# Toward Intuitive and Accessible Machine Learning Education: A Structured Pythonic Pseudocode Approach

Hadrian Lazic

April, 28, 2025

**Abstract**

The traditional reliance on compressed mathematical notation in machine learning, particularly in calculus-intensive domains such as optimization, presents significant cognitive barriers for both practitioners and students. Core methodologies, such as the Adam optimizer, are widely used in applied settings but are often introduced through dense symbolic expressions that obscure foundational intuitions and hinder practical comprehension.

This work proposes a rearticulation of foundational machine learning concepts—including derivatives, gradients, and optimization algorithms, using an intuitive, Pythonic pseudocode paradigm, supported by annotated visual exemplars. By replacing abstract mathematical formalism with code-oriented, formula, inspired explanations, the proposed framework enhances conceptual transparency, operational clarity, and pedagogical accessibility.

The overarching goal is to empower developers—particularly those without formal training in advanced mathematics—to internalize, implement, and extend key machine learning constructs with confidence and rigor. In democratizing theoretical understanding, this work seeks to broaden participation in machine learning research and development, fostering a more diverse and interdisciplinary technical community.

## 1 Introduction

Mathematics, at its core, constitutes a formal language, an abstract system developed to represent, generalize, and communicate universal patterns and phenomena. Within the field of machine learning, this language is predominantly instantiated through the framework of calculus. Foundational machine learning algorithms are frequently expressed as calculus-based formulations, reflecting a mathematical tradition originating in the 17th century to model continuous and dynamic systems.

In contemporary practice, however, machine learning is primarily instantiated through programming. This evolution has introduced a pronounced disjunction: while many practitioners possess substantial expertise in software engineering and systems design, comparatively few have received formal instruction in the calculus underpinning the algorithms they implement. As a result, the mathematical logic foundational to machine learning often remains opaque, inhibiting developers, even those with advanced technical proficiency, from meaningfully engaging with, modifying, or innovating upon algorithmic structures.

More broadly, this disconnect exacerbates systemic inequities in technical education. A substantial cohort of students, particularly those from under, resourced educational environments, may develop exceptional

programming skills yet lack exposure to advanced mathematics. Consequently, access to the deeper strata of machine learning research and development remains uneven, constraining individual opportunity and diminishing the diversity of perspectives essential to innovation in the field.

While calculus remains indispensable for theoretical advancements in machine learning, there is a critical need to construct more inclusive and operationally transparent bridges between mathematical theory and computational practice. To this end, this paper introduces *Formulaic Pythonic Pseudocode*: a structured notation system designed to translate complex calculus, based formulations into accessible, code-oriented representations. This approach preserves mathematical rigor while mitigating the cognitive overhead associated with simultaneously navigating distinct formal systems.

Importantly, the proposed framework is grounded not in ephemeral syntactic conventions, which vary across programming languages, but rather in enduring computational paradigms, such as iterative loops, that mirror fundamental mathematical constructs, including the summation operator ($\Sigma$). As such, Formulaic Pythonic Pseudocode constitutes a robust, durable bridge between mathematical abstraction and practical implementation, rendering machine learning theory more transparent, intuitive, and accessible to a broader and more diverse community of developers.

## 2    Motivation and Background

The development of a machine learning framework in the Rust programming language constituted a significant undertaking in my trajectory as a student researcher. This project necessitated not only the construction of computationally efficient algorithms but also sustained engagement with the advanced mathematical formalism underpinning contemporary machine learning methodologies. The process of translating theoretical formulations into executable systems revealed a series of critical challenges, particularly in navigating the dense and abstract notation prevalent throughout the existing literature.

A focal point of this experience was the implementation of foundational optimization algorithms, such as Adam. The canonical formulations of such algorithms often employ compact symbolic representations that, while mathematically elegant, present formidable barriers to interpretation, even among students with considerable mathematical training. In particular, the standard expression of the Adam optimizer exemplifies this opacity, with many constituent symbols, assumptions, and implied computational steps remaining inaccessible without extensive auxiliary investigation. Through a methodical process of independent research, critical analysis, and iterative refinement, I ultimately derived a fully operational implementation of Adam in Rust, comprising approximately thirty-five lines of source code. The choice of Rust was deliberate, given its capacity to deliver performance comparable to C-based environments while preserving memory safety guarantees and offering a robust ecosystem for build and dependency management.

This experience illuminated a broader structural phenomenon: the persistent disconnect between theoretical macdetailshine learning research and practical software engineering. Although pseudocode is frequently included in academic publications, such representations often approximate condensed calculus more than they resemble executable algorithmic specifications. Essential details, explicit control flow structures, data handling strategies, and operational semantics, are frequently omitted or under-specified, rendering direct translation into code nontrivial and error-prone.

Motivated by these observations, this work seeks to contribute toward narrowing the gap between theoretical exposition and practical implementation. By offering rigorously validated yet accessible algorithmic translations, the project aspires to facilitate deeper engagement with machine learning methodologies among

students, software developers, and early-career researchers. It is hoped that these contributions will foster a more integrative understanding of machine learning as both a mathematical discipline and an engineering practice.

*The overarching goal is to construct a more inclusive and operationally transparent pathway into the field, beginning with my own contributions.*

# 3   Core Principles

In the development of this work, several core principles were adopted to ensure the accessibility, transparency, and pedagogical rigor of complex Machine Learning concepts. These guiding methodologies are intended to promote a deeper and more intuitive understanding before engaging with formal mathematical expressions.

- **Precise Definitions of Foundational Terms:** All critical terminology (e.g., *weights*, *gradients*) is clearly and systematically defined to establish a consistent conceptual framework.

- **Graphical and Diagrammatic Representation:** Visual aids, including graphs and diagrams, are utilized extensively to illustrate data flow, transformations, and algorithmic processes, thereby enhancing conceptual clarity.

- **Formulaic Pseudocode as a Bridging Mechanism:** Pseudocode formulations are employed to systematically bridge abstract mathematical expressions and practical programming logic, facilitating a more intuitive transition from theory to implementation.

- **Prioritization of Intuitive Explanations:** Initial emphasis is placed on cultivating intuitive understanding before introducing formal mathematical structures, with the goal of reducing cognitive barriers and promoting more effective knowledge transfer.

These principles collectively inform the structure of the paper and reflect a commitment to making Machine Learning education more accessible, particularly for students and practitioners without an extensive background in advanced mathematics.

# 4   Target Audience

This work is intended for a diverse audience committed to advancing their understanding of Machine Learning, regardless of prior formal exposure to advanced mathematics. The target readership includes:

- High school students embarking on the study of Machine Learning, particularly those seeking accessible and intuitive learning resources.

- Coding bootcamp graduates aiming to deepen their conceptual and practical grasp of Machine Learning beyond standard curricula.

- Software developers and practitioners who possess limited exposure to calculus or higher-level mathematics, yet aspire to implement and innovate within Machine Learning systems.

- Individuals who have encountered challenges with traditional Machine Learning literature, often characterized by excessive formalism and limited pedagogical accessibility.

- Current and future educators interested in re-envisioning Machine Learning instruction models, with an emphasis on enhancing accessibility, transparency, and engagement for emerging generations of developers and researchers.

By explicitly addressing the needs of these groups, this work seeks to bridge existing educational gaps and promote a more inclusive and effective approach to Machine Learning education.

It is **not** intended for those already fluent in math-heavy Machine Learning theory, rather, it fills the gap for coders who want to go deeper without waiting for university-level math.

# 5    Vision for the Future of Machine Learning Education

The traditional paradigm for teaching machine learning has historically privileged abstract mathematical formalism, often at the expense of accessibility and inclusivity. While mathematical rigor remains indispensable for the theoretical advancement of the field, an overreliance on abstraction risks alienating a substantial cohort of capable students, practitioners, and developers who might otherwise contribute meaningfully to the discipline.

This work advocates for an alternative educational framework: prioritizing intuition-driven, structured pseudocode alongside annotated visual representations to convey complex machine learning concepts. Such an approach seeks to preserve theoretical depth while significantly enhancing accessibility and conceptual clarity.

I envision a future in which educators, curriculum designers, and self-directed learners adopt intuition-first pedagogical models that lower barriers to entry without sacrificing mathematical integrity. By explicitly bridging the domains of computational logic and mathematical abstraction, it becomes possible to cultivate educational environments wherein intuitive understanding and formal reasoning reinforce one another rather than stand in opposition.

This paper aspires to serve both as a practical resource for learners and as a call to action for educators committed to reimagining the future of machine learning education. A shift toward more inclusive, intuitively structured pedagogy promises not only to democratize access to machine learning expertise but also to accelerate innovation and diversify contributions across the field.

# 6    Contribution Summary

This paper presents a novel approach to enhancing machine learning accessibility by systematically translating complex, calculus-based concepts into code-oriented pseudocode, designed to mirror real-world programming logic. Central to this contribution is the introduction of *Formulaic Pythonic Pseudocode*, a notation system specifically engineered to render traditional mathematical formulations accessible to developers without advanced calculus backgrounds by leveraging conventions recognizable across modern programming languages.

Unlike prevailing practices, where pseudocode often mirrors mathematical formalism rather than executable logic, this approach emphasizes programmer-readable structures, enabling comprehension and practical implementation. By facilitating earlier engagement among high school students, coding bootcamp graduates, and early-career developers, the framework aims to expand participation in machine learning model development and foster a more inclusive talent pipeline.

Additionally, this work seeks to establish a new pedagogical standard for machine learning education resources. The proposed framework systematically organizes core machine learning concepts, providing simplified linguistic descriptions, visual representations (such as graphs), detailed breakdowns of constants and hyperparameters, and corresponding Formulaic Pseudocode exemplars. The practical viability of the framework has been demonstrated through its successful application within an independently developed Rust-based machine learning library.

Through these contributions, this paper endeavors to bridge the gap between theoretical exposition and practical engineering, thereby broadening the accessibility, transparency, and operational clarity of machine learning education.

# Contents

# 7   Fundations to Machine Learning

Machine Learning is a subfield of artificial intelligence (AI) that focuses on enabling computational systems to automatically identify patterns and structures within large datasets. The core principle of machine learning involves the development of algorithms that allow machines to learn from data, adapt to new information, and make predictions or decisions without being explicitly programmed for every task. In this context, the patterns discovered by the machine are considered a form of intelligence, wherein the system demonstrates the ability to infer relationships, make generalizations, and apply learned knowledge to novel situations.

Machine learning models, often referred to as networks, are typically structured as complex systems of interconnected components. These components, such as neurons and layers, are arranged in various architectures, depending on the specific learning task and algorithm used. The learning process involves a sequence of computational steps, where the system iteratively adjusts internal parameters based on the input data to minimize errors in its predictions or classifications. This process, commonly referred to as training, is essential for the model's ability to generalize from the training data to unseen data, ensuring its applicability in real-world scenarios.

- **Gradient:** Refers to a tensor or scalar value that represents the rate of change of a function with respect to its inputs, typically used in optimization algorithms to adjust parameters during training.

- **Scalar:** A scalar is a tensor with a single parameter or value, often representing a simple quantity such as a real number or a single element within a dataset.

- **Vector:** A vector is a one-dimensional tensor, typically represented as an array or list of numbers, each corresponding to a distinct quantity within a multi-dimensional space.

- **Tensor:** A tensor is a multi-dimensional array that generalizes the concept of scalars (zero-dimensional tensors) and vectors (one-dimensional tensors) to higher dimensions, used to represent and store multi-dimensional data structures.

- **Target:** The target refers to the desired gradient or output value that the model aims to approximate or predict, typically specified in the output layer of a neural network during the training process.

- **Output:** The output refers to the gradient or value that is produced by the network's final layer, commonly referred to as the output layer. This value represents the model's prediction after processing the input through all preceding layers

- **Weights, and Bias:** Weights and biases are the parameters of each layer within a neural network. Weights represent the strength of the connections between neurons, while biases provide additional offsets to the activation functions. Together, these parameters define the transformation applied to input data as it propagates through the network. Analogous to synapses in the human brain, they control the flow of information within the model.

- **Underfitting and Overfitting:** Underfitting occurs when a model fails to capture the underlying patterns within the training data, often due to insufficient complexity or inadequate training. Conversely, overfitting arises when a model excessively memorizes the training data, capturing noise and irrelevant details instead of generalizable patterns. This often results in poor performance on unseen data.

Figure 1: Neural networks follow the standard Machine Learning pipeline:starting with data collection and preprocessing, followed by model training, validation, and testing, in order to iteratively learn patterns and improve performance over time.

- *Further sections of this paper provide a more in-depth exploration of these concepts and their implications in the context of neural network training and optimization.*

## 7.1 What can a Tensor store

Tensors can be conceptualized as multi-dimensional arrays that generalize matrices to higher dimensions. Specifically, in the case of two-dimensional tensors, they can be represented as matrices defined by an X and Y axis, effectively forming a grid structure composed of values. In the context of machine learning, it is common to utilize 32-bit floating point representations (float32 or f32) for tensor storage. In certain cases, even lower-precision formats such as half-precision (e.g., float16) may be employed, as the increased precision provided by 64-bit floating point formats (float64 or f64) is often deemed unnecessary for the majority of tasks. The adoption of lower-precision formats offers a significant advantage in reducing memory consumption and computational overhead, while maintaining the performance of the model in most practical applications. This trade-off between precision and efficiency is a critical consideration in the design and implementation of machine learning models, especially when scaling to large datasets or complex architectures.

## 7.2 Layer in a Machine Learning Context

A neural network is inherently composed of multiple layers, each playing a critical role in the transformation and processing of input data. Within each layer, there exist two fundamental components: the weight tensor and the bias tensor. The weight tensor governs the scaling of input signals, determining the strength and direction of the connections between neurons. Conversely, the bias tensor facilitates the network's ability to model complex, non-linear relationships, enabling the network to capture patterns that cannot be represented through weighted inputs alone. Together, these components allow the network to approximate highly intricate functions, making it capable of learning from data in a manner that transcends simple linear transformations.

Figure 2: Basic structure of a fully connected neural network, comprising the Input, Hidden, and Output layers. Each layer is represented as a set of neurons, with each neuron in one layer being connected to every neuron in the subsequent layer.

## 7.3    Neural Network

### 7.3.1    Understanding the Neural Network

**Notice how each weight/neuron within the network is connected to nearby neurons in the adjacent top and bottom layers.**

1. **Input:** Accepts raw data, which can be normalized to prevent gradient issues.

2. **Hidden:** Applies matrix multiplication to propagate signal strengths throughout the network.

3. **Output:** Produces the final prediction, which can be compared to the target to compute the network error.

**Training Steps Overview:** The training of neural networks involves a systematic procedure that ensures the model learns effectively from the data. The steps outlined below describe the critical stages in this process:

1. **Data Collection and Preprocessing:** A fundamental prerequisite for training a neural network is the acquisition of a clean, well-prepared dataset. This involves identifying and removing significant outliers and ensuring that the input features are correctly paired with their corresponding target values. The integrity of the data must be preserved to avoid introducing noise that could impair model performance.

2. **Model Architecture Definition:** The next step is to design the architecture of the neural network. This includes the selection of layer types (such as Dense, Convolutional, or Recurrent layers), determining the number of layers, and specifying the number of units within each layer. Additionally, appropriate activation functions must be chosen for each layer to enable the model to learn complex relationships within the data.

3. **Loss Function Selection:** The choice of a loss function is crucial for quantifying the error between the predicted outputs and the actual target values. Common loss functions include Mean Squared Error

(MSE) for regression tasks and Mean Absolute Error (MAE), among others. The loss function serves as the objective that the network aims to minimize during training.

4. **Optimizer Selection:** To optimize the learning process, an optimizer must be chosen to adjust the network's weights. Algorithms such as Adam are commonly used due to their adaptive learning rates and effective convergence properties. The optimizer adjusts the learning rate dynamically to ensure efficient convergence based on the gradients computed during backpropagation.

5. **Foward Pass:** During the forward pass, input data is propagated through the network. Each layer computes an activation function, transforming the input data at each stage. The final layer produces the predicted output, which is compared to the true labels or target values.

6. **Loss Computation:** Once the forward pass is complete, the predicted output is evaluated against the true target values using the selected loss function. This step quantifies the difference between the predicted and actual values, providing a measure of the model's performance.

7. **Backward Pass (Backpropagation):** The backward pass involves calculating the gradients of the loss function with respect to the network's weights. This is achieved by applying the Chain Rule of calculus to propagate the error backward through the network, layer by layer, to compute the necessary gradients.

8. **Weight Update:** The optimizer is then used to update the network's weights and biases based on the gradients computed during the backward pass. This iterative process continues until the model converges to a solution that minimizes the loss function.

This work provides dedicated sections that systematically examine each of these processes in greater detail. The *Forward Pass*, *Loss Computation*, and *Backward Pass* are presented as essential phases within the broader Backpropagation framework, which constitutes the core mechanism enabling neural networks to iteratively refine their parameters and thereby learn from data.

# 8   Tensors

In the context of Machine Learning, tensors serve as a foundational mathematical structure for representing and manipulating data across multiple dimensions. Although tensors are frequently utilized to encode gradients the derivatives of functions with respect to their parameters their utility extends far beyond this application. More generally, tensors provide a unified and systematic framework for modeling multi-dimensional data arrays, irrespective of dimensional complexity.

This intrinsic versatility enables tensors to facilitate the efficient storage, transformation, and computation of complex datasets within modern computational models. As such, tensors are indispensable not only for the theoretical formulation of Machine Learning algorithms but also for their practical implementation across a wide range of applications. Their role as a bridge between abstract mathematical formalism and applied algorithmic processes underscores their centrality in contemporary Machine Learning research and practice.

1. **Constituent Components of Tensors:** Tensors are composed of several interrelated structural elements that define their behavior and application within computational frameworks.

2. **Dimension:** In most Machine Learning contexts, tensors are utilized in their two-dimensional (2D) form. This paper primarily focuses on two-dimensional tensors; however, the conceptual principles discussed are extensible to tensors of higher dimensions.

3. **Shape:** The shape of a two-dimensional tensor (matrix) is characterized by two parameters: the number of rows (*X-axis*) and the number of columns (*Y-axis*). This structure determines the organizational layout of the data within the tensor.

4. **Parameters:** The total number of parameters within a tensor is defined as the product of its dimensions, calculated as *rows $\times$ columns*. This quantity corresponds to the total number of individual elements or connections represented.

5. **Weights:** In neural network architectures, the primary tensor associated with each layer is referred to as the *weights tensor*. These tensors encode the learnable parameters that are updated during the training process.

6. **Bias:** Complementing the weights tensor, the *bias tensor* is introduced at each layer to enable additional model complexity and flexibility, allowing the network to better capture patterns in the data.

7. **Scalar vs. Gradient:** A tensor consisting of a single element with a shape of $(1, 1)$ is classified as a *scalar*. In contrast, gradients tensors representing derivatives with respect to model parameters generally possess higher-dimensional shapes and necessitate iteration mechanisms (e.g., loops) for computation across their elements.

8. **Logits:** Logits refer to the raw, unnormalized output values produced by a neural network layer, typically represented within tensors prior to the application of activation functions such as softmax or sigmoid.

9. **Element:** An element denotes a single numerical value contained within a tensor, identified by its unique position according to the tensor's dimensional structure.

$$[ [1.0, 2.1, 3.0],$$
$$[4.0, 5.0, 6.9],$$
$$[7.4, 8.2, 9.6] ]$$

Figure 3: Example: Tensor/Matrix/Gradient, with shape: (3,3), Parameters: 9, Sum: 47.2

Certain tensor operations necessitate the use of *broadcasting*, a process by which tensor shapes are automatically expanded to ensure compatibility for element-wise computations. Broadcasting is particularly critical for operations such as matrix multiplication, where the participating tensors must conform to specific dimensional requirements to permit valid arithmetic operations.

The type and complexity of tensor operations applied— ranging from scalar interactions to vector and matrix (two-dimensional tensor) manipulations are determined by both the architectural design of the network and the underlying computational objectives. Careful management of tensor shapes and operations is therefore essential to maintaining the mathematical integrity and efficiency of Machine Learning models.

## 8.1 Key Gradient Pathologies

During the Machine Learning training pipeline, gradients are subjected to repeated multiplication across multiple layers and operations. As a consequence, their magnitudes can either grow uncontrollably or diminish to near-zero values, leading to severe numerical instabilities. Two primary types of gradient pathologies are typically observed:

1. **Exploding Gradients:** This phenomenon occurs when the magnitudes of tensor elements increase exponentially during backpropagation, potentially exceeding the representational limits of the specified data type and causing numerical overflow or instability.

2. **Vanishing Gradients:** Conversely, vanishing gradients arise when the magnitudes of tensor elements diminish exponentially, approaching zero and resulting in numerical underflow. This leads to stagnation in learning, as weight updates become negligibly small.

Mitigating these issues is crucial for ensuring stable and effective neural network training, and various architectural and optimization strategies have been developed to address them.

## 8.2    Tensor Rank Versus Tensor Shape

Although often colloquially interchanged, the concepts of *rank* and *shape* possess distinct meanings within the context of Machine Learning and tensor operations:

1. **Rank (Order):** The rank of a tensor refers to the number of dimensions it possesses. For example, a scalar has rank 0, a vector has rank 1, and a matrix has rank 2.

2. **Shape:** The shape of a tensor specifies the explicit size of each dimension. For instance, a tensor with a shape of $(3, 5)$ has three rows and five columns in its two-dimensional representation.

A clear understanding of the distinction between rank and shape is essential for designing and debugging Machine Learning models, as errors in dimensional reasoning frequently underlie operational failures.

## 8.3    Axis Operations

In tensor algebra, operations along *Axis 0* are performed row-wise, whereas operations along *Axis 1* are performed column-wise. For tensors of higher dimensionality (i.e., rank greater than two), subsequent axes correspond to progressively deeper structural layers, following the same organizational pattern. A clear understanding of axis conventions is critical for correctly implementing tensor manipulations and avoiding shape-related computational errors.

## 8.4    Dot Product

The dot product is a fundamental operation in linear algebra, particularly relevant to Machine Learning applications involving vector and matrix computations. It is defined as the sum of the element-wise products of two vectors occupying parallel positions within their respective structures.

```
Function DotProduct(VectorA, VectorB) -> ScalarResult:

    // Verify that both input vectors have identical lengths.
    if VectorA.len() != VectorB.len():
        return Error("Input vectors must be of the same length.")

    // Initialize the accumulator for the dot product result.
    ScalarResult = 0

    // Perform element-wise multiplication and accumulate the sum.
    for i in range(VectorA.len()):
        ScalarResult = ScalarResult + (VectorA[i] * VectorB[i])

    return ScalarResult
```

Listing 1: Formulaic Pythonic Pseudocode for Dot Product Computation

**Use Cases of the Dot Product**

1. Aggregating features within neural network architectures, particularly in fully connected (dense) layers.

2. Computing the angle between two vectors, which provides insights into their directional alignment.

3. Serving as the basis for cosine similarity, a common metric for comparing vector representations in tasks such as information retrieval and recommendation systems.

**Interpreting the Results**

1. A large positive result indicates strong agreement between the vectors (i.e., vectors oriented in similar directions).

2. A negative result reflects opposing directions, suggesting that the vectors are substantially misaligned.

3. A result close to zero indicates orthogonality (perpendicularity) between the vectors, implying the absence of linear correlation.

A thorough understanding of the dot product is foundational for operations such as fully connected layers, similarity-based retrieval methods, and various forms of vector space modeling within Machine Learning.

$$\begin{bmatrix} 0,0,1,0 \\ 0,2,34 \\ 4,2,4,5 \end{bmatrix} + \begin{bmatrix} 1,1,1,1 \\ 0,1,31 \\ 4,2,4,0 \end{bmatrix} = \begin{bmatrix} 1,1,2,1 \\ 0,3,64 \\ 8,4,8,5 \end{bmatrix}$$

Figure 4: Tensor Addition

## 8.5   Tensor Addition

Tensor addition is a fundamental operation wherein two tensors are combined through element-wise summation. Given two tensors of compatible shapes, the addition operation is performed by adding corresponding elements from each tensor.

Unlike tensor multiplication and division, which impose stricter constraints on dimensional compatibility, tensor addition can be applied between tensors of different shapes through broadcasting, provided that the shapes align according to broadcasting rules. This flexibility makes tensor addition particularly useful for constructing and manipulating complex computational structures within Machine Learning models.

```
Function Add(Tensor1, Tensor2) -> ResultTensor:

    // Initialize an empty tensor with the same shape as Tensor1.
    ResultTensor = CreateEmptyTensor(shape=Tensor1.shape)

    // Perform element-wise addition.
    For rowIndex in range(Tensor1.rowCount()):
        For colIndex in range(Tensor1.colCount()):

            value1 = Tensor1[rowIndex][colIndex]
            value2 = Tensor2[rowIndex][colIndex]
            ResultTensor[rowIndex][colIndex] = value1 + value2

    Return ResultTensor
```

Listing 2: Formulaic Pythonic Pseudocode for Tensor Addition

Tensor addition is widely utilized in Machine Learning workflows for a variety of purposes. It can be employed to merge gradients or feature representations during training, as well as to incorporate bias vectors into the outputs of neural network layers.

Additionally, tensor addition plays a critical role in broadcasting operations, enabling tensors of differing shapes to be expanded to a common shape. This alignment is essential for performing matrix-wise operations that require shape compatibility between operands.

## 8.6   Tensor Subtraction

Tensor subtraction constitutes an element-wise operation wherein corresponding elements from two tensors are subtracted. Structurally, it mirrors the process of tensor addition, differing only in the arithmetic

operation applied.

Given two tensors of compatible shapes, subtraction is performed by subtracting each element of the second tensor from the corresponding element of the first tensor. This operation maintains dimensional consistency and is fundamental in a variety of Machine Learning contexts, such as calculating residuals, adjusting feature maps, or implementing optimization updates.

```
1    Function Subtract(Tensor1, Tensor2) -> Tensor:
2
3        ResultTensor = create empty tensor with same shape as Tensor1
4
5        For row in Tensor1:
6            For col in row:
7
8                value1 = Tensor1[rowIndex][colIndex]
9                value2 = Tensor2[rowIndex][colIndex]
10               ResultTensor[rowIndex][colIndex] = value1 - value2
11
12       Return ResultTensor
```

Listing 3: Formulaic Pythonic Code for Tensor Subtraction

## 8.7  Scalar Addition

Scalar addition in the context of tensors typically involves an element-wise operation, where a scalar value is added individually to each element of the tensor. This approach results in a uniform transformation across all entries, preserving the relative structure of the original tensor while uniformly shifting its values.

Alternatively, scalar addition can be implemented via tensor broadcasting, wherein the scalar is implicitly expanded to match the shape of the tensor. While broadcasting enables computational efficiency, the scalar's impact may be perceived as minimal, particularly at specific coordinates (e.g., $(x_1, y_2)$), where its relative contribution may be diminished by the broader dimensional context.

```
1   Function Add_Scalar(Tensor, Scalar) -> Tensor:
2     For row in Tensor:
3        For col in row:
4           col += Scalar
```

Listing 4: Formulaic Pythonic Code for Scalar Addition

## 8.8  Scalar Subtraction

One common way to perform Scalar Subtraction is through element-wise operations across an entire tensor, with the result updating the element. In this approach, the scalar value is substracted individually to each element of the tensor. This ensures a uniform transformation throughout the tensor. Alternatively, a scalar can be subtracted using tensor broadcasting during tensor subtraction, but this often results in the scalar having minimal impact, especially at specific coordinates like (x1, y2), where its effect may be diluted by the dimensional context.

```
1   Function Sub_Scalar(Tensor, Scalar) -> Tensor:
2     For row in Tensor:
3        For col in row:
4           col -= Scalar
```

Listing 5: Formulaic Pythonic Code for Scalar Subtraction

## 8.9  Scalar Multiplication

Scalar multiplication constitutes an element-wise operation in which a scalar value is multiplied with each individual element of a tensor. Structurally, it mirrors tensor addition and subtraction in that the operation is applied uniformly across all tensor indices; however, the arithmetic operation performed is multiplication.

This technique enables uniform scaling of tensor values and is fundamental in a variety of Machine Learning applications, including feature normalization, gradient scaling, and adjusting activation outputs.

```
1
2   Function MultiplyScalar(Tensor, Scalar) -> ResultTensor:
3
4       // Initialize an empty tensor with the same shape as the input tensor.
5       ResultTensor = CreateEmptyTensor(shape=Tensor.shape)
6
7       // Perform element-wise scalar multiplication.
8       For rowIndex in range(Tensor.rowCount()):
9           For colIndex in range(Tensor.colCount()):
10
```

```
11            ResultTensor[rowIndex][colIndex] = Tensor[rowIndex][colIndex] * Scalar
12
13      Return ResultTensor
```

Listing 6: Formulaic Pythonic Pseudocode for Scalar Multiplication

## 8.10   Scalar Division

Scalar division constitutes an element-wise operation analogous to scalar multiplication, differing only in the arithmetic operation applied. Instead of multiplying each element of a tensor by a scalar, each element is divided by the scalar value.

This operation is commonly employed for normalization tasks, scaling tensor values to a desired range, and adjusting model parameters during optimization procedures.

```
1
2   Function DivideScalar(Tensor, Scalar) -> ResultTensor:
3
4      // Initialize an empty tensor with the same shape as the input tensor.
5      ResultTensor = CreateEmptyTensor(shape=Tensor.shape)
6
7      // Perform element-wise scalar division.
8      For rowIndex in range(Tensor.rowCount()):
9          For colIndex in range(Tensor.colCount()):
10
11            ResultTensor[rowIndex][colIndex] = Tensor[rowIndex][colIndex] / Scalar
12
13      Return ResultTensor
```

Listing 7: Formulaic Pythonic Pseudocode for Scalar Division

## 8.11   Use Cases of Tensor-Scalar Operations

Tensor-scalar operations are critical components in Machine Learning workflows, serving to reduce complex multi-dimensional structures into singular, interpretable scalar values. These operations underpin several fundamental aspects of modern computational models, including:

- **Loss Functions:** Quantifying prediction errors by aggregating differences between predicted and actual outputs into a single optimization target.

- **Evaluation Metrics:** Assessing model performance through scalar scores such as accuracy, precision, recall, or F1 score.

- **Regularization Terms:** Penalizing model complexity by incorporating scalar-valued norms (regularization) into the loss function to encourage simpler models.

- **Gradient Aggregation:** Summarizing parameter updates across batches or layers to stabilize and guide the optimization process.

- **State Monitoring:** Tracking statistical properties of tensors, such as mean, variance, or norm, to facilitate model diagnostics and performance tuning.

The ability to seamlessly transition between tensorial representations and scalar summaries is foundational for both the theoretical formulation and the practical implementation of Machine Learning systems.

Figure 5: Example of Matrix Multiplication.

## 8.12  Matrix Multiplication

Matrix multiplication constitutes one of the more computationally intricate operations performed on tensors. In biological systems, such as the human brain, neurons operate in parallel, dynamically adjusting the strengths of their connections to facilitate learning and adaptation. This dynamic behavior can be abstractly modeled in Machine Learning through matrix multiplication (`MatMul`), enabling the simulation of neural activity in a form amenable to the linear and deterministic processing style of digital computing hardware.

In formal terms, matrix multiplication involves the systematic combination of two matrices by multiplying corresponding elements and summing the results across aligned dimensions. Crucially, for matrix multiplication to be well-defined, the number of columns in the first matrix must equal the number of rows in the second matrix. Furthermore, efficient computation of matrix multiplication relies on optimized access patterns to large, contiguous sections of memory, underscoring its centrality to both the mathematical formulation and the hardware acceleration of Machine Learning algorithms.

```
1  Function MatMul(TensorA, TensorB) -> TensorC:
2
3      // Retrieve matrix dimensions.
4      rowsA = TensorA.rowCount()
5      colsA = TensorA.colCount()
6      colsB = TensorB.colCount()
7
8      // Initialize the output tensor with zeros, shaped (rowsA, colsB).
9      TensorC = CreateEmptyTensor(shape=(rowsA, colsB))
10
11     // Perform matrix multiplication.
12     For rowIndex in range(rowsA):
13         For colIndex in range(colsB):
14             sum = 0
15             For k in range(colsA):
16                 sum += TensorA[rowIndex][k] * TensorB[k][colIndex]
17             TensorC[rowIndex][colIndex] = sum
18
```

Listing 8: Formulaic Pythonic Pseudocode for Matrix Multiplication

It is important to note that the index **k** represents the shared dimension between the two matrices being multiplied. During the matrix multiplication process, **k** iterates over the corresponding elements of the first matrix's row and the second matrix's column, accumulating their products to compute each element of the resulting matrix. This structure constitutes the fundamental dot product operation underlying matrix multiplication and is critical for enabling efficient linear parallelization in modern computational implementations.

### 8.12.1   Misconception: Matrix Division?

There is no Matrix Division, while Scalars could apply division to a whole Tensor, a Matrix applyig division, is not a defined general operation. As reverse transformations is too computational expansive.

## 8.13  Broadcasting

Broadcasting is a fundamental mechanism that enables element-wise operations between tensors of differing shapes by automatically expanding their dimensions to achieve compatibility. The *Tensor Broadcasting Rules* formally define the conditions under which such expansions are permissible, ensuring that operations can be applied without explicit replication of data.

By adhering to these rules, Machine Learning frameworks can efficiently perform arithmetic operations across tensors of varying ranks and shapes, facilitating flexible model design and reducing computational overhead.

1. **Broadcasting Rules for Two-Dimensional Tensors:**

2. Let Tensor A and Tensor B be the two tensors involved in the operation.

3. Broadcasting is permissible along the row dimension if the number of rows in A equals the number of rows in B, or if one of the tensors has exactly one row.

4. Broadcasting is permissible along the column dimension if the number of columns in A equals the number of columns in B, or if one of the tensors has exactly one column.

5. If either the row or column dimension is 1 in one tensor, it is virtually expanded (broadcast) to match the corresponding dimension of the other tensor.

In basic terms, if a tensor dimension equals one, it can be broadcast by expanding that dimension to match the corresponding dimension of the other tensor. However, if the dimensions are neither equal nor one, additional techniques must be employed to achieve shape compatibility. These techniques include **zero-padding**, **reshaping**, **unsqueezing**, and **transposition**, each of which modifies the tensor's structure to enable valid operations.

```
1   Function broadcastable(TensorA, TensorB) -> Tensor:
2
3     if TensorA.rows == TensorB.rows or TensorA.rows == 1 or TensorB.rows == 1:
4       if TensorA.cols == TensorB.cols or TensorA.cols == 1 or TensorB.cols == 1:
5           //broadcast is valid
6     else:
7         //broadcast is invalid
```

Listing 9: Formulaic Pythonic Code for Tensor Subtraction



Figure 6: Example of correct broadcasting.

This verification step is particularly critical for **matrix multiplication**, as the shapes of the input tensors must either be inherently compatible or made compatible through broadcasting in order for the operation to be correctly applied.

## 8.14 Tensor Raised to a Scalar Power

Raising a tensor to a scalar power constitutes an element-wise operation, wherein each individual element of the tensor is exponentiated independently by the scalar value. Unlike matrix multiplication, this operation does not involve any interaction between different elements or dimensions; each computation is performed locally at the element level.

```
1  Function RaiseByScalar(Tensor, Scalar) -> ResultTensor:
2
3    // Initialize an empty tensor with the same shape as the input tensor.
4    ResultTensor = CreateEmptyTensor(shape=Tensor.shape)
5
6    // Perform element-wise exponentiation.
7    For rowIndex in range(Tensor.rowCount()):
8        For colIndex in range(Tensor.colCount()):
9
10           ResultTensor[rowIndex][colIndex] = Tensor[rowIndex][colIndex] ** Scalar
11
12   Return ResultTensor
```

Listing 10: Formulaic Pythonic Pseudocode for Element-wise Tensor Exponentiation

### 8.14.1 Use Cases

1. **Activation Functions:** Certain activation functions involve raising tensor elements to a scalar power, particularly in custom or experimentally-derived non-linearities.

2. **Regularization:** Element-wise exponentiation can be employed as part of regularization strategies to control model complexity and mitigate overfitting.

3. **Normalization Techniques:** Several normalization methods involve squaring tensor elements, such as in variance normalization or RMS (Root Mean Square) computations.

4. **Probability Distributions:** In probabilistic modeling, tensor elements are often raised to powers when computing likelihoods, especially in formulations involving exponential families or in certain Bayesian inference procedures.

### 8.14.2 Important Notes

1. **Element-wise Application:** Exponentiation by a scalar applies strictly in an element-wise manner. It does not constitute matrix multiplication, even if the tensor is two-dimensional.

2. **Scalar Flexibility:** The scalar exponent can be any real number, including negative values and fractions, provided that the base tensor elements support the operation (e.g., non-negative bases for real-valued outputs when using fractional exponents).

3. **Differentiability:** Element-wise exponentiation is differentiable almost everywhere and is thus suitable for use in gradient-based optimization algorithms.

4. **Numerical Stability:** The use of non-integer scalars in combination with zero or negative tensor elements may lead to undefined or unstable results, depending on the specific computational context.

5. **Nonlinear Transformation:** Raising tensor elements to a scalar power constitutes a nonlinear transformation, significantly altering the data distribution and affecting downstream computations.



Figure 7: Illustration of how scalar exponentiation transforms tensor values in a non-linear fashion. The graph plots several power operations—specifically $x^2$, $x^1$, $x^{0.5}$, $x^{-1}$, and $x^{-0.5}$— applied to an original input range of $[0, 2]$. While the linear function $x^1$ preserves proportionality, the other transformations introduce significant non-linear distortions: quadratic ($x^2$) expands larger magnitudes, square root ($x^{0.5}$) compresses them, and negative powers ($x^{-1}$, $x^{-0.5}$) produce inverse relationships with rapid growth near zero. This behavior emphasizes that scalar exponentiation inherently constitutes a non-linear operation, often introducing asymmetries, discontinuities, or sensitivity to small inputs within tensor transformations.

Figure 8: Diagram of a Tensor with Transposed applied, demonstrating the rotation.

# 9 Tensor Shaping and Broadcasting Operations

In Machine Learning, tensors frequently require structural transformations to ensure compatibility with model architectures. This section presents key tensor shaping and broadcasting operations that facilitate seamless data flow across different layers and model components, while preserving the integrity of the underlying data representations.

## 9.1 Transpose

The transpose operation enables tensors of differing shapes to achieve matrix-wise compatibility, particularly in contexts where broadcasting is insufficient. Transposition reorients a tensor by swapping its dimensions, effectively rotating the structure to align with the dimensional requirements of another tensor.

For example, a tensor with shape (5,1) can be transposed to shape (1,5), facilitating element-wise or matrix operations with tensors of compatible shapes.

```
Function Transpose(Tensor) -> TransposedTensor:

  // Retrieve original dimensions.
  rows = Tensor.rowCount()
  cols = Tensor.colCount()

  // Initialize an empty tensor with swapped dimensions.
  TransposedTensor = CreateEmptyTensor(shape=(cols, rows))

  // Perform transposition by swapping rows and columns.
  For rowIndex in range(rows):
      For colIndex in range(cols):
          TransposedTensor[colIndex][rowIndex] = Tensor[rowIndex][colIndex]

  Return TransposedTensor
```

Listing 11: Formulaic Pythonic Pseudocode for Tensor Transposition

This operation is particularly useful for aligning tensor shapes to enable compatibility in operations such as matrix multiplication.

## 9.2   Padding (Zero Padding)

When tensors do not possess the required dimensions for a given operation, padding can be applied to adjust their shapes accordingly. Padding typically operates in two modes: either by appending zero-valued elements to extend rows or columns, or by trimming excess elements to align the tensor with the desired target shape.

Zero padding is particularly common in Machine Learning applications where maintaining spatial or dimensional consistency across layers is critical, such as in convolutional neural networks (CNNs) and structured input pipelines.

```
1  Function PadTensor(Tensor, TargetShape, ZeroValue) -> PaddedTensor:
2
3    // Initialize a copy of the original tensor to modify.
4    PaddedTensor = CopyTensor(Tensor)
5
6    // Pad rows if necessary.
7    While PaddedTensor.rowCount() < TargetShape.rows:
8        PaddedTensor.appendRow(FillArray(length=TargetShape.columns, value=ZeroValue))
9
10   // Pad columns if necessary.
11   For rowIndex in range(PaddedTensor.rowCount()):
12       While PaddedTensor[rowIndex].length() < TargetShape.columns:
13           PaddedTensor[rowIndex].append(ZeroValue)
14
15   Return PaddedTensor
```

Listing 12: Formulaic Pythonic Pseudocode for Tensor Padding

## 9.3   Flatten

In Machine Learning, particularly within neural network architectures, the *flatten* operation is used to convert multi-dimensional data structures (e.g., two-dimensional images or three-dimensional feature maps) into a one-dimensional vector.

This transformation is essential because fully connected (dense) layers require one-dimensional input for processing. Flattening serves as a critical bridge between feature extraction layers—such as convolutional or pooling layers—and subsequent classification or output layers, enabling the network to transition from complex spatial representations to scalar predictions or decisions.

```
1  Function Flatten(Tensor) -> FlattenedTensor:
2
3    // Initialize an empty one-dimensional tensor.
4    FlattenedTensor = CreateEmptyTensor(shape=(1, 0))
5
6    // Append each element in row-major order.
7    For rowIndex in range(Tensor.rowCount()):
8        For colIndex in range(Tensor.colCount()):
9            FlattenedTensor[0].append(Tensor[rowIndex][colIndex])
10
11   Return FlattenedTensor
```

Listing 13: Formulaic Pythonic Pseudocode for Tensor Flattening

Flattening preserves the original element order in a row-major sequence, which is essential for maintaining the structural consistency of features within Machine Learning models.

## 9.4 Dimensional Products

In tensor analysis, the dimensional product corresponds to the total number of discrete elements instantiated within a tensor structure. It is formally derived as the product of the cardinalities along each axis of the tensor's shape tuple. This invariant is paramount in validating reshape operations, as the conservation of element cardinality is a necessary condition for the consistency and reversibility of tensor transformations.

```
Function Product(Shape_Array) -> Integer:
  If Shape_Array is empty:
      Return 1

  Partition Shape_Array into subsets for concurrent processing
  Initialize partial_products as empty list

  For each subset concurrently:
      local_product = 1
      For each dim in subset:
          local_product *= dim
      Append local_product to partial_products

  Initialize final_result = 1
  For each local_product in partial_products:
      final_result *= local_product

  Return final_result
```

Listing 14: Concurrent Formulaic Pythonic Code for Dimensional Product Calculation

## 9.5 Reshape and Resize Operations

During the construction and execution of neural network architectures, it is often encountered that a tensor's intrinsic shape does not conform to the dimensional requirements of specific operations, notably matrix multiplication and related algebraic transformations. In such instances, it becomes necessary to perform reshape or resize operations, wherein the tensor's structure is reconfigured to align with the expected dimensionalities of subsequent computational stages. Crucially, these transformations must preserve the total cardinality of the tensor elements, thereby maintaining data integrity while ensuring compatibility across computational layers.

```
Function Reshape(Tensor, New_Shape) -> Tensor:
  // Step 1: Flatten the original Tensor into a contiguous 1D array
  flat_data = Flatten(Tensor)

  // Step 2: Validate that reshaping is feasible
  if Product(New_Shape) != Length(flat_data):
      Throw Error("Shape mismatch: total element count inconsistent")

  // Step 3: Initialize an empty Tensor structure with the target shape
  Reshaped_Tensor = Empty_Tensor(New_Shape)

```

```
12    // Step 4: Concurrently populate the new Tensor structure
13    Parallel for index in Range(0, Length(flat_data)):
14        destination_indices = Compute_MultiDimensional_Indices(index, New_Shape)
15        Reshaped_Tensor[destination_indices] = flat_data[index]
16
17    Return Reshaped_Tensor
```

Listing 15: Concurrent Formulaic Pythonic Code for Tensor Reshape/Resize

Reshape operations reconfigure the memory layout of tensors while preserving the total cardinality of elements. Crucially, reshaping is a valid and safe transformation only when the product of the target dimensions exactly equals the original tensor's element count. Failure to satisfy this invariant results in structural inconsistency and must be guarded against to avoid runtime violations. Concurrent reshape implementations utilize independent index mappings to facilitate parallel population of the target tensor without contention.

# 10  Fundamental Mathematical Concepts for Machine Learning

A common misconception among aspiring practitioners is that effective engagement with Machine Learning (ML) requires a comprehensive background in advanced calculus or higher mathematics. While such knowledge may enrich one's theoretical understanding, the foundational practices of ML rely primarily on a small subset of core mathematical principles. This section delineates those essential constructs, providing both conceptual clarity and pedagogical accessibility.

## 10.1  Calculus: A Descriptive Framework, Not a Barrier

Calculus serves as a formal mathematical framework for describing continuous change, and historically underpins the derivation of many algorithmic foundations in Machine Learning, including gradient descent, optimization, and activation dynamics. In essence, it provides the language in which many core ML operations were originally formulated. However, for practical implementation, these expressions are routinely translated into discrete algorithmic procedures—what we term *formulaic pseudocode*—rendering deep calculus fluency unnecessary for initial engagement. Practitioners can thus implement state-of-the-art models by understanding the intuition behind derivatives and gradients, even without engaging in symbolic manipulation.

## 10.2  Non-Linearity: The Essential Break from Linearity

In the context of Machine Learning, non-linearity refers to transformations wherein the output is not a direct proportional mapping of the input. Mathematically, this property manifests in functions whose graphs are not straight lines—such as polynomials, exponentials, or activation functions like ReLU and GELU. The introduction of non-linear transformations within models is critical, as it enables neural networks to approximate complex, non-trivial patterns in high-dimensional data. Without non-linearity, models would be restricted to learning only affine transformations, severely limiting their expressive capacity.

## 10.3  Euler's Number: A Fundamental Limit in Continuous Growth

Euler's number, denoted as $e \approx 2.71828$, emerges as a foundational constant in both calculus and Machine Learning. It characterizes the asymptotic behavior of processes involving continuous exponential growth and decay, as well as the natural logarithmic base in optimization landscapes. Formally, Euler's number may be defined as the infinite limit of either a summation or a compounding product. The most common series representation is given by:

$$e = \sum_{n=0}^{\infty} \frac{1}{n!}$$

This identity expresses $e$ as the sum of the reciprocals of factorials—a rapidly converging series foundational to both exponential function approximations and gradient-based methods involving activation functions such as softmax.

The following pseudocode outlines an iterative approach to approximate $e$ using its factorial series expansion:

```
function calculate_e(terms):
    e = 0
    factorial = 1
```

```
4
5    for n from 0 to terms - 1:
6        if n > 0:
7            factorial = factorial * n
8        e = e + (1 / factorial)
9
10   return e
```

Listing 16: Formulaic Pythonic Code for Euler's Number Approximation

Most modern programming languages provide Euler's number as a built-in constant within their standard mathematical libraries (e.g., `math.e` in Python, `std::numbers::e` in C++20). However, explicit calculation via factorial expansion serves as a valuable pedagogical exercise for understanding its convergence behavior and underlying mathematical structure.

## 10.4 Mean of a Tensor: A Measure of Central Tendency

In Machine Learning, the mean—or arithmetic average—of a tensor's elements provides a fundamental measure of central tendency. This scalar value summarizes the overall distribution of the data and is frequently employed in normalization schemes, statistical analyses, and loss function diagnostics. For a tensor $T \in \mathbb{R}^{d_1 \times d_2 \times \cdots \times d_n}$, the mean is defined as:

$$\mu = \frac{1}{N} \sum_{i=1}^{N} T_i$$

where $N = \prod_{k=1}^{n} d_k$ denotes the total number of elements in the tensor, and $T_i$ represents the $i$-th flattened element.

The operation is both commutative and associative, enabling safe concurrent or parallelized computation. The following pseudocode illustrates a basic implementation:

```
1  function Tensor_Mean(Tensor) -> f32:
2      data_sum = 0                    // Initialize accumulator
3      data_points = Product(Tensor.Shape)  // Compute total elements
4
5      for row in Tensor:
6          for col in row:
7              data_sum += col
8
9      return data_sum / data_points
```

Listing 17: Formulaic Pythonic Code for Tensor Mean

In practice, most numerical computing libraries (e.g., NumPy, TensorFlow, PyTorch) provide optimized implementations that support high-dimensional tensors and exploit vectorized operations or GPU acceleration. Nevertheless, the explicit formulation serves as an instructive model for understanding the mechanics of aggregation in tensor algebra.

## 10.5 Median of a Tensor: A Robust Measure of Central Tendency

In statistical analysis of tensor data, the *median* serves as a robust estimator of central tendency, particularly effective in scenarios involving skewed distributions or the presence of outliers. Unlike the mean, which

aggregates all values and may be influenced by extreme values, the median captures the central value of the dataset after ordering the elements by magnitude.

Formally, given a tensor $T$ of arbitrary rank with $n$ total elements, the median is defined by first flattening the tensor into a one-dimensional array and then sorting its entries in non-decreasing order. The median value $m$ is then determined as follows:

$$
m = \begin{cases} x_{\frac{n+1}{2}}, & \text{if } n \text{ is odd} \\ \frac{1}{2}(x_{\frac{n}{2}} + x_{\frac{n}{2}+1}), & \text{if } n \text{ is even} \end{cases}
$$

where $x_i$ denotes the $i$-th smallest element in the sorted array.

This procedure ensures that 50% of the data lies below and above the computed median. The median is particularly advantageous in Machine Learning workflows for tasks such as feature scaling, distributional normalization, and robust data imputation.

```
1  function Tensor_Median(Tensor) -> f32:
2      flat = Flatten(Tensor)
3      sorted_data = Sort(flat)
4      n = Length(sorted_data)
5
6      if n % 2 == 1:
7          return sorted_data[(n - 1) // 2]
8      else:
9          mid1 = sorted_data[(n // 2) - 1]
10         mid2 = sorted_data[n // 2]
11         return (mid1 + mid2) / 2
```

Listing 18: Formulaic Pythonic Code for Tensor Median

While computationally more expensive than the mean, the median offers resilience against noisy data and extreme values, making it a preferred choice in adversarial or corrupted input environments. Optimized implementations often leverage partial sorting or selection algorithms (e.g., QuickSelect) to achieve linear-time complexity in practice.

## 10.6  Element-wise Absolute Value of a Tensor

The element-wise absolute value operation is a fundamental transformation in numerical computing and Machine Learning, frequently employed in contexts such as error computation, signal processing, and gradient stabilization. For a given tensor $T \in \mathbb{R}^{d_1 \times d_2 \times \cdots \times d_n}$, the absolute value operator $\text{abs}(T)$ is defined as:

$$
\text{abs}(T)_{i_1, i_2, \ldots, i_n} = |T_{i_1, i_2, \ldots, i_n}|
$$

where $|\cdot|$ denotes the scalar absolute value. This operation ensures that all resulting elements are non-negative, as it transforms each negative entry into its positive counterpart while preserving non-negative values.

In vectorized frameworks, this transformation is typically applied in a fully parallelized manner, operating independently on each tensor element. Its utility spans across applications such as the L1 norm, absolute error metrics (e.g., MAE), and robust loss functions that penalize deviations symmetrically regardless of sign.

```
1  function Tensor_Absolute(Tensor) -> Tensor:
2      result = Empty_Tensor(Shape = Tensor.Shape)
3
4      for i from 0 to Tensor.Shape.x - 1:
5          for j from 0 to Tensor.Shape.y - 1:
6              result[i][j] = abs(Tensor[i][j])
7
8      return result
```

Listing 19: Formulaic Pythonic Code for Element-wise Tensor Absolute Value

Most high-performance libraries provide optimized implementations of this operation using SIMD instructions or GPU kernels. Nonetheless, the conceptual clarity of the element-wise absolute function makes it a foundational operation in both algorithm design and educational formulations.

## 10.7  Tensor Standard Deviation (std): A Measure of Dispersion

The standard deviation is a fundamental statistical metric used to quantify the dispersion or spread of data relative to its mean. Within the context of Machine Learning, computing the standard deviation of a tensor is essential for normalization techniques, statistical diagnostics, and sensitivity analysis in gradient-based optimization.

Formally, given a tensor $T$ of $n$ elements, with mean $\mu = \frac{1}{n} \sum_{i=1}^{n} T_i$, the population standard deviation $\sigma$ is defined as:

$$\sigma = \sqrt{\frac{1}{n} \sum_{i=1}^{n} (T_i - \mu)^2}$$

In practice, the sample standard deviation may be used instead, where the denominator is $n - 1$ to provide an unbiased estimator of variance when working with sample data:

$$s = \sqrt{\frac{1}{n-1} \sum_{i=1}^{n} (T_i - \mu)^2}$$

The standard deviation thus captures the average magnitude of deviation from the mean and provides insight into the distributional shape of tensor values.

The following pseudocode outlines a generalized computation of standard deviation for a flattened tensor:

```
1  function Tensor_Standard_Deviation(Tensor, sample = False) -> f32:
2      flattened = Tensor.flatten()
3      n = Length(flattened)
4
5      if n == 0:
6          return 0
7
8      mean = Sum(flattened) / n
9      sum_squared_diff = 0
10
11     for x in flattened:
12         sum_squared_diff += (x - mean) ** 2
13
14     if sample:
```

```
15        variance = sum_squared_diff / (n - 1)
16    else:
17        variance = sum_squared_diff / n
18
19    return sqrt(variance)
```

Listing 20: Formulaic Pythonic Code for Tensor Standard Deviation

This operation is numerically stable and can be parallelized across elements, particularly when computing in GPU-accelerated environments. Most tensor libraries (e.g., PyTorch, TensorFlow, NumPy) offer optimized implementations that leverage fused operations to compute the mean and variance with minimal passes over memory.

## 10.8   Hyperparameters: Configurable Variables in Model Design and Optimization

In the architecture and training of Machine Learning models, *hyperparameters* are external configuration variables whose values are not learned from data but are instead set prior to the training process. These values govern both the structural properties of the model and the dynamics of the optimization procedure. Unlike parameters (e.g., weights and biases), which are optimized via backpropagation, hyperparameters are tuned manually or through meta-optimization strategies.

Hyperparameters can be broadly classified into the following categories:

**(1) Model Architecture Hyperparameters** define the structural configuration of the neural network:

- **Number of layers:** The depth of the network, affecting its representational capacity.

- **Neurons per layer:** The dimensionality of each transformation, controlling expressive power.

- **Activation functions:** Non-linear transformations (e.g., ReLU, Sigmoid, GELU) applied to neurons, introducing model non-linearity.

**(2) Optimization Hyperparameters** govern how the model is trained using gradient descent:

- **Learning rate:** The step size at each iteration of optimization; critical for convergence behavior.

- **Batch size:** Number of training samples processed per update; affects memory usage and gradient stability.

- **Epochs:** The number of complete passes over the training dataset.

- **Optimizer:** Algorithm used to update weights (e.g., SGD, Adam, RMSProp), each with its own internal dynamics.

**(3) Regularization Hyperparameters** mitigate overfitting and improve generalization:

- **Dropout rate:** Fraction of neurons randomly deactivated during training to prevent co-adaptation.

- **L1/L2 regularization coefficients:** Penalize large weights to encourage sparsity or smoothness.

- **Non-linearity inclusion:** The use of activation functions contributes indirectly to regularization by introducing model flexibility.

**(4) Hyperparameter Tuning Strategies** determine how hyperparameters are selected:

- **Grid search:** Exhaustive evaluation across a Cartesian product of predefined hyperparameter values.

- **Random search:** Random sampling from defined distributions over hyperparameter values.

- **Bayesian optimization:** Model-based approach that builds a probabilistic surrogate function to efficiently explore the search space.

- **Advanced methods:** Includes algorithms like Hyperband, Optuna, and Population-Based Training (PBT) for scalable and adaptive tuning.

Hyperparameter tuning remains a central component of model performance optimization and generalization. Automated hyperparameter search frameworks are increasingly integrated into modern ML workflows to streamline this process while balancing exploration and computational cost.

## 10.9 Function Hyperparameters: Tunable Modifiers of Algorithmic Behavior

In the formulation of many mathematical functions used within Machine Learning algorithms, certain parameters—termed *function hyperparameters*—are introduced to modulate the function's internal dynamics. These hyperparameters are not learned from data, but are manually specified or tuned to adapt the function's behavior to the characteristics of a given task.

Function hyperparameters often control properties such as smoothing, weighting, or curvature, and are integral to operations including optimization, regularization, and activation. Although many such functions provide empirically derived default values, adjusting these hyperparameters can significantly impact model convergence, performance, and stability.

**Example: Exponential Moving Averages in the Adam Optimizer**

The Adam optimizer, widely used for stochastic gradient-based training, maintains exponential moving averages of past gradients and squared gradients. It introduces two key hyperparameters:

- **Beta$_1$** ($\beta_1$) controls the decay rate of the moving average of the first moment (gradient). A typical default is $\beta_1 = 0.9$, meaning recent gradients are weighted more heavily than earlier ones.

- **Beta$_2$** ($\beta_2$) governs the second moment (squared gradients), often set to $\beta_2 = 0.999$ to ensure long-term memory of variance.

These hyperparameters effectively shape the optimizer's sensitivity to recent versus historical gradient information, influencing convergence speed and stability. Improper tuning can lead to oscillation or stagnation, particularly in the presence of noisy gradients or sparse features.

Function hyperparameters appear across the Machine Learning pipeline, including:

- **Dropout rate** in regularization functions

- **Alpha** in leaky ReLU or ELU activations

- **Temperature** in softmax-based sampling

- **Kernel width** in Gaussian functions or RBF kernels

Understanding and tuning these hyperparameters is crucial for aligning theoretical function behavior with empirical model dynamics.

## 10.10  Derivatives: Quantifying Change in Continuous Functions

The derivative is a foundational concept in calculus that quantifies the instantaneous rate of change of a function with respect to one of its input variables. In the context of Machine Learning, derivatives are central to optimization procedures such as gradient descent, where model parameters are iteratively updated based on the slope of the loss function.

Formally, the derivative of a scalar-valued function $f(x)$ with respect to its input $x$ is defined as the limit:

$$f'(x) = \lim_{\Delta x \to 0} \frac{f(x + \Delta x) - f(x)}{\Delta x}$$

This expression captures the rate at which $f(x)$ changes as $x$ varies, and provides the tangent slope of the function at a given point. For differentiable functions, this allows local linear approximations and facilitates analytical optimization.

**Example: Power Function Derivative**

Consider the function:

$$f(x) = x^2$$

Its derivative, computed using standard calculus rules, is:

$$f'(x) = 2x$$

This indicates that the rate of change of the function at any point $x$ is proportional to $x$ itself. Evaluating the derivative at $x = 3$, we obtain:

$$f'(3) = 2 \times 3 = 6$$

Thus, at $x = 3$, the function is increasing at a rate of 6 units for every unit increase in $x$.

In Machine Learning, derivatives are used to compute gradients—vector-valued generalizations of partial derivatives across multivariate functions—which are instrumental in parameter updates. The efficient computation of such derivatives, often via automatic differentiation, is critical to the scalability of modern deep learning frameworks.

## 10.11  Epsilon ($\varepsilon$): Thresholds for Numerical Stability and Convergence

In numerical analysis and Machine Learning algorithms, the symbol $\varepsilon$ (epsilon) conventionally denotes a small positive quantity used as a threshold for precision, stability, or convergence criteria. Unlike dynamic variables or learned parameters, epsilon is a manually specified hyperparameter that defines tolerable bounds for terminating iterative procedures or preventing division-by-zero instabilities.

Formally, in optimization algorithms such as gradient descent, an iteration may be halted once the magnitude of the gradient norm falls below a predefined threshold $\varepsilon$, indicating that further progress is negligible:

$$\|\nabla f(\mathbf{x})\| < \varepsilon$$

Similarly, in algorithms involving floating-point divisions, epsilon is often added to denominators to avoid undefined behavior:

$$\text{Adjusted division:} \quad \frac{1}{x + \varepsilon}$$

**Illustrative Example: Iterative Stopping Criterion**

Consider an iterative optimization algorithm where updates continue until:

$$|f(x_{k+1}) - f(x_k)| < \varepsilon$$

This ensures that the relative change between successive iterations is sufficiently small, implying convergence to a (local) optimum within acceptable tolerance.

**Important Clarification**

Epsilon ($\varepsilon$) should not be confused with Euler's number ($e \approx 2.71828$), which is a mathematical constant describing natural exponential growth. While $e$ is fundamental to continuous mathematical modeling, $\varepsilon$ serves as a practical device for managing numerical precision and algorithmic termination.

## 10.12   Geometric Symbols and Mathematical Constants in Machine Learning

Mathematical constants and geometric symbols play critical yet often understated roles in the formulation and analysis of Machine Learning algorithms. Although some constants, such as $\pi$, are more prominent in pure mathematics and physics, they appear in subtle yet essential ways within statistical modeling, optimization theory, and loss function derivations.

**(1) The Constant $\pi$**

The constant $\pi \approx 3.14159$ arises primarily in Machine Learning through its connections to probability theory, particularly in the context of continuous probability distributions. For instance, the probability density function of the normal (Gaussian) distribution involves $\pi$ explicitly:

$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

Here, $\pi$ ensures proper normalization of the distribution such that the total probability integrates to one. Applications of Gaussian distributions are ubiquitous in Machine Learning, including modeling noise, regularization (e.g., Gaussian priors), and understanding loss landscapes in stochastic optimization.

Moreover, $\pi$ occasionally appears in more advanced operations involving Fourier analysis, kernel methods, and information theory metrics, albeit less frequently than constants like $e$.

**(2) The Constant $e$**

Euler's number $e \approx 2.71828$ is far more central to Machine Learning, underpinning exponential functions, logarithmic transformations, and the softmax function. $e$ governs the behavior of probability distributions, learning rate decays, and loss formulations involving cross-entropy.

**(3) The Symbol $\theta$**

The Greek letter $\theta$ is widely used to denote model parameters, particularly in optimization and learning contexts. Conceptually, $\theta$ may represent:

- The parameter vector in linear models (e.g., in logistic regression, $\theta$ encapsulates weights and biases).

- The orientation or slope of a function with respect to its gradient during optimization.

In gradient-based methods such as Adam, SGD, or RMSprop, $\theta$ evolves iteratively via update rules of the form:

$$\theta_{t+1} = \theta_t - \alpha \nabla_\theta \mathcal{L}(\theta_t)$$

where $\alpha$ denotes the learning rate and $\mathcal{L}(\theta)$ the loss function. Here, the gradient $\nabla_\theta \mathcal{L}$ provides the direction and rate of steepest ascent, and updates to $\theta$ occur in the opposite direction to minimize loss.

While $\theta$ originally symbolizes angular quantities in geometry and trigonometry, its usage in Machine Learning reflects a broader abstraction: the idea of navigating parameter spaces and controlling optimization trajectories.

In summary, geometric constants and symbols, although sometimes subtle in appearance, are deeply embedded within the theoretical infrastructure of modern Machine Learning. Their understanding facilitates deeper insights into algorithmic behavior, statistical assumptions, and optimization landscapes.

## 10.13 The Error Function (erf): A Fundamental Link to Gaussian Distributions

The error function, denoted as $\text{erf}(x)$, is a special function of fundamental importance in probability theory, numerical analysis, and Machine Learning. It quantifies the probability that a random variable from a standard normal distribution falls within a given symmetric interval around the mean.

Formally, the error function is defined as:

$$\text{erf}(x) = \frac{2}{\sqrt{\pi}} \int_0^x e^{-t^2} \, dt$$

This expression describes the normalized integral of the Gaussian function $e^{-t^2}$ from 0 to $x$, scaled appropriately to map asymptotically to $\pm 1$ as $x$ tends toward $\pm\infty$.

**Applications in Machine Learning**

The erf function arises in several contexts within Machine Learning and related disciplines:

- **Activation functions:** The Gaussian Error Linear Unit (GELU) activation function, used prominently in Transformer architectures, incorporates erf for smooth non-linear transformations.

- **Optimization algorithms:** Some adaptive optimizers and smoothing functions leverage erf for stabilizing updates or defining soft thresholds.

- **Probabilistic modeling:** In Gaussian processes, diffusion models, and statistical physics-inspired ML methods, erf naturally emerges due to its relation to cumulative Gaussian distributions.

**Numerical Approximation**

Given that the error function has no closed-form expression in terms of elementary functions, numerical approximations are often employed. The following pseudocode illustrates a basic numerical integration approach using the trapezoidal rule:

```
1  function erf(x) -> f32:
2      // Constants
3      pi = std.pi
4      sqrt_pi = sqrt(pi)
5
6      // Numerical integration settings
7      N = 1000
8      dx = x / N
9
10     // Initialize integration result
11     result = 0
12
13     // Apply trapezoidal rule for numerical integration
14     for i from 0 to N - 1:
15         t1 = i * dx
16         t2 = (i + 1) * dx
17
18         result += (exp(-t1**2) + exp(-t2**2)) * (t2 - t1) / 2
19
20     // Scale according to erf definition
21     return (2 / sqrt_pi) * result
```

Listing 21: Formulaic Pythonic Code for Numerical Approximation of erf

Most modern programming environments, such as Python (`math.erf`), C++ (`std::erf`), and numerical libraries (e.g., SciPy), provide highly optimized native implementations of the error function. When available, these should be preferred for computational efficiency and numerical stability.

# 11 Activation Functions: Modeling Complex Neuronal Dynamics

## 11.1 Meaning of Activation

In the context of artificial neural networks, an *activation* refers to the output response of a neuron given its input signal. Conceptually, it quantifies the extent to which a neuron "fires" based on the strength and characteristics of incoming information. Mathematically, activation functions determine whether a neuron transmits its signal onward within the network architecture.

## 11.2 Purpose and Biological Motivation

Biological neurons are interconnected through synapses—specialized junctions where chemical and electrical signals determine whether the post-synaptic neuron becomes activated. The decision to fire is influenced not only by the magnitude of incoming signals but also by complex, nonlinear biochemical processes occurring within the synaptic cleft.

Simulating these intricate processes at a molecular level is computationally intractable for large-scale models. Consequently, artificial neural networks approximate neuronal dynamics through the use of mathematically defined *activation functions*. These functions capture the essence of biological activation by introducing nonlinear transformations that allow networks to approximate complex, high-dimensional mappings between inputs and outputs.

**Mathematical Role of Activation Functions**

Activation functions serve two primary roles:

- **Non-linearity introduction:** Without non-linear activations, a network composed solely of affine transformations (matrix multiplications and bias additions) would collapse into a single equivalent linear transformation, severely limiting representational capacity. Non-linear activations (e.g., ReLU, sigmoid, GELU) enable networks to approximate arbitrary continuous functions, a property formally guaranteed by the Universal Approximation Theorem.

- **Information shaping:** Activations modulate the flow of information through the network layers, allowing selective amplification, suppression, or transformation of learned features. This shaping mechanism is crucial for the emergence of hierarchical feature representations in deep architectures.

Thus, activation functions are indispensable both biologically and computationally. They enable artificial systems to mirror key principles of biological cognition while simultaneously expanding their mathematical expressive power beyond that of simple linear models.

Figure 9: Comparison of common activation functions (ReLU, Leaky ReLU, GELU, Swish, Tanh, Sigmoid, Softmax)

## 11.3 Activation Function Progression and Comparative Analysis

The selection and ordering of activation functions in this section is designed to introduce increasingly sophisticated transformations, beginning with foundational piecewise-linear functions and progressing toward smooth, probabilistic activations optimized for modern deep learning architectures. The sequence culminates with the Gaussian Error Linear Unit (GELU), widely adopted in Transformer-based models.

### 11.3.1 Comparative Overview of Activation Functions

| Activation | Range | Non-linearity | Key Notes |
|---|---|---|---|
| ReLU | $[0, \infty)$ | Sharp | Simple, sparse activation; efficient for deep convolutional networks. |
| Leaky ReLU | $(-\infty, \infty)$ | Sharp | Introduces a small negative slope to address neuron inactivity. |
| Swish | $(-\infty, \infty)$ | Smooth | Self-gating; improves performance, especially in very deep networks. |
| GELU (erf) | $(-\infty, \infty)$ | Smooth | Probabilistic activation; superior performance in Transformer architectures. |
| tanh | $(-1, 1)$ | Smooth | Zero-centered; common in recurrent neural networks. |
| GELU (tanh) | $(-\infty, \infty)$ | Smooth | Computationally efficient approximation of GELU. |

Table 1: Summary Comparison of Activation Functions

### 11.3.2 Expanded Insights on Activation Functions

| Activation | Additional Insights |
|---|---|
| ReLU | Prone to the "dying ReLU" phenomenon, where neurons output zero for all inputs. Dominant in CNNs due to computational simplicity. |
| Leaky ReLU | Commonly uses a negative slope coefficient $\alpha = 0.01$; mitigates dying neuron problems without introducing significant complexity. |
| Swish | Developed by Google researchers; exhibits non-monotonicity, enabling richer feature representations. |
| GELU (erf) | Softly gates inputs based on magnitude; combines ReLU's sparsity with smoothness near the origin. |
| tanh | Saturates at the extremes, leading to vanishing gradients; more stable than sigmoid for zero-centered outputs. |
| GELU (tanh) | Provides a near-equivalent functional behavior to GELU with reduced computational overhead. |

Table 2: Detailed Observations on Activation Functions

### 11.3.3 Practical Tips and Cautions

| Activation | Usage Tips and Warnings |
|---|---|
| ReLU | Recommended as a default; monitor for neuron death, especially under aggressive learning rates. |
| Leaky ReLU | Useful when ReLU leads to stagnant training; helps maintain gradient flow even for negative inputs. |
| Swish | Yields marginally better performance in deep architectures at the cost of slightly increased computational burden. |
| GELU (erf) | Preferred for Transformer-based architectures (e.g., BERT, GPT) due to its smooth probabilistic gating behavior. |
| tanh | Best suited for recurrent architectures where centered activations improve gradient dynamics. |
| GELU (tanh) | Ideal for hardware-constrained environments seeking GELU-like performance with faster evaluation. |

Table 3: Practical Guidelines for Activation Function Selection

## 11.4 Rectified Linear Unit (ReLU): Definition, Properties, and Practical Considerations

### 11.4.1 Intuitive Overview

The Rectified Linear Unit (ReLU) activation function introduces non-linearity by transforming all negative inputs to zero while preserving positive values. Conceptually, ReLU "filters" input signals, suppressing negative activations and allowing positive signals to propagate. This operation enhances the expressivity of neural networks without introducing significant computational overhead.

### 11.4.2 Mathematical Definition

Formally, ReLU is defined as an element-wise transformation:

$$\text{ReLU}(x) = \max(0, x)$$

for each input $x$. This operation is applied independently to each element in the input tensor, requiring no global context or inter-element dependencies.

**Formulaic Pseudocode Implementation:**

```
Function ReLU(Tensor) -> Tensor:
    result = Empty_Tensor(Shape = Tensor.Shape)

    for i from 0 to Tensor.Shape.x - 1:
        for j from 0 to Tensor.Shape.y - 1:
            if Tensor[i][j] > 0:
                result[i][j] = Tensor[i][j]
            else:
                result[i][j] = 0

    return result
```

Listing 22: Formulaic Pythonic Code for ReLU Activation

This element-wise independence allows efficient parallelization on modern hardware architectures, such as GPUs and TPUs.

### 11.4.3 Advantages and Use Cases

ReLU activation is widely adopted in deep learning due to the following properties:

1. **Introduction of Non-linearity:** ReLU introduces non-linear transformations while retaining simplicity.

2. **Computational Efficiency:** It requires only a simple thresholding operation, making it extremely fast to compute.

3. **Alleviation of Vanishing Gradient Problem:** Unlike saturating activations (e.g., sigmoid, tanh), ReLU maintains stronger gradient magnitudes, improving optimization in deep networks.

4. **Sparsity Induction:** ReLU naturally produces sparse activations by zeroing out negative values, reducing the effective computational burden and promoting representational efficiency.

5. **Positive Constraint:** By filtering out negative values, ReLU enforces non-negative activations, which can be advantageous in certain architectures and regularization schemes.

6. **Faster Convergence:** Empirically, networks employing ReLU converge faster during training compared to those using traditional saturating activations.

### 11.4.4 Limitations and Cautions

Despite its advantages, ReLU presents certain limitations:

1. **Limited Representational Scope:** By entirely discarding negative activations, ReLU can hinder the model's ability to learn functions where negative outputs are semantically meaningful.

2. **Dying ReLU Phenomenon:** Neurons can become permanently inactive if their inputs consistently yield negative values during training, resulting in zero gradients and loss of learning capacity for those neurons.

3. **Suitability for Complex Relationships:** For highly intricate or symmetric data patterns involving negative domains, alternative activations such as Leaky ReLU, ELU, or GELU may offer superior performance.

### 11.4.5 Why is it faster to train?

You're trying to capture a simpler relationship, with a very efficient computation formula.

### 11.4.6 Alternative for ReLU, Leaky ReLU

1. Solves the Negative relationship to keep neurons out of the inactive state.

2. Able to calculate a tiny bit more relationships.

Figure 10: ReLU graphed mainly to show the curve.

### 11.4.7 ReLU graph

Notice how the graph remains flat and then abruptly shifts at a sharp angle. This discontinuity in slope is a hallmark of a non-smooth function. For example, ReLU exhibits this behavior with a sharp kink at zero. In contrast, smooth activation functions like GELU produce a more gradual polynomial-like curve with continuous derivatives. This smoothness allows the model to capture more nuanced patterns and learn more complex relationships in data, whereas non-smooth functions may limit expressiveness in certain cases.

1. **Slope, less than 0 region:** Remain at 0

2. **Slope, greater than 0 region:** Regular linear increase

## 11.5 Leaky Rectified Linear Unit (Leaky ReLU, LReLU): Definition, Properties, and Applications

### 11.5.1 Intuitive Overview

The Leaky Rectified Linear Unit (Leaky ReLU or LReLU) extends the standard ReLU activation by allowing a small, non-zero gradient when the input is negative. Instead of mapping all negative values to zero, LReLU scales them by a small constant factor, thereby maintaining a minimal signal flow even for negative activations.

### 11.5.2 Mathematical Definition

Formally, the Leaky ReLU function is defined as:

$$\text{LReLU}(x) = \begin{cases} x, & \text{if } x > 0 \\ \alpha x, & \text{otherwise} \end{cases}$$

where $\alpha$ is a small positive hyperparameter (commonly set to $\alpha = 0.01$) that controls the slope for negative inputs.

### 11.5.3 Formulaic Pseudocode Implementation

```
1  Function LReLU(Tensor, negative_slope) -> Tensor:
2      result = Empty_Tensor(Shape = Tensor.Shape)
3
4      for i from 0 to Tensor.Shape.x - 1:
5          for j from 0 to Tensor.Shape.y - 1:
6              if Tensor[i][j] > 0:
7                  result[i][j] = Tensor[i][j]
8              else:
9                  result[i][j] = negative_slope * Tensor[i][j]
10
11     return result
```

Listing 23: Formulaic Pythonic Code for Leaky ReLU Activation

Like ReLU, LReLU is an element-wise operation, enabling efficient parallelization on modern hardware architectures.

### 11.5.4 Advantages and Use Cases

Leaky ReLU offers several practical benefits:

1. **Introduction of Non-linearity:** Similar to ReLU, LReLU introduces essential non-linear transformations into the network.

2. **Mitigation of Vanishing Gradients:** By maintaining a non-zero gradient for negative inputs, LReLU reduces the risk of gradient vanishing compared to saturating functions like sigmoid or tanh.

3. **Neuron Survival:** Allows neurons to maintain small gradients even when weights initially result in negative pre-activations, preventing complete neuron inactivity ("dying ReLU" problem).

4. **Efficient Computation:** Simple to implement and computationally inexpensive, making it suitable for real-time or resource-constrained applications.

5. **Improved Convergence Speed:** Networks employing LReLU often converge faster than those using pure ReLU in scenarios where data distributions produce many negative activations.

6. **Alternative to ReLU in Hidden Layers:** Particularly beneficial in architectures where ReLU leads to excessive neuron death or sparse representations become detrimental.

### 11.5.5 Limitations and Cautions

Despite its improvements over ReLU, LReLU presents specific limitations:

1. **Limited Expressivity for Complex Patterns:** While suitable for basic feature extraction, LReLU may not sufficiently capture highly intricate or symmetric relationships that require more sophisticated activations.

2. **Retention of Negative Values:** Although small, the leakage of negative signals may not always align with tasks that benefit from strict non-negative representations.

### 11.5.6 Training Efficiency Considerations

LReLU accelerates training primarily by preserving gradient flow throughout the network, even when activations are negative. This ensures that learning signals are not entirely extinguished, facilitating more stable and faster convergence compared to ReLU in challenging data regimes.

### 11.5.7 Graphical Interpretation of Leaky ReLU



Figure 11: Visualization of Leaky ReLU Activation Function.

As illustrated in Figure 11, the behavior of LReLU can be characterized as:

- **Negative Region:** For inputs $x < 0$, the function exhibits a slight downward tilt determined by the negative slope $\alpha$, allowing a small gradient (leak) to persist.

- **Positive Region:** For inputs $x > 0$, the function behaves as a standard identity mapping ($\text{LReLU}(x) = x$), identical to ReLU behavior.

## 11.6 Swish Activation Function: Definition, Properties, and Practical Considerations

### 11.6.1 Intuitive Overview

The Swish activation function introduces a smooth, non-monotonic transformation that enhances the representational capacity of neural networks. Unlike ReLU, which abruptly zeroes out negative inputs, Swish softly suppresses negative values without completely discarding them. This property enables networks to learn more nuanced and complex patterns across the input domain.

### 11.6.2 Mathematical Definition

Formally, the Swish activation function is defined as:

$$\text{Swish}(x) = x \cdot \sigma(x)$$

where $\sigma(x)$ denotes the sigmoid function:

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

Thus, the Swish function smoothly gates input values based on their own magnitude, combining linear and non-linear behaviors adaptively.

### 11.6.3 Formulaic Pseudocode Implementation

```
Function Swish(Tensor) -> Tensor:
    result = Empty_Tensor(Shape = Tensor.Shape)

    for i from 0 to Tensor.Shape.x - 1:
        for j from 0 to Tensor.Shape.y - 1:
            result[i][j] = Tensor[i][j] * Sigmoid(Tensor[i][j])

    return result
```

Listing 24: Formulaic Pythonic Code for Swish Activation

The operation is element-wise and inherently parallelizable, making it suitable for large-scale tensor computations.

### 11.6.4 Advantages and Use Cases

Swish activation offers several compelling properties:

1. **Non-linearity Introduction:** Provides a smooth, non-monotonic mapping that enhances feature expressivity.

2. **Preservation of Negative Information:** Unlike ReLU, Swish retains small negative values, facilitating richer learning dynamics.

3. **Improved Gradient Flow:** The continuous derivative of Swish helps mitigate both vanishing and exploding gradient problems during backpropagation.

4. **Superior Performance:** Empirical studies show that Swish often outperforms ReLU on deep convolutional networks (CNNs), Transformer architectures, and large-scale classification tasks.

5. **Suitability for Deep Architectures:** Particularly effective in very deep models where smooth transitions improve optimization dynamics.

6. **Alternative to Sigmoid:** Provides smoothness without suffering from the severe saturation issues of pure sigmoid activations.

### 11.6.5 Limitations and Practical Considerations

Despite its advantages, Swish has certain drawbacks:

1. **Increased Computational Cost:** Evaluation of the sigmoid function adds overhead compared to simpler activations like ReLU.

2. **Training Latency:** Networks utilizing Swish may exhibit slower per-epoch training times relative to ReLU.

3. **Compatibility Limitations:** Some machine learning frameworks may not natively support Swish, requiring custom implementation.

4. **Tradeoff with Simplicity:** In some scenarios, the benefits of Swish may be marginal relative to its additional complexity and computational burden.

### 11.6.6 Graphical Interpretation of Swish



Figure 12: Graph of the Swish activation function illustrating its behavior across input domains.

As shown in Figure 12, the behavior of Swish can be characterized as follows:

- **Negative Region ($x < 0$):** For negative inputs, Swish exhibits a small positive slope, allowing modest signal flow. Unlike ReLU, which zeroes negative inputs, Swish softly attenuates them, preserving some gradient and information flow during training.

- **Positive Region ($x \geq 0$):** For positive inputs, Swish behaves increasingly like the identity function, with the slope approaching 1 for large positive values. This ensures that strong activations pass through with minimal distortion, similar to ReLU but with smoother transitions.

## 11.7 Sigmoid Activation Function: Definition, Properties, and Practical Considerations

### 11.7.1 Intuitive Overview

The Sigmoid activation function transforms input values into a bounded range between 0 and 1, introducing smooth non-linearity into the network. This mapping is particularly useful for interpreting outputs as probabilities or confidence scores in binary classification tasks. Sigmoid activation ensures that the output is easily interpretable while compressing extreme input values toward asymptotic limits.

### 11.7.2 Mathematical Definition

Formally, the Sigmoid function is defined as:

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

where $e$ denotes Euler's number ($e \approx 2.71828$). As $x$ tends toward positive infinity, $\sigma(x)$ approaches 1; as $x$ tends toward negative infinity, $\sigma(x)$ approaches 0. The function is continuous, differentiable, and monotonic.

### 11.7.3 Formulaic Pseudocode Implementation

```
1  Function Sigmoid(Tensor) -> Tensor:
2      result = Empty_Tensor(Shape = Tensor.Shape)
3
4      for i from 0 to Tensor.Shape.x - 1:
5          for j from 0 to Tensor.Shape.y - 1:
6              result[i][j] = 1 / (1 + exp(-Tensor[i][j]))
7
8      return result
```

Listing 25: Formulaic Pythonic Code for Sigmoid Activation

The Sigmoid operation is element-wise and highly parallelizable, enabling efficient implementation on modern computational hardware.

### 11.7.4 Advantages and Use Cases

Sigmoid activation provides several practical benefits:

1. **Smooth Non-linearity:** Introduces a continuous and differentiable mapping, essential for gradient-based optimization methods.

2. **Probabilistic Interpretation:** Outputs can be directly interpreted as probabilities, facilitating binary classification tasks.

3. **Output Layer in Binary Classification:** Commonly used in the final layer of binary classifiers to squash raw model outputs into interpretable confidence scores.

4. **Computer Vision and Signal Processing:** Used when outputs need to represent bounded, normalized measurements (e.g., pixel brightness).

### 11.7.5 Limitations and Practical Considerations

Despite its usefulness, Sigmoid presents several challenges:

1. **Vanishing Gradient Problem:** Gradients diminish near the asymptotic regions ($x \gg 0$ or $x \ll 0$), hindering learning in deep networks.

2. **Computational Overhead:** Evaluation of the exponential function incurs higher computational cost relative to piecewise-linear activations like ReLU.

3. **Non-zero-centered Output:** Since outputs are strictly positive, sigmoid activations can introduce bias in gradient updates, slowing convergence.

4. **Poor Suitability for Hidden Layers:** Using sigmoid in hidden layers can severely impair training efficiency; modern practice favors ReLU, GELU, or tanh for intermediate layers.

### 11.7.6 Graphical Interpretation and Application in Binary Classification



Figure 13: Graph of the Sigmoid activation function. Although it visually appears to saturate at 0 and 1, these values are approached asymptotically but never exactly attained.

As depicted in Figure 13, the Sigmoid function:

- Compresses input values into the $(0, 1)$ interval.

- Exhibits near-linear behavior around $x = 0$, enhancing sensitivity to small input changes.

- Flattens for large positive or negative inputs, where gradient magnitudes shrink dramatically.

**Example Application:**
When applied to the output layer of a neural network for binary classification:

- An output close to 1 indicates high confidence in the positive class.

- An output close to 0 indicates high confidence in the negative class.

## 11.8  Gaussian Error Linear Unit (GELU, `erf`-based): Definition, Properties, and Practical Applications

### 11.8.1  Intuitive Overview

The Gaussian Error Linear Unit (GELU) activation function introduces smooth, stochastic-like non-linearity by blending input information based on a probabilistic model. Instead of deterministically zeroing or scaling inputs, as in ReLU or Leaky ReLU, GELU softly gates them using the Gaussian error function (erf). This property enables networks to learn more nuanced, expressive representations, especially in large-scale architectures.

### 11.8.2  Mathematical Definition

The original GELU activation, based on the error function, is defined as:

$$\text{GELU}(x) = 0.5x \left( 1 + \text{erf}\left( \frac{x}{\sqrt{2}} \right) \right)$$

where $\text{erf}(\cdot)$ is the Gauss error function:

$$\text{erf}(x) = \frac{2}{\sqrt{\pi}} \int_0^x e^{-t^2} dt$$

This formulation approximates the cumulative distribution function (CDF) of a standard normal distribution, blending the input magnitude with its probability of being positive under Gaussian assumptions.

### 11.8.3  Formulaic Pseudocode Implementation

```
Function GELU_erf(Tensor) -> Tensor:
    result = Empty_Tensor(Shape = Tensor.Shape)

    for i from 0 to Tensor.Shape.x - 1:
        for j from 0 to Tensor.Shape.y - 1:
            result[i][j] = 0.5 * Tensor[i][j] * (1 + erf(Tensor[i][j] / sqrt(2)))

    return result
```

Listing 26: Formulaic Pythonic Code for GELU (erf) Activation

As with other activation functions, GELU is applied in an element-wise fashion, facilitating efficient parallelization.

### 11.8.4  Advantages and Use Cases

GELU (erf-based) provides several critical advantages:

1. **Smooth Non-linearity:** Balances the sparsity introduced by ReLU with the smoothness of sigmoid-like activations.

2. **Superior Performance in Deep Networks:** Empirically shown to improve training stability and accuracy, particularly in Transformer-based architectures.

3. **Stable Gradient Flow:** The smooth, differentiable structure ensures robust gradient propagation even across very deep layers.

4. **Model Expressiveness:** Enables the capture of subtle relationships in high-dimensional data spaces.

5. **High Performance-to-Complexity Ratio:** Especially effective in scaling models to billions of parameters (e.g., BERT, GPT architectures).

### 11.8.5   Limitations and Practical Considerations

Despite its advantages, GELU (erf-based) introduces specific computational tradeoffs:

1. **Increased Computational Complexity:** Evaluation of the error function is significantly more expensive than basic ReLU or even Swish activations.

2. **Hardware Efficiency:** Less suited for deployment on edge devices or real-time applications due to the higher evaluation cost.

3. **Potential Overfitting Risk:** Due to its ability to closely model complex data patterns, GELU may overfit small or noisy datasets.

4. **Training Latency:** Slower per-epoch training times compared to simpler activations.

5. **Alternative Approximations:** Faster approximations, such as GELU(tanh), are often preferred in production environments to reduce inference time.

### 11.8.6 Graphical Interpretation of GELU (erf)



Figure 14: Graph of the GELU (erf-based) activation function, illustrating its smooth gating behavior across input domains.

As illustrated in Figure 14, the behavior of GELU can be interpreted as:

- **Negative Region ($x < 0$):** For negative inputs, the output decays smoothly toward zero without abrupt cutoffs. This soft suppression enables minor negative information to contribute to feature learning, while preserving stable gradient flow.

- **Positive Region ($x \geq 0$):** For positive inputs, the activation increases approximately linearly, asymptotically approaching the identity line. This behavior maintains strong gradient magnitudes for positive signals, supporting efficient learning.

## 11.9 Hyperbolic Tangent (tanh): Definition, Properties, and Practical Applications

### 11.9.1 Intuitive Overview

The hyperbolic tangent (tanh) activation function introduces smooth, zero-centered non-linearity by mapping real-valued inputs into a bounded range between $-1$ and $1$. Conceptually, tanh compresses inputs while preserving their sign and relative magnitude, enabling neural networks to learn balanced representations where both positive and negative activations are meaningful.

### 11.9.2 Mathematical Definition

Formally, the tanh function is defined as:

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

where $e$ denotes Euler's number ($e \approx 2.71828$). As $x$ approaches positive infinity, $\tanh(x) \to 1$; as $x$ approaches negative infinity, $\tanh(x) \to -1$. The function is smooth, continuous, and differentiable everywhere.

### 11.9.3 Formulaic Pseudocode Implementation

```
1  Function tanh(Tensor) -> Tensor:
2      result = Empty_Tensor(Shape = Tensor.Shape)
3
4      for i from 0 to Tensor.Shape.x - 1:
5          for j from 0 to Tensor.Shape.y - 1:
6              numerator = exp(Tensor[i][j]) - exp(-Tensor[i][j])
7              denominator = exp(Tensor[i][j]) + exp(-Tensor[i][j])
8              result[i][j] = numerator / denominator
9
10     return result
```

Listing 27: Formulaic Pythonic Code for tanh Activation

As with other activations discussed, tanh operates element-wise, enabling straightforward parallelization in high-dimensional tensor computations.

### 11.9.4 Advantages and Use Cases

The tanh activation function remains relevant in various Machine Learning applications due to the following properties:

1. **Zero-centered Output:** Unlike sigmoid, tanh outputs are centered around zero, facilitating more balanced gradient updates during optimization.

2. **Smooth Non-linearity:** Provides smooth, continuous gradients conducive to stable learning.

3. **Use in Recurrent Architectures:** Commonly employed in Recurrent Neural Networks (RNNs) where controlling positive and negative signal propagation is beneficial.

4. **Legacy Neural Architectures:** Historically favored in early neural networks for its ability to handle both positive and negative data distributions.

### 11.9.5 Limitations and Practical Considerations

Despite its advantages, tanh presents certain challenges in deep learning applications:

1. **Vanishing Gradient Problem:** Saturation at extreme input values leads to diminished gradients, impairing learning in deep networks.

2. **Slower Convergence:** Networks using tanh may train slower compared to ReLU-based architectures.

3. **Limited Suitability for Deep CNNs and Transformers:** In modern deep convolutional and Transformer models, activations like ReLU, GELU, or Swish are preferred due to better scaling properties.

4. **Higher Computational Cost Compared to Piecewise Activations:** Requires exponential operations per evaluation.

**11.9.6 Graphical Interpretation of tanh**



Figure 15: Graph of the tanh activation function, illustrating its smooth transition from $-1$ to 1.

As illustrated in Figure 15, the behavior of tanh can be interpreted as:

- **Negative Region ($x < 0$):** For negative inputs, $\tanh(x)$ smoothly approaches $-1$. As inputs become more negative, the function flattens, resulting in smaller gradients and weaker sensitivity to input variations.

- **Positive Region ($x \geq 0$):** For positive inputs, $\tanh(x)$ increases toward 1 with a similar flattening effect. Near the origin ($x \approx 0$), tanh behaves approximately linearly, providing strong sensitivity to small changes.

This characteristic "S-shaped" curve enables neural networks to suppress extreme input values while maintaining smooth, interpretable non-linearity, particularly valuable in recurrent learning settings.

## 11.10 Gaussian Error Linear Unit Approximation (GELU, `tanh`-based): Definition, Properties, and Practical Applications

### 11.10.1 Intuitive Overview

The Gaussian Error Linear Unit (GELU) activation function, when approximated using the hyperbolic tangent (tanh) function, provides a computationally efficient variant of the original GELU formulation based on the error function (erf). The tanh-based approximation introduces smooth non-linearity while maintaining a high-fidelity approximation of the probabilistic gating behavior of GELU, enabling deep networks to model complex patterns effectively with reduced computational overhead.

### 11.10.2 Mathematical Definition

The tanh-based approximation of GELU is defined as:

$$\text{GELU}_{\text{tanh}}(x) = 0.5x \left( 1 + \tanh \left( \sqrt{\frac{2}{\pi}} \left( x + 0.044715x^3 \right) \right) \right)$$

Here, the cubic term inside the tanh introduces mild curvature, aligning the approximation closely to the exact GELU curve while significantly accelerating computation.

### 11.10.3 Formulaic Pseudocode Implementation

```
1  Function GELU_tanh(Tensor) -> Tensor:
2      result = Empty_Tensor(Shape = Tensor.Shape)
3
4      for i from 0 to Tensor.Shape.x - 1:
5          for j from 0 to Tensor.Shape.y - 1:
6              x = Tensor[i][j]
7              inner = sqrt(2 / pi) * (x + 0.044715 * x**3)
8              result[i][j] = 0.5 * x * (1 + tanh(inner))
9
10     return result
```

Listing 28: Formulaic Pythonic Code for GELU (tanh) Approximation

As with all activation functions discussed, GELU(tanh) is applied element-wise across the tensor.

### 11.10.4 Advantages and Use Cases

GELU(tanh) retains many of the core advantages of GELU while offering additional practical benefits:

1. **Smooth Non-linearity:** Maintains smooth gradient flow critical for training deep networks.

2. **Efficient Computation:** Significantly faster to evaluate than the erf-based GELU, making it suitable for large-scale deployments.

3. **High-Fidelity Approximation:** Provides a close match to the probabilistic behavior of the original GELU formulation.

4. **Transformer Architectures:** Frequently employed in Transformer-based models such as BERT and GPT to balance computational cost with model expressivity.

### 11.10.5 Limitations and Practical Considerations

Despite its efficiency, GELU(tanh) entails certain trade-offs:

1. **Approximation Error:** Although close, GELU(tanh) is still an approximation and may introduce minor deviations from true probabilistic behavior.

2. **Higher Cost Compared to ReLU:** Still computationally heavier than simple piecewise activations like ReLU or Leaky ReLU.

3. **Potential Overfitting Risk:** As with the exact GELU, its fine-grained sensitivity to input variations can lead to overfitting on small or noisy datasets.

### 11.10.6 Graphical Interpretation of GELU (tanh)



Figure 16: Graph of the GELU activation function using the tanh-based approximation.

As illustrated in Figure 16, the behavior of GELU(tanh) can be interpreted as:

- **Negative Region ($x < 0$):** For negative inputs, the output softly decays toward zero without abrupt cutoffs, preserving minor negative activations. This enables richer feature learning compared to ReLU-like activations.

- **Positive Region ($x \geq 0$):** For positive inputs, the function behaves almost linearly for large magnitudes, closely tracking the identity function while maintaining smooth differentiability.

The combination of these behaviors contributes to enhanced learning stability, improved optimization dynamics, and superior generalization in deep architectures.

## 11.11   Softmax Activation: Definition, Properties, and Practical Applications

### 11.11.1   Intuitive Overview

The Softmax function transforms a vector of arbitrary real-valued scores into a probability distribution over mutually exclusive classes. By applying exponentiation and normalization, Softmax ensures that each output value is non-negative and that the outputs collectively sum to one. This allows raw model outputs (logits) to be interpreted as human-readable confidence scores, making Softmax a crucial component in multi-class classification tasks.

### 11.11.2   Step-by-Step Breakdown

The Softmax transformation can be broken into the following conceptual stages, applied to each row or vector independently:

1. **Stabilization (Optional):** To avoid numerical overflow during exponentiation, subtract the maximum value from each element in the vector.

2. **Exponentiation:** Apply the exponential function to each stabilized value, emphasizing larger inputs.

3. **Summation:** Compute the sum of all exponentiated values within the vector.

4. **Normalization:** Divide each exponentiated value by the local sum to produce a normalized probability distribution.

### 11.11.3   Mathematical Definition

Given an input vector $\mathbf{x} = (x_1, x_2, \ldots, x_n)$, the Softmax function outputs a vector $\mathbf{p} = (p_1, p_2, \ldots, p_n)$ where:

$$p_i = \frac{e^{x_i}}{\sum_{j=1}^{n} e^{x_j}} \quad \text{for each} \quad i \in \{1, 2, \ldots, n\}$$

When stabilization is applied, the formula becomes:

$$p_i = \frac{e^{x_i - \max(\mathbf{x})}}{\sum_{j=1}^{n} e^{x_j - \max(\mathbf{x})}}$$

which prevents overflow when $x_i$ values are large.

### 11.11.4   Formulaic Pseudocode Implementation

```
1  Function Softmax(Tensor) -> Tensor:
2      Result = []   # Final normalized tensor
3
4      for row in Tensor:
5          # Step 0: Stabilize
6          max_val = max(row)
7
8          # Step 1: Exponentiation
9          exps = []
10         for x in row:
11             exps.append(exp(x - max_val))
```

```
12
13        # Step 2: Summation
14        sum_exps = sum(exps)
15
16        # Step 3: Normalization
17        softmax_row = []
18        for value in exps:
19            softmax_row.append(value / sum_exps)
20
21        Result.append(softmax_row)
22
23    return Result
```

Listing 29: Formulaic Pythonic Code for Softmax Activation

Unlike prior element-wise activations (e.g., ReLU, tanh), Softmax is a **vector-wise** operation — its output for a given element depends on all elements within the same row or vector.

### 11.11.5 Advantages and Use Cases

Softmax offers several crucial advantages:

1. **Probabilistic Interpretation:** Outputs can be interpreted as class confidence scores.

2. **Compatibility with Cross-Entropy Loss:** When combined with cross-entropy loss, Softmax enables efficient optimization for multi-class classification.

3. **Human-Readable Outputs:** Converts complex logits into easily interpretable probability distributions.

4. **Multi-Class Decision Making:** Facilitates decisions between mutually exclusive outcomes.

### 11.11.6 Limitations and Practical Considerations

Despite its utility, Softmax is not universally applicable:

1. **Binary Classification:** For binary problems, the sigmoid activation is typically preferred over Softmax.

2. **Regression Tasks:** In regression problems predicting continuous values, Softmax is inappropriate.

3. **Hidden Layers:** Softmax is rarely used in hidden layers, as it constrains activations unnecessarily early in the network.

### 11.11.7 Application in Multi-Class Output Layers

When a model is tasked with choosing among multiple classes, applying Softmax at the output layer enables probabilistic interpretation:

- Each class receives a probability score between 0 and 1.

- All class scores collectively sum to 1.

- The class with the highest probability is typically selected as the predicted label.

This behavior contrasts with sigmoid activations, which are best suited for independent binary predictions rather than mutually exclusive classification tasks.

### 11.11.8 Graphical Interpretation of Softmax



Figure 17: Graph of Softmax output distributions across three example classes.

As illustrated in Figure 17:

- Each curve corresponds to a different class output.

- As one input becomes dominant, its associated probability sharply rises toward 1, while the probabilities of competing classes decline toward 0.

- The resulting probability vectors remain normalized, ensuring interpretability and consistency.

Conceptually, Softmax generalizes the behavior of the sigmoid function to multi-class settings, offering a smooth and differentiable mapping from real-valued scores to probabilistic decisions.

# 12 Derivatives of Activation Functions for Backpropagation

In supervised learning with neural networks, the backpropagation algorithm relies fundamentally on the derivatives of activation functions, applying the chain rule of calculus to compute gradients with respect to model parameters. The derivative of an activation function quantifies the sensitivity of a neuron's output to infinitesimal changes in its input, directly influencing the flow of learning signals during gradient descent optimization.

**Note:** In all cases presented, it is assumed that the variable $x$ refers to the **input** to the activation function during the forward pass, as required by the backpropagation procedure.

## 12.1 Derivative of the Rectified Linear Unit (ReLU)

The ReLU derivative is defined piecewise. During backpropagation:

- If the pre-activation input $x > 0$, the gradient flows unchanged (multiplied by 1).

- If $x \leq 0$, the gradient is zeroed out (multiplied by 0).

Formally:

$$\frac{d}{dx}\text{ReLU}(x) = \begin{cases} 1, & \text{if } x > 0 \\ 0, & \text{otherwise} \end{cases}$$

**Formulaic Pseudocode for ReLU Derivative:**

```
1  Function relu_derivative(x):
2      if x > 0:
3          return 1
4      else:
5          return 0
```

Listing 30: Formulaic Pythonic Code for ReLU Derivative

The behavior of the ReLU derivative is a major contributor to the "dying ReLU" phenomenon in deep networks, where neurons become permanently inactive.

## 12.2 Derivative of the Leaky ReLU

The Leaky ReLU derivative modifies the ReLU behavior by allowing a small, non-zero gradient when $x \leq 0$, controlled by a hyperparameter $\alpha$ (e.g., $\alpha = 0.01$). During backpropagation:

- If $x > 0$, the derivative is 1.

- If $x \leq 0$, the derivative is $\alpha$.

Formally:

$$\frac{d}{dx}\text{LeakyReLU}(x) = \begin{cases} 1, & \text{if } x > 0 \\ \alpha, & \text{otherwise} \end{cases}$$

**Formulaic Pseudocode for Leaky ReLU Derivative:**

```
1  Function leaky_relu_derivative(x, alpha):
2      if x > 0:
3          return 1
4      else:
5          return alpha
```

Listing 31: Formulaic Pythonic Code for Leaky ReLU Derivative

By preserving a small gradient for negative inputs, Leaky ReLU mitigates the vanishing gradient issues associated with standard ReLU, improving the robustness of learning dynamics, particularly in very deep or unbalanced networks.

## 12.3 Derivatives of Advanced Activation Functions

For effective training via backpropagation, the derivatives of activation functions are crucial. They govern the flow of gradients across layers and directly influence optimization dynamics. Below, we formalize the derivatives of Swish, Sigmoid, GELU (erf), and tanh activations.

**Note:** All derivatives assume $x$ refers to the input to the activation function unless otherwise specified.

### 12.3.1 Derivative of Swish Activation

The Swish activation, defined as $\text{Swish}(x) = x \cdot \sigma(x)$, has a derivative:

$$\frac{d}{dx}\text{Swish}(x) = \sigma(x) + x\sigma(x)(1 - \sigma(x))$$

where $\sigma(x)$ is the sigmoid function.

**Formulaic Pseudocode for Swish Derivative:**

```
1  Function swish_derivative(x):
2      s = sigmoid(x)
3      return s + x * s * (1 - s)
```

Listing 32: Formulaic Pythonic Code for Swish Derivative

Alternatively, if only computing $\sigma(x)$ is desired for efficiency:

```
1  Function sigmoid(x):
2      return 1 / (1 + exp(-x))
```

Listing 33: Sigmoid Helper Function for Swish Derivative

### 12.3.2 Derivative of Sigmoid Activation

The derivative of the Sigmoid activation, $\sigma(x) = \frac{1}{1+e^{-x}}$, is:

$$\frac{d}{dx}\sigma(x) = \sigma(x)(1 - \sigma(x))$$

where $\sigma(x)$ denotes the sigmoid output.

**Formulaic Pseudocode for Sigmoid Derivative:**

```
1  Function sigmoid_derivative_from_output(sigmoid_output):
2      return sigmoid_output * (1 - sigmoid_output)
```

Listing 34: Formulaic Pythonic Code for Sigmoid Derivative

This compact form allows reusing the forward-pass sigmoid result during backpropagation, improving computational efficiency.

### 12.3.3 Derivative of GELU (erf-based)

The GELU activation, defined via the Gaussian error function, has a derivative:

$$\frac{d}{dx}\text{GELU}(x) = 0.5\left(1 + \text{erf}\left(\frac{x}{\sqrt{2}}\right)\right) + \frac{xe^{-\frac{x^2}{2}}}{\sqrt{2\pi}}$$

where $\text{erf}(\cdot)$ is the error function.

**Formulaic Pseudocode for GELU (erf) Derivative:**

```
Function gelu_derivative(x):
    SQRT_2PI = sqrt(2 * pi)   # or use a standard constant if available
    erf_input = x / sqrt(2)

    derivative = 0.5 * (1 + erf(erf_input)) + (x * exp(-0.5 * x**2)) / SQRT_2PI
    return derivative
```

Listing 35: Formulaic Pythonic Code for GELU (erf) Derivative

The derivative captures both the smoothness and the probabilistic gating effect inherent in GELU activations.

### 12.3.4 Derivative of tanh Activation

The derivative of the hyperbolic tangent activation, $\tanh(x)$, is:

$$\frac{d}{dx}\tanh(x) = 1 - \tanh^2(x)$$

where $\tanh(x)$ is the forward activation output.

**Formulaic Pseudocode for tanh Derivative:**

```
Function tanh_derivative(tanh_output):
    return 1 - tanh_output**2
```

Listing 36: Formulaic Pythonic Code for tanh Derivative

This structure, similar to the sigmoid derivative, facilitates efficient computation during backpropagation by leveraging cached forward-pass outputs.

## 12.4 Derivative of GELU (tanh-based Approximation)

The `tanh`-based approximation of the Gaussian Error Linear Unit (GELU) is computationally efficient and widely used in Transformer-based architectures. Given the approximation:

$$\text{GELU}_{\text{tanh}}(x) = 0.5x \left( 1 + \tanh\left( \sqrt{\frac{2}{\pi}}(x + 0.044715x^3) \right) \right)$$

its derivative with respect to $x$ is:

$$\frac{d}{dx}\text{GELU}_{\text{tanh}}(x) = 0.5\left(1 + \tanh(u)\right) + 0.5x\left(1 - \tanh^2(u)\right) \cdot \frac{d}{dx}u \quad \text{where} \quad u = \sqrt{\frac{2}{\pi}}(x + 0.044715x^3)$$

This derivative smoothly modulates gradient flow and is more computationally tractable than the exact erf-based form.

**Formulaic Pseudocode for GELU(tanh) Derivative:**

```
Function gelu_tanh_derivative(x):
    # GELU approximation uses a cubic term inside tanh
    u = sqrt(2 / pi) * (x + 0.044715 * x**3)
    tanh_u = tanh(u)
```

```
5     sech_sq = 1 - tanh_u**2
6     du_dx = sqrt(2 / pi) * (1 + 3 * 0.044715 * x**2)
7
8     return 0.5 * (1 + tanh_u) + 0.5 * x * sech_sq * du_dx
```
Listing 37: Formulaic Pythonic Code for GELU (tanh) Derivative

This form ensures compatibility with auto-differentiation frameworks and provides smooth, stable gradients during backpropagation.

## 12.5  Derivative of Softmax (and Its Practical Simplification)

The Softmax function is defined as:

$$\text{Softmax}(x_i) = \frac{e^{x_i}}{\sum_j e^{x_j}}$$

Its derivative is formally given by the Jacobian matrix:

$$\frac{\partial \text{Softmax}(x_i)}{\partial x_j} = \begin{cases} \text{Softmax}(x_i)(1 - \text{Softmax}(x_i)), & \text{if } i = j \\ -\text{Softmax}(x_i) \cdot \text{Softmax}(x_j), & \text{if } i \neq j \end{cases}$$

This results in a dense Jacobian matrix that is seldom computed directly in practice.

**Practical Note:** In modern machine learning pipelines, Softmax is almost exclusively paired with cross-entropy loss. The gradient of their composition simplifies considerably:

$$\frac{\partial \mathcal{L}}{\partial x_i} = \text{Softmax}(x_i) - y_i$$

where $y_i$ is the true label (one-hot encoded). This expression avoids computing the full Jacobian and is numerically stable, making it standard in classification models.

Therefore, while the theoretical derivative of Softmax is well-defined, it is typically bypassed via this simplification in the context of supervised learning.

# 13  Loss Functions: Formal Definitions and Computational Structure

In supervised learning, loss functions quantify the divergence between a model's predictions and the true target outputs. During training, the optimization algorithm—typically a gradient-based method—uses these loss values to iteratively adjust the model's parameters with the goal of minimizing prediction error.

Loss is not to be confused with error itself. While error often refers to the raw difference between predicted and actual values, loss functions translate that discrepancy into a single scalar quantity optimized during learning.

**Note:** Unless otherwise specified, all functions assume that input tensors are of equal shape and appropriately batched.

## 13.1 Mean Absolute Error (MAE)

MAE computes the mean of the absolute differences between predicted and actual values:

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^{n} |y_i - \hat{y}_i|$$

**Formulaic Pseudocode for MAE:**

```
function mae_loss(actual: Tensor, predicted: Tensor) -> f32:
    diff = actual - predicted
    abs_diff = diff.abs()
    return abs_diff.mean()
```

Listing 38: Formulaic Pythonic Code for Mean Absolute Error (MAE)

## 13.2 Mean Squared Error (MSE)

MSE penalizes larger deviations more strongly than MAE by squaring the errors:

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{y}_i)^2$$

**Formulaic Pseudocode for MSE:**

```
function mse_loss(actual: Tensor, predicted: Tensor) -> f32:
    squared_error = (actual - predicted).powered(2)
    return squared_error.mean()
```

Listing 39: Formulaic Pythonic Code for Mean Squared Error (MSE)

## 13.3 Hinge Loss

Hinge loss is used primarily for binary classification with Support Vector Machines (SVMs). It encourages a decision boundary with a margin:

$$\text{HingeLoss} = \max(0, 1 - y_i \cdot \hat{y}_i) \quad \text{where } y_i \in \{-1, 1\}$$

**Formulaic Pseudocode for Hinge Loss:**

```
Function HingeLoss(predicted, actual):
    return max(0, 1 - predicted * actual)
```

Listing 40: Formulaic Pythonic Code for Hinge Loss

## 13.4 Huber Loss

Huber loss combines the advantages of MAE and MSE. For small errors, it behaves quadratically (like MSE), while for large errors, it behaves linearly (like MAE), thus improving robustness to outliers:

$$\text{Huber}_\delta(y, \hat{y}) = \begin{cases} \frac{1}{2}(y - \hat{y})^2, & \text{if } |y - \hat{y}| \leq \delta \\ \delta(|y - \hat{y}| - \frac{1}{2}\delta), & \text{otherwise} \end{cases}$$

**Formulaic Pseudocode for Huber Loss:**

```
1  Function HuberLoss(predicted, actual, delta):
2      error = predicted - actual
3      if abs(error) <= delta:
4          return 0.5 * error * error
5      else:
6          return delta * (abs(error) - 0.5 * delta)
```

Listing 41: Formulaic Pythonic Code for Huber Loss

## 13.5 Binary Cross-Entropy (BCE)

BCE is commonly used in binary classification. It quantifies the divergence between predicted probabilities and actual binary labels:

$$\text{BCE} = - \left[ y \cdot \log(\hat{y}) + (1 - y) \cdot \log(1 - \hat{y}) \right]$$

To avoid logarithmic instability, predictions are typically clamped using a small constant $\varepsilon$.

**Formulaic Pseudocode for BCE:**

```
1  Function BinaryCrossEntropy(predicted, actual):
2      predicted = clamp(predicted, epsilon, 1 - epsilon)
3      return - (actual * log(predicted) + (1 - actual) * log(1 - predicted))
```

Listing 42: Formulaic Pythonic Code for Binary Cross-Entropy (BCE)

## 13.6 Categorical Cross-Entropy (CCE)

CCE generalizes BCE to multi-class settings. It penalizes incorrect class probabilities based on a one-hot encoded label:

$$\text{CCE} = - \log(p_y) \quad \text{where } p_y \text{ is the predicted probability for the correct class}$$

**Formulaic Pseudocode for CCE:**

```
1  Function CategoricalCrossEntropy(predicted_probs, actual_label):
2      predicted_probs = clamp_each(predicted_probs, epsilon, 1 - epsilon)
3      return -log(predicted_probs[actual_label])
```

Listing 43: Formulaic Pythonic Code for Categorical Cross-Entropy (CCE)

| Loss Function | Key Behavior |
|---|---|
| Mean Absolute Error (MAE) | Treats all errors equally by measuring the absolute difference between actual and predicted values. |
| Mean Squared Error (MSE) | Penalizes larger errors more severely by squaring them, making the model more sensitive to outliers. |
| Hinge Loss | Builds a margin of safety around classes, encouraging confident classification boundaries (commonly used in Support Vector Machines). |
| Huber Loss | Smoothly combines MAE and MSE behavior: acts like MSE for small errors and MAE for large errors, making it robust to outliers. |
| Binary Cross Entropy (BCE) | Measures the difference between two probability distributions in binary classification tasks (two classes: 0 or 1). |
| Categorical Cross Entropy (CCE) | Extends Cross Entropy to multi-class classification problems, comparing predicted probability distributions over multiple categories. |

Table 4: Summary of Common Loss Functions and Their Behavior

# 14 Optimizers: Theory and Role in Training Dynamics

In machine learning, particularly in gradient-based neural network training, an **optimizer** is an algorithm that governs how a model's parameters are updated during learning. At the core of this process lies the concept of the *learning rate*—a hyperparameter that controls the size of the updates applied to model weights in response to the loss gradient.

Optimizers aim to minimize the loss function by iteratively adjusting the model's parameters in the direction of steepest descent (or approximations thereof). While basic methods apply a fixed learning rate and simple updates, more advanced optimizers adaptively modify the learning rate over time, incorporate past gradients, and regularize the updates to improve convergence speed and stability.

The choice of optimizer plays a critical role in determining:

- The speed of convergence toward a solution.

- The stability of training across varying data distributions.

- The model's ability to escape local minima or saddle points.

- Generalization performance on unseen data.

This section presents both classical and modern optimization algorithms, each accompanied by its update rule, key characteristics, and usage considerations.

Traditional optimizers such as Stochastic Gradient Descent (SGD) update parameters by stepping proportionally to the negative gradient. However, more sophisticated optimizers like Adam introduce adaptive learning rates and momentum, greatly improving convergence speed and stability."

## 14.1 Adam Optimizer (Adaptive Moment Estimation)

The Adam optimizer is a first-order gradient-based optimization algorithm that adaptively adjusts the learning rate for each parameter based on estimates of the first and second moments of the gradients.

| Aspect | Adam Optimizer | SGD Optimizer |
|---|---|---|
| Learning Rate Adjustment | Adaptive learning rate for each parameter | Fixed or manually decayed learning rate |
| Speed of Convergence | Faster convergence, especially early on | Slower convergence without momentum |
| Memory Usage | Higher (stores moment estimates) | Lower (only stores gradients) |
| Sensitivity to Learning Rate | Less sensitive to initial learning rate choice | Highly sensitive to learning rate choice |
| Use Case | Good for complex networks, sparse gradients | Good for simple or very large datasets |
| Mathematical Complexity | More complex (uses moving averages of gradients) | Simpler (straightforward gradient update) |

Table 5: Comparison of Adam and SGD Optimizers

Introduced by Kingma and Ba (2015), Adam merges the benefits of Momentum (smoothing using first moments) and RMSProp (normalization using second moments) into a single, robust method.

Formally, Adam maintains two exponentially decaying averages:

- $m_t$ – the exponentially weighted moving average of past gradients (first moment estimate)

- $v_t$ – the exponentially weighted moving average of squared gradients (second moment estimate)

To correct bias introduced during initialization, bias-corrected estimates are computed:

$$\hat{m}_t = \frac{m_t}{1 - \beta_1^t}, \quad \hat{v}_t = \frac{v_t}{1 - \beta_2^t}$$

The parameter update rule becomes:

$$\theta_{t+1} = \theta_t - \eta \cdot \frac{\hat{m}_t}{\sqrt{\hat{v}_t} + \varepsilon}$$

where:

- $\theta_t$ – current parameter value

- $\eta$ – learning rate (default: 0.001)

- $\varepsilon$ – small constant for numerical stability (e.g., $10^{-8}$)

- $\beta_1, \beta_2$ – decay rates for the first and second moment estimates (typically $\beta_1 = 0.9$, $\beta_2 = 0.999$)

**Formulaic Pythonic Pseudocode for Adam Optimizer:**

```
Function adam_optimizer(Tensor, m, v, learning_rate, timestep):
    beta1 = 0.9
    beta2 = 0.999
    epsilon = 1e-8

    for i in range(Tensor.rows):
        for j in range(Tensor.cols):
            grad = Tensor[i][j]

```

```
10        # Update biased first moment estimate
11        m[i][j] = beta1 * m[i][j] + (1 - beta1) * grad
12
13        # Update biased second raw moment estimate
14        v[i][j] = beta2 * v[i][j] + (1 - beta2) * (grad ** 2)
15
16        # Compute bias-corrected moment estimates
17        m_hat = m[i][j] / (1 - beta1 ** timestep)
18        v_hat = v[i][j] / (1 - beta2 ** timestep)
19
20        # Update parameter
21        Tensor[i][j] -= learning_rate * m_hat / (sqrt(v_hat) + epsilon)
```

Listing 44: Formulaic Pythonic Code for Adam Optimizer

**Advantages of Adam**

- Adaptive learning rate per parameter

- Requires minimal tuning and is robust to sparse gradients

- Combines fast convergence with stability in non-stationary settings

- Well-suited for large models and noisy objective functions

**Limitations and Practical Notes**

- Adam may sometimes lead to overfitting or poor generalization compared to SGD with momentum in certain vision tasks.

- Learning rate warm-up, weight decay (AdamW), or switching to SGD mid-training are often employed for better results.

**Commentary:** While basic stochastic gradient descent (SGD) is inherently integrated into the backpropagation update rule, optimizers like Adam offer substantial performance gains in training time and convergence speed—particularly in deep and high-dimensional architectures.

## 14.2   Stochastic Gradient Descent (SGD)

Stochastic Gradient Descent (SGD) is the most fundamental optimization algorithm in machine learning and serves as the theoretical basis for more advanced methods. SGD seeks to minimize a loss function $\mathcal{L}(\theta)$ by iteratively updating model parameters $\theta$ in the direction of the negative gradient.

Unlike full-batch gradient descent, which computes gradients over the entire training dataset at each step, SGD performs updates using a single sample (or mini-batch), defined as:

$$\theta_{t+1} = \theta_t - \eta \cdot \nabla_\theta \mathcal{L}(\theta_t; x_i, y_i)$$

where:

- $\eta$ – learning rate

- $(x_i, y_i)$ – the $i^{th}$ training example or mini-batch

- $\nabla_\theta \mathcal{L}$ – gradient of the loss with respect to model parameters

By introducing stochasticity through random sampling, SGD helps models escape shallow local minima and saddle points, particularly in high-dimensional, non-convex optimization landscapes.

**Formulaic Pseudocode for SGD:**

```
1  for param in parameters:
2      param.value -= learning_rate * param.gradient
```

Listing 45: Formulaic Pythonic Code for SGD Optimizer

**Advantages of SGD**

- Simple to implement and computationally efficient.

- Well-suited for very large datasets (via mini-batch training).

- Can generalize well when combined with regularization and proper learning rate scheduling.

**Limitations and Practical Notes**

- Convergence may be slow and noisy without momentum or learning rate decay.

- Requires careful tuning of the learning rate to avoid divergence or stagnation.

- Prone to getting stuck in sharp local minima or saddle points in complex loss surfaces.

**Note:** SGD is essentially a direct implementation of the backpropagation update rule and is thus already integrated into the core mechanics of neural network training. Its simplicity makes it a standard benchmark and starting point for optimizer comparisons.

# 15   Normalization: Mitigating Gradient Instabilities in Deep Learning

One of the primary challenges in training deep neural networks is the instability of gradients—specifically, the phenomena of *vanishing gradients* and *exploding gradients*. These arise when the magnitude of gradients becomes excessively small or large, causing either stalled learning or numerical overflow, particularly in networks with many layers or poorly scaled inputs.

Normalization techniques are designed to mitigate these issues by rescaling data into numerically manageable ranges. This preserves essential distributional characteristics while reducing the risk of pathological gradients during backpropagation.

Two of the most common normalization strategies are:

- **Min-Max Normalization** rescales input features to a predefined range, usually $[0, 1]$, preserving relative proportions.

- **Z-Score Standardization (Standard Scaling)** centers the data around zero and scales it to unit variance, promoting symmetry in gradient flow.

**Comparison of Normalization Methods:**

| Feature | Min-Max Normalization | Z-Score Standardization |
|---|---|---|
| Range after scaling | Typically $[0,1]$; customizable bounds | Centered around 0 with unit variance |
| Sensitive to outliers? | Highly sensitive | More robust due to variance scaling |
| Use case examples | Pixel intensities, bounded sensors | ML models assuming Gaussian features |
| Effect on distribution | Linear rescaling; preserves original shape | Normalizes spread and location |
| When to use? | When known min/max bounds exist | When data is approximately normal |

Table 6: Comparison of Min-Max Normalization and Z-Score Standardization

**Practical Note:** Normalization is not strictly required for all models, but it plays a pivotal role in stabilizing gradient-based training processes. For networks prone to exploding/vanishing gradients, appropriate normalization can dramatically improve convergence rates and model robustness.

## 15.1 Z-Score Normalization (Standardization)

Z-score normalization, also referred to as *standardization*, transforms data such that the resulting distribution has zero mean and unit variance. This technique is particularly effective when the underlying data distribution is approximately Gaussian.

Given a tensor $X$ with elements $x_i$, the standardized form is:

$$x_i^{\text{norm}} = \frac{x_i - \mu}{\sigma} \quad \text{where} \quad \mu = \frac{1}{n}\sum_{i=1}^{n} x_i, \quad \sigma = \sqrt{\frac{1}{n-1}\sum_{i=1}^{n}(x_i - \mu)^2}$$

**Formulaic Pythonic Pseudocode for Standardization:**

```
function normalize(Tensor) -> Tensor:
    mean = Tensor.mean()
    std = Tensor.std(sample=True)  # Use sample std (n-1 denominator)
    return (Tensor.clone() - mean) / std
```

Listing 46: Formulaic Pythonic Code for Z-Score Normalization

**Use Cases**

- Effective when model assumptions include normally distributed features.

- Common in linear regression, logistic regression, SVMs, and neural networks.

## 15.2 Min-Max Normalization

Min-Max normalization rescales each element in the dataset to a fixed range, typically $[0,1]$, by subtracting the minimum and dividing by the range. It preserves the shape of the original distribution and is highly interpretable but sensitive to outliers.

The transformation is defined as:

$$x_i^{\text{norm}} = \frac{x_i - \min(X)}{\max(X) - \min(X)}$$

**Formulaic Pythonic Pseudocode for Min-Max Normalization:**

```
function min_max_normalize(Tensor):
    min_value = find_minimum(Tensor)
    max_value = find_maximum(Tensor)

    for i in range(Tensor.rows):
        for j in range(Tensor.cols):
            Tensor[i][j] = (Tensor[i][j] - min_value) / (max_value - min_value)
```

Listing 47: Formulaic Pythonic Code for Min-Max Normalization

**Use Cases**

- Image pixel scaling (e.g., from $[0, 255]$ to $[0, 1]$).

- Features with known and bounded value ranges.

**Practical Note:** Min-max normalization assumes a stable range. In the presence of outliers or shifting distributions (e.g., in online learning), Z-score standardization or robust scaling may be preferred.

# 16 Tensor/Weights and Biases Initialization

Proper initialization of weights and biases plays a pivotal role in stabilizing training dynamics and promoting efficient convergence in deep learning models. In particular, poor initialization can cause gradients to vanish or explode as they are propagated through layers—issues that severely hinder or prevent learning.

To address this, various initialization schemes have been developed to align the variance of activations and gradients across layers, depending on the choice of activation function and model architecture.

**General Principles of Initialization**

- **Variance Preservation:** Maintaining activation variance across layers prevents gradient shrinkage or amplification.

- **Activation Function Dependency:** Initialization should match the non-linearity used (e.g., ReLU vs. tanh).

- **Symmetry Breaking:** All weights must be initialized independently to allow neurons to learn distinct features. Biases may be initialized to zero.

The following table summarizes widely adopted initialization strategies:

**Practical Note**

Weight initialization should always be selected in tandem with the activation function and network architecture. For deep networks using ReLU or GELU activations, He initialization is typically preferred. For sigmoid or tanh-based networks, Xavier initialization is more appropriate.

Biases, unless specifically tuned (e.g., for LSTMs or BatchNorm), are typically initialized to zero.

## 16.1   Zero Initialization

Zero initialization sets all weights in a neural network layer to zero. While this method is computationally efficient and trivial to implement, it introduces a fundamental flaw: all neurons in a given layer receive identical gradients during backpropagation and thus learn the same features. This symmetry prevents the network from effectively breaking degeneracy in weight space.

**Key Drawbacks:**

- Fails to break symmetry between neurons.

- Leads to vanishing gradients and stalled learning.

- Only suitable for initializing biases.

## 16.2   Pseudo-Random Initialization (Uniform)

Random uniform initialization generates weights using a uniform distribution within a specified range (e.g., $[-0.5, 0.5]$). While it avoids the symmetry problem seen in zero initialization, it does not scale based on network depth or activation function, and thus may lead to unstable training in deep architectures.

**Formulaic Pythonic Pseudocode:**

```
function Random_Tensor(shape, min_val, max_val) -> Tensor:
    Tensor = []
    for i in range(shape.x):
        row = []
        for j in range(shape.y):
            row.append(Pseudo_Random_Range_float32(min_val, max_val))
        Tensor.append(row)
    return Tensor
```

Listing 48: Random Uniform Tensor Initialization

**Note:** For uniform initialization, a common range is $[-0.5, 0.5]$. This should be adjusted depending on layer size and depth.

## 16.3   Xavier (Glorot) Initialization

Xavier (or Glorot) initialization is designed to preserve the variance of activations across layers when using `tanh` or `sigmoid` activations. It draws weights from a uniform distribution with bounds based on the number of input and output units.

$$\text{bound} = \frac{\sqrt{6}}{\sqrt{n_{\text{in}} + n_{\text{out}}}} \quad \Rightarrow \quad W \sim \mathcal{U}(-\text{bound}, +\text{bound})$$

**Formulaic Pythonic Pseudocode:**

```
1  fan_in = number_of_input_units
2  fan_out = number_of_output_units
3
4  for weight in layer:
5      bound = sqrt(6) / sqrt(fan_in + fan_out)
6      weight = Pseudo_Random_Range_float32(-bound, +bound)
```

Listing 49: Formulaic Pythonic Code for Xavier Initialization

**Best suited for:** Activation functions such as `sigmoid`, `tanh`

## 16.4   He Initialization

He initialization (also called Kaiming initialization) is tailored for layers using ReLU or Leaky ReLU activations. It initializes weights by sampling from a normal distribution with zero mean and a standard deviation of:

$$\text{std\_dev} = \sqrt{\frac{2}{n_{\text{in}}}} \quad \Rightarrow \quad W \sim \mathcal{N}(0, \text{std\_dev}^2)$$

**Formulaic Pythonic Pseudocode:**

```
1  number_of_inputs = number of input connections to this layer
2
3  for weight in layer:
4      weight = Random_Normal(mean=0.0, std_dev=sqrt(2 / number_of_inputs))
```

Listing 50: Formulaic Pythonic Code for He Initialization

**Best suited for:** ReLU, Leaky ReLU, ELU, GELU
**Advantages:**

- Helps avoid vanishing gradients.

- Promotes stable signal propagation in deep networks.

- Generally converges faster in deep architectures with ReLU-based activations.

| Initialization Method | Key Characteristics and Use Cases |
|---|---|
| **Xavier (Glorot)** | • Balances the variance of activations between input and output layers.<br><br>• Designed for use with `tanh` and `sigmoid` activations.<br><br>• Helps mitigate vanishing/exploding gradients in shallow to moderately deep networks.<br><br>• Draws from a uniform or normal distribution with variance $\frac{2}{n_{\text{in}}+n_{\text{out}}}$.<br><br>• Not recommended for ReLU-family activations due to underscaling. |
| **He (Kaiming)** | • Tailored for ReLU and Leaky ReLU activations.<br><br>• Preserves forward-pass variance with variance $\frac{2}{n_{\text{in}}}$.<br><br>• Enables stable training of very deep networks by avoiding vanishing gradients.<br><br>• Can lead to exploding gradients if incorrectly scaled or used with incompatible activations. |
| **Zero Initialization** | • Initializes all weights to zero, creating symmetry where all neurons learn identical features.<br><br>• Disables gradient flow in hidden layers—rendering the network non-functional.<br><br>• Acceptable only for biases; strictly discouraged for weights.<br><br>• Often used as a baseline to demonstrate the necessity of variance. |
| **Random Uniform** | • Weights are drawn from a uniform distribution without scale normalization.<br><br>• Simple and historically common, but problematic in deep architectures.<br><br>• Offers baseline functionality in shallow networks or in reinforcement learning policies.<br><br>• May lead to slow convergence or instability due to unregulated variance. |

Table 7: Summary of Common Weight Initialization Methods and Their Core Characteristics

# 17    Backpropagation: The Core Learning Algorithm in Neural Networks

Backpropagation is the central algorithm responsible for training neural networks via supervised learning. It refines a model's internal parameters—namely, weights and biases—by propagating the error signal backward through the network, layer by layer, and computing gradients that guide each parameter update.

This process is built upon the chain rule of calculus, which allows the computation of how the loss changes with respect to every parameter in the network, even deep within multiple layers. When paired with gradient-based optimization algorithms such as Stochastic Gradient Descent (SGD) or Adam, backpropagation enables neural networks to approximate complex, nonlinear functions from data.

The training process typically consists of four key stages:

1. **Forward Pass** – Computes activations layer by layer using current weights and biases.

2. **Loss Function Evaluation** – Measures the discrepancy between the model's prediction and the true target.

3. **Backward Pass** – Propagates the loss gradient backward through the network, computing parameter-wise gradients.

4. **Gradient Descent Step** – Updates weights and biases using gradients and a selected optimization algorithm.

## 17.1    Key Computational Tools for Backpropagation

To implement backpropagation from first principles, the following computational primitives are essential:

- **Matrix Multiplication** – Enables dot products between weights and activations.

- **Transposition** – Required for aligning gradients across layer interfaces.

- **Element-wise Operations** – Applied to activation outputs and error signals.

- **Activation Functions and Their Derivatives** – Nonlinearities such as ReLU, tanh, or sigmoid must support differentiation.

## 17.2    Forward Pass

The forward pass processes inputs sequentially through the network's layers. Each layer computes:

$$z^{(l)} = W^{(l)}a^{(l-1)} + b^{(l)}, \quad a^{(l)} = f(z^{(l)})$$

Where:

- $z^{(l)}$ is the pre-activation output of layer $l$,

- $W^{(l)}$ and $b^{(l)}$ are the weight matrix and bias vector,

- $f(\cdot)$ is the activation function (e.g., ReLU, tanh),

- $a^{(l)}$ is the post-activation output, passed to the next layer.

This process continues through all hidden layers and concludes at the output layer.

## 17.3 Loss Function

The loss function quantifies prediction error. Common examples include:

$$\text{MSE} = \frac{1}{n}\sum_{i=1}^{n}(\hat{y}_i - y_i)^2, \quad \text{CrossEntropy} = -\sum_{i}y_i\log(\hat{y}_i)$$

The choice of loss function depends on the task (e.g., regression vs. classification) and impacts gradient computation.

## 17.4 Backward Pass (Backpropagation)

The backward pass computes the gradient of the loss function with respect to each weight and bias in the network by traversing layers in reverse. This relies on repeated application of the chain rule:

$$\frac{\partial \mathcal{L}}{\partial W^{(l)}} = \delta^{(l)}(a^{(l-1)})^{\top} \quad \text{where} \quad \delta^{(l)} = \left(\frac{\partial \mathcal{L}}{\partial a^{(l)}}\right) \odot f'(z^{(l)})$$

**Steps:**

1. Compute error gradient at the output layer.

2. Propagate this error backwards through hidden layers.

3. Multiply by activation function derivative.

4. Use the gradients to update parameters.

**Formulaic Pythonic Pseudocode:**

```
for neuron in layer:
    error_gradient = 2.0 * (outputs[neuron] - targets[neuron])   # MSE derivative
    activation_gradient = activation_derivative(neuron)
    total_gradient = error_gradient * activation_gradient

    for i in range(len(weights[neuron])):
        weights[neuron][i] -= learning_rate * total_gradient * inputs[i]

    biases[neuron] -= learning_rate * total_gradient
```

Listing 51: Formulaic Pythonic Code for Backpropagation

## 17.5 Gradient Descent Step

Once gradients are computed, they are used to update model parameters:

$$\theta^{(l)} \leftarrow \theta^{(l)} - \eta \cdot \nabla \mathcal{L}$$

This step may be implemented using SGD or more sophisticated optimizers (e.g., Adam, RMSProp) discussed earlier.
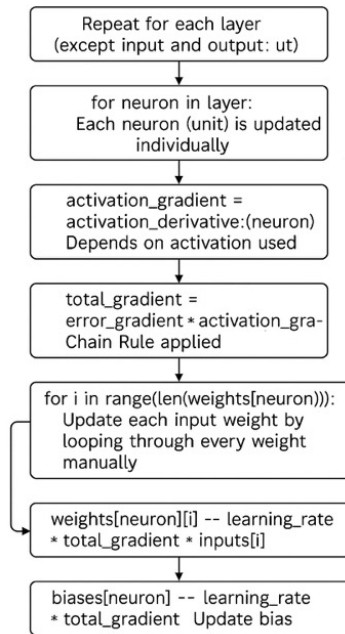
Figure 18: Breakdown of the backward pass in formulaic pseudocode.

## 17.6 Batch Training and Epochs

Training is typically performed in mini-batches for computational efficiency and generalization. A full dataset pass is called an **epoch**. The backpropagation cycle is repeated across batches and epochs to progressively reduce the loss.

## 17.7 Summary: The Backpropagation Pipeline

**Backpropagation** is the engine that drives learning in neural networks. It enables the network to systematically refine its parameters to minimize a loss function, translating gradient signals into corrective updates. The algorithm's structure, forward pass, loss computation, backward pass, and parameter update, forms the training loop at the heart of nearly all modern deep learning architectures.

Mastering backpropagation offers both a practical and theoretical foundation for understanding how complex models generalize, optimize, and ultimately learn to represent abstract patterns in data.

# Recent Advances

In recent years, substantial progress has been made in machine learning techniques, optimization methods, and model architectures. Innovations such as the Adam optimizer [1] and its improved variant AdamW [2] have reshaped training dynamics. Transformer architectures [3] introduced self-attention mechanisms that are now foundational across domains. Improvements in activation functions, including GELU [4] and Swish [5], have enhanced non-linearity handling. Regularization strategies like Cutout [6], Mixup [7], and RICAP [8] have strengthened model robustness.

New approaches to optimization such as RAdam [9], Lookahead [10], and AdaBelief [11] have emerged, refining convergence stability. Advances in normalization techniques like LayerNorm [12], GroupNorm [13],

and EvoNorm [14] have further improved the training of deep models. Weight initialization strategies like He initialization [15] and Fixup initialization [16] addressed challenges in very deep networks.

Recent network architectures, such as EfficientNet [17], ResNeSt [18], and ConvNeXt [19], have pushed the frontier of computer vision benchmarks. Self-supervised learning techniques like SimCLR [20] and BYOL [21] have unlocked new paradigms for representation learning without extensive labels.

These developments reflect an active and rapidly evolving field, providing a strong foundation for further innovation and accessibility in machine learning education.

# 18   Conclusion

This work has aimed to demystify and recontextualize the mathematical foundations of Machine Learning (ML) in order to reduce the barriers that often exclude highly capable developers, particularly those without formal training in calculus or theoretical mathematics. By translating core ML concepts into formulaic pseudocode and structured algorithmic logic, we advocate for a more inclusive framework, one that recognizes programming expertise as an equally valid and powerful entry point into the field.

Many developers today possess extensive knowledge in low-level systems engineering, embedded computing, Linux kernel development, and WebAssembly. These skillsets are often disconnected from ML discourse due to the pervasive assumption that mastery in advanced calculus is a prerequisite for meaningful participation. This dual requirement, proficiency in both abstract mathematics and complex coding, creates what may be termed the **double-genius barrier**, which inhibits broader innovation.

By reducing this barrier through pedagogical reform and accessible explanations, we can accelerate cross-disciplinary collaboration and invite contributions from a much wider talent pool. In doing so, the Machine Learning community stands to benefit from a more diverse, global influx of ideas, tools, and methodologies.

## 18.1   Guidelines for Increasing Accessibility in Machine Learning Education

To support this vision, the following principles are proposed:

- **Avoid exclusive reliance on calculus-based notation.** While mathematically rigorous exposition remains vital for theoretical development, educational materials should accompany such expressions with structured pseudocode or executable logic wherever possible.

- **Recognize programming fluency as foundational.** ML is ultimately implemented through code. Prioritizing algorithmic clarity ensures that practitioners from software engineering and systems programming domains can meaningfully engage with and apply ML principles.

- **Promote hybrid explanations.** Ideal instructional content should bridge formal mathematics and operational intuition—using diagrams, modular code, and relatable metaphors.

Such approaches will empower self-taught programmers, bootcamp graduates, and systems engineers alike, unlocking contributions from communities historically sidelined by overly theoretical teaching methods.

## 18.2   Machine Learning Without Walls: A Vision for Global Innovation

Imagine a future in which:

- Mathematicians focus exclusively on developing faster and more biologically realistic activation functions, such as novel generalizations of the Hodgkin–Huxley model.

- Systems engineers embed neural inference directly into operating system kernels to achieve ultra-low-latency pattern recognition.

- Educators adopt accessible frameworks for teaching ML logic, allowing students to implement and experiment before ever seeing formal integrals.

Such a future is only possible if we stop requiring every contributor to split their energy equally between two orthogonal skill sets. Let mathematicians build better theory. Let programmers build faster, safer systems. Let educators build bridges between them. In combination, not in isolation, these communities can accelerate the future of Machine Learning far beyond what any single discipline could accomplish alone.

*To the mathematicians: You do not have to carry the future of Machine Learning alone.*
*To the programmers: Your intuition, structure, and systems thinking are just as vital.*
*Together, we can dismantle the barriers that no one discipline can overcome in isolation.*

*To mathematicians: you don't have to carry the future of Machine Learning alone. To programmers: your skills are just as vital. Together, we can break barriers that no single discipline could ever overcome.*

## 18.3   A Complete Machine Learning Network Example

To demonstrate how the foundational principles discussed in this paper can be applied in practice, Figure 19 presents a complete example of a machine learning pipeline. This network is intentionally designed with transparency and educational clarity, enabling learners to follow each computational step.
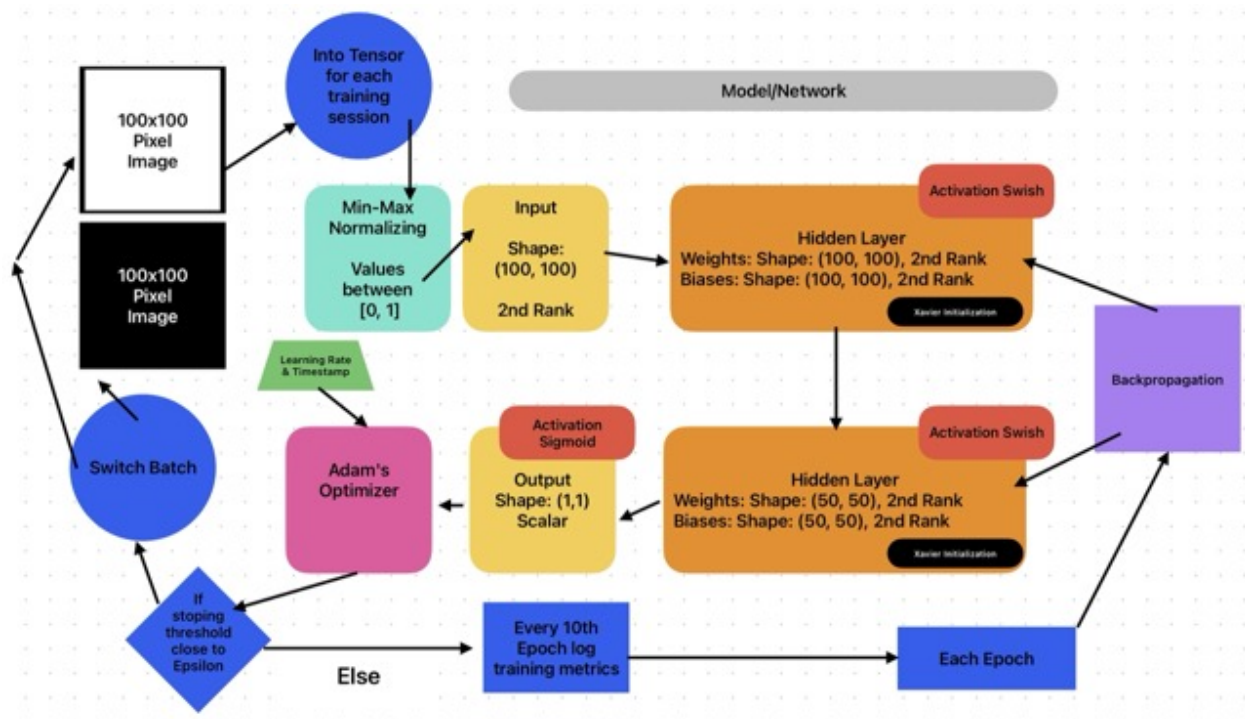
Figure 19: Overview of a complete image-based machine learning workflow. The process begins with preprocessing and normalization of 100×100 pixel grayscale images, which are then converted into tensors for input into the model. The network includes two hidden layers with Swish activations and Xavier weight initialization, followed by a sigmoid output layer. Adam's Optimizer guides training via backpropagation, iterating across epochs with periodic logging and batch switching based on convergence thresholds.

```
1  // 1. Preprocessing
2  For each image:
3      Normalize pixel values to [0, 1]
4
5  // 2. Initialization
6  Initialize weights and biases for each layer with Xavier Initialization
7
8  // 3. Training Loop
9  For each epoch:
10     For each image in batch:
11         // Forward Pass
12         Input = image tensor
13         Hidden1 = Swish(Weights1 x Input + Biases1)
14         Hidden2 = Swish(Weights2 x Hidden1 + Biases2)
15         Output = Sigmoid(Weights3 x Hidden2 + Bias3)
16
17         // Loss Calculation
18         Loss = (Output - Target).raised2
19
20         // Backward Pass (Backpropagation)
21         Compute gradients of Loss w.r.t Weights and Biases
22         Update Weights and Biases using Adam Optimizer
23
24     // Logging
```

```
25    If epoch mod 10 == 0:
26        Log training metrics
27
28    // Stopping Condition
29    If Loss < Epsilon:
30        Break training loop
```

Listing 52: Formulaic Pythonic Code for Example Complete machine learning network flow.

## 19   Final Reflection: The Power of Formulaic Thinking

To conclude, this example shows that the entire machine learning process can be captured cleanly through formulaic pseudocode, without hiding behind walls of pure math. From image preprocessing to network layer transitions, activations, optimization, and backpropagation, every step remains clear and accessible.

While mathematical notation like:

$$h_1 = \text{Swish}(W_1 x + b_1) \quad \text{and} \quad L = (\hat{y} - y)^2$$

is powerful, it often skips the operational thinking that real-world engineering demands.

Formulaic pseudocode, however, preserves the logic path step-by-step, making it easier to bridge understanding between mathematics and practical implementation.

This bridge is crucial.

By empowering both mathematicians and programmers to see through each other's worlds more clearly, we can push forward a future where neural networks are not only more powerful but also more understandable and accessible to all learners.

## License & Attribution

# References

[1]  Diederik P. Kingma and Jimmy Ba, *Adam: A Method for Stochastic Optimization*, ICLR, 2015.

[2]  Ilya Loshchilov and Frank Hutter, *Decoupled Weight Decay Regularization*, ICLR, 2019.

[3]  Ashish Vaswani et al., *Attention Is All You Need*, NeurIPS, 2017.

[4]  Dan Hendrycks and Kevin Gimpel, *Gaussian Error Linear Units (GELUs)*, arXiv:1606.08415, 2016.

[5]  Prajit Ramachandran et al., *Searching for Activation Functions*, arXiv:1710.05941, 2017.

[6]  Terrance DeVries and Graham W. Taylor, *Improved Regularization of Convolutional Neural Networks with Cutout*, arXiv:1708.04552, 2017.

[7]  Hongyi Zhang et al., *mixup: Beyond Empirical Risk Minimization*, ICLR, 2018.

[8]  Ryo Takahashi, Takashi Matsubara, Koichiro Yoshino, *RICAP: Random Image Cropping and Patching Data Augmentation for Deep CNNs*, ACML, 2018.

[9]  Liyuan Liu et al., *On the Variance of the Adaptive Learning Rate and Beyond*, arXiv:1908.03265, 2019.

[10]  Michael R. Zhang et al., *Lookahead Optimizer: k steps forward, 1 step back*, NeurIPS, 2019.

[11]  Juntang Zhuang et al., *AdaBelief Optimizer: Adapting Stepsizes by the Belief in Observed Gradients*, NeurIPS, 2020.

[12]  Jimmy Ba, Jamie Ryan Kiros, Geoffrey Hinton, *Layer Normalization*, arXiv:1607.06450, 2016.

[13]  Yuxin Wu and Kaiming He, *Group Normalization*, ECCV, 2018.

[14]  Zhuang Liu et al., *EvoNorm: Evolutionary Normalization for Deep Learning*, arXiv:2004.02967, 2020.

[15]  Kaiming He et al., *Delving Deep into Rectifiers: Surpassing Human-Level Performance*, ICCV, 2015.

[16]  Hongyi Zhang et al., *Fixup Initialization: Residual Learning Without Normalization*, ICLR, 2019.

[17]  Mingxing Tan and Quoc V. Le, *EfficientNet: Rethinking Model Scaling*, ICML, 2019.

[18]  Hang Zhang et al., *ResNeSt: Split-Attention Networks*, arXiv:2004.08955, 2020.

[19]  Zhuang Liu et al., *A ConvNet for the 2020s*, CVPR, 2022.

[20]  Ting Chen et al., *A Simple Framework for Contrastive Learning*, ICML, 2020.

[21]  Jean-Bastien Grill et al., *Bootstrap Your Own Latent (BYOL)*, NeurIPS, 2020.