]. These operations allow for scalable optimization, which is vital when training large models with millions of parameters. The key idea is that derivatives of composite functions can be decomposed using the chain rule, with matrix multiplication serving as the primary computational mechanism.

**Theorem 1.** Let  $f(x) = g(h(x))$  be a composition of two differentiable functions. Then, by the chain rule: If  $h : \mathbb{R}^n \rightarrow \mathbb{R}^m$  and  $g : \mathbb{R}^m \rightarrow \mathbb{R}$ , then the Jacobian of  $f$  is: where  $J_g$  and  $J_h$  are Jacobian matrices. This matrix product forms the basis of backpropagation.

### 4 Tensor Computations in PyTorch and TensorFlow

Applying Theorem 1 to Definitions 1 and 2, we understand how gradient descent relies on vector and matrix operations at its core. Deep learning libraries like TensorFlow and PyTorch utilize tensors—generalizations of matrices to higher dimensions—to represent data and model parameters. As Sankaran et al. (2022) demonstrate, the linear algebra awareness of these frameworks is critical to performance [?]. Both frameworks implement operations such as batched matrix multiplications, element-wise activation functions, and broadcasting rules that facilitate efficient training of deep networks. They also exploit sparse representations and GPU acceleration to speed up linear algebra operations, highlighting the importance of numerical efficiency.

### 5 Principal Component Analysis and Dimensionality Reduction

Principal Component Analysis (PCA) is a technique used in data preprocessing to reduce dimensionality while preserving the variance of the data. PCA relies heavily on the spectral decomposition of the covariance matrix and involves eigenvalues and eigenvectors—core concepts from linear algebra.

Given a zero-centered data matrix  $X \in \mathbb{R}^{n \times d}$ , the covariance matrix is:

The eigenvectors of  $C$  point in the directions of maximum variance, and the corresponding eigenvalues indicate the magnitude of variance in each direction. PCA selects the top  $k$  eigenvectors with the highest eigenvalues and projects the original data onto the subspace they span:

where  $W_k$  is a matrix containing the top  $k$  eigenvectors as columns.

This linear transformation reduces computational complexity and is useful in training more efficient neural networks, particularly for image data, text embeddings, and high-dimensional sensor readings.

### 6 The Moore-Penrose Pseudoinverse in Learning

In many learning tasks, especially when the system of equations  $Ax = b$  does not have an exact solution, we seek a solution in the least squares sense. This is done using the Moore-Penrose pseudoinverse  $A^+$ . For a matrix  $A \in \mathbb{R}^{m \times n}$ , the pseudoinverse is defined such that:

This allows us to compute:

as the best approximation to a solution. In machine learning, this technique underpins closed-form solutions for linear regression and is used for initializing parameters in deep models.

gradients were not discussed previously, so new term.

can this be a citation?

wow, these summaries are amazingly clear.

maybe a proof here? backprop is also a little poorly contextualized

thanks I finally kind of understand PCA

## 7 Solving Inverse Problems with Deep Learning

Linear algebra also plays a central role in solving inverse problems, where the goal is to recover an input from an output. Bai et al. (2020) explore how deep learning methods can be adapted to solve linear inverse problems common in imaging and signal processing [?]. These methods often involve learning an approximate inverse of a linear transformation, leveraging neural networks trained via gradient descent. Here, understanding the condition number of matrices and the stability of solutions becomes crucial, again drawing on core linear algebra concepts.

*Proof.* Let  $A \in \mathbb{R}^{n \times n}$  be an invertible matrix. Then the inverse  $A^{-1}$  satisfies:

where  $I$  is the identity matrix. If we assume  $A$  is ill-conditioned, small perturbations  $\Delta b$  in the solution  $Ax = b$  can result in large changes in  $x$ . This instability is quantified by the condition number  $\kappa(A) = |A| \cdot |A^{-1}|$ , which affects the convergence of learning-based solutions.  $\square$

## Conclusion

Linear algebra is far more than a mathematical tool in deep learning—it is the very fabric of how models are defined, trained, and optimized. From fundamental operations like matrix multiplication to advanced topics such as tensor decompositions and automatic differentiation, linear algebra enables the scalability and expressiveness of modern neural networks. As deep learning continues to evolve, so too will the sophistication of its mathematical underpinnings, with linear algebra remaining at the core of this transformation.

over all well written, but as I got further in I lost track of what a lot of this means. Also where are the proofs?