

Logistic Regression vs. Perceptron

Geometric Objective: Both algorithms aim to find a hyperplane that separates two classes in the feature space.

Logistic Regression:

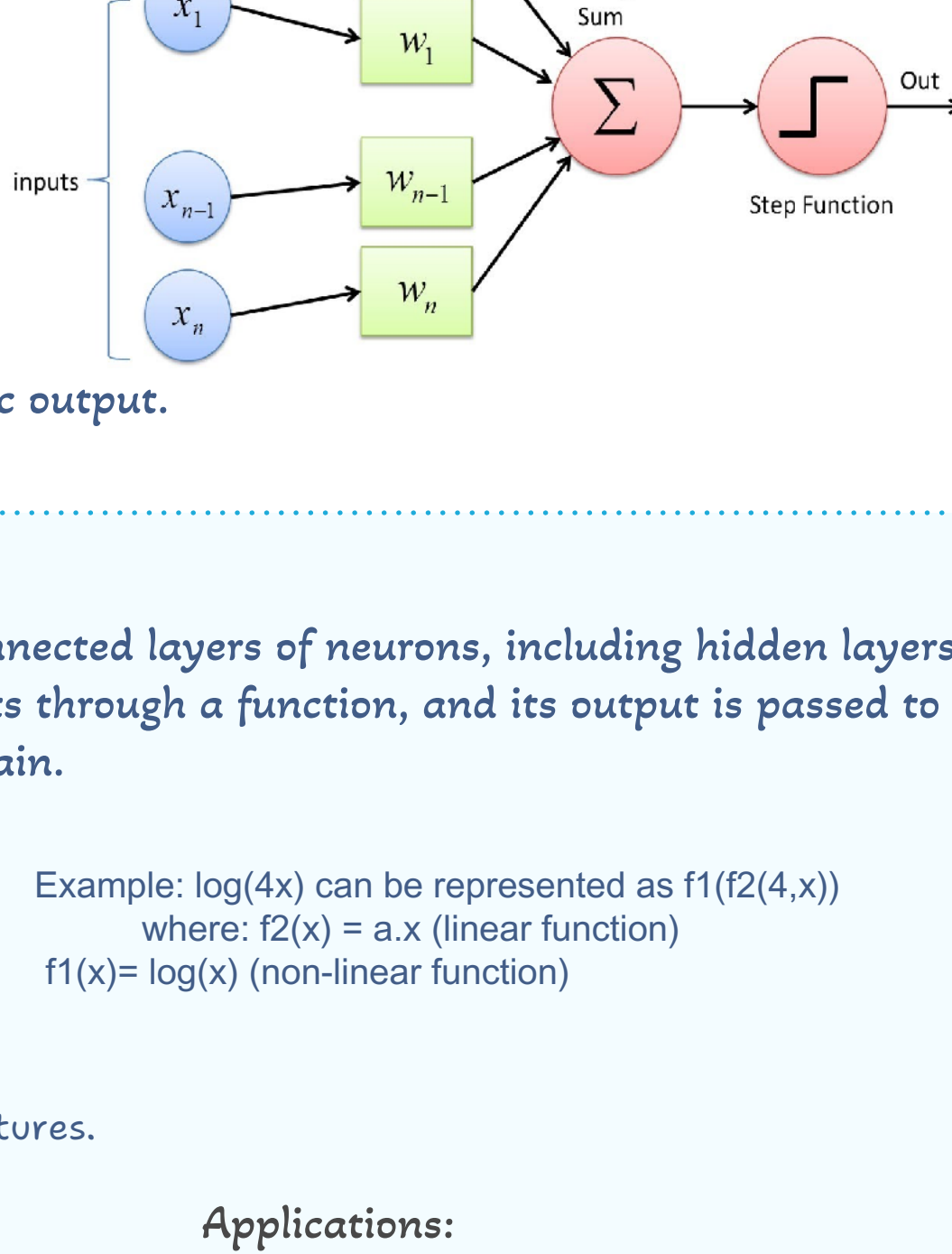
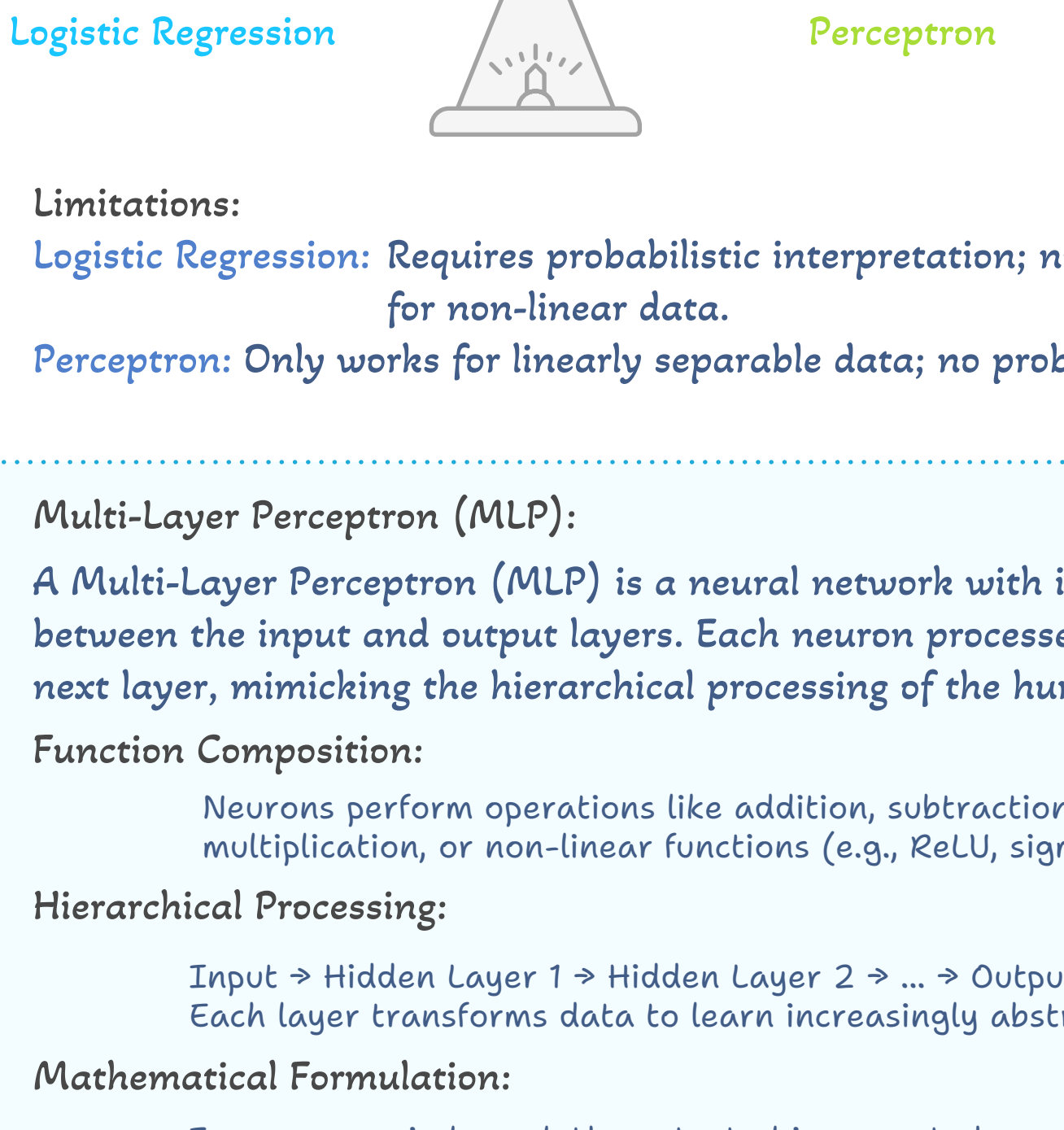
- Prediction Function:** $\hat{y}_i = \sigma(Wx + b)$ where $\sigma(z) = \frac{1}{1 + e^{-z}}$
- Learning:** Adjusts weights (W) and bias (b) using maximum likelihood estimation.
- Output:** Probabilistic (values between 0 and 1).

Perceptron:

- Prediction Function:** $f(x_i) = \sum w_j x_{ij} + b$; **Activation:** Step function (e.g., 1 if $f(x) > 0$, else 0)
- Learning:** If misclassified, update weights: $w_j \leftarrow w_j + \eta (y_i - \hat{y}_i) x_{ij}$
- Output:** Binary (0 or 1).

Sigmoid Function

Key Differences



Limitations:

- Logistic Regression:** Requires probabilistic interpretation; not ideal for non-linear data.
- Perceptron:** Only works for linearly separable data; no probabilistic output.

Multi-Layer Perceptron (MLP):

A Multi-Layer Perceptron (MLP) is a neural network with interconnected layers of neurons, including hidden layers between the input and output layers. Each neuron processes inputs through a function, and its output is passed to the next layer, mimicking the hierarchical processing of the human brain.

Function Composition:

- Neurons perform operations like addition, subtraction, multiplication, or non-linear functions (e.g., ReLU, sigmoid).

Example: $\log(4x)$ can be represented as $f_1(f_2(4, x))$ where: $f_2(x) = a \cdot x$ (linear function) $f_1(x) = \log(x)$ (non-linear function)

Hierarchical Processing:

- Input \rightarrow Hidden Layer 1 \rightarrow Hidden Layer 2 $\rightarrow \dots \rightarrow$ Output.
- Each layer transforms data to learn increasingly abstract features.

Mathematical Formulation:

- For a neuron in layer l , the output a_l is computed as:

$$a_l = \sigma(W^l a_{l-1} + b^l)$$
 Where:
 σ : Activation function (e.g., sigmoid, ReLU)
 W^l : Weight matrix for layer l
 b^l : Bias vector for layer l

Applications:

- Classification tasks (e.g., image recognition)
- Regression analysis
- Feature learning

Training a Single-Neuron Model:

1. Model Structure:

A single-neuron model consists of:

- Inputs (x_1, x_2, \dots, x_n): Features from the data.
- Weights (w_1, w_2, \dots, w_n): Learnable parameters.
- Bias (b): Shifts the decision boundary.
- Activation Function (σ): Converts the weighted sum into an output (e.g., sigmoid, ReLU).
- Output (\hat{y}): Prediction based on inputs and weights.

2. Forward Propagation:

- Compute the output \hat{y}_i : $z = \sum_{i=1}^n w_i x_i + b$ (Weighted sum)
- $\hat{y}_i = \sigma(z)$ (Activation)

Example: For binary classification, use the sigmoid activation: $\sigma(z) = \frac{1}{1 + e^{-z}}$

3. Loss Function:

Measure prediction error using a loss function:

- Binary Cross-Entropy Loss (classification):

$$L = -\frac{1}{n} \sum_{i=1}^n [y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)]$$
- Mean Squared Error (regression):

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

4. Backpropagation:

Update weights and bias using gradient descent:

- Compute gradients:

$$\frac{\partial L}{\partial w_j} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i) x_{ij}$$
- Update parameters:

$$w_j \leftarrow w_j - \eta \frac{\partial L}{\partial w_j}$$

$$b \leftarrow b - \eta \frac{\partial L}{\partial b}$$
 Where η is the learning rate.

5. Key Considerations:

Learning Rate (η): Controls step size during weight updates (too high \rightarrow divergence; too low \rightarrow slow training).

Activation Choice: Sigmoid for classification, ReLU/linear for regression.

Linearity: Single neurons can only model linear decision boundaries.

Training an MLP: Chain Rule for Backpropagation:

1. Overview:

Training a Multi-Layer Perceptron (MLP) involves adjusting weights and biases using the chain rule to compute gradients during backpropagation. This allows the network to minimize the loss function by propagating errors backward from the output layer to the input layer.

3. Chain Rule in Backpropagation:

- To update weights (W_1, W_2) and biases (b_1, b_2), compute gradients using the chain rule:

Output Layer: $\frac{\partial L}{\partial W_2} = (y - \hat{y}) \cdot \sigma'(W_2 h + b_2) \cdot h^T$

Hidden Layer: $\frac{\partial L}{\partial W_1} = [W_2^T \cdot (y - \hat{y}) \cdot \sigma'(W_2 h + b_2)] \cdot \sigma'(W_1 x + b_1) \cdot x^T$

5. Key Equations:

Output Layer Weights: $\frac{\partial L}{\partial W_2} = \delta_{out} \cdot h^T$

Hidden Layer Weights: $\frac{\partial L}{\partial W_1} = \delta_{hidden} \cdot x^T$

Output Layer Bias: $\frac{\partial L}{\partial b_2} = \delta_{out}$

Hidden Layer Bias: $\frac{\partial L}{\partial b_1} = \delta_{hidden}$

2. Mathematical Framework:

Consider an MLP with:

- Input: $x \in \mathbb{R}^d$
- Hidden: $h = \sigma(W_1 x + b_1)$ where σ is the activation function (e.g., ReLU).
- Output: $\hat{y} = \sigma(W_2 h + b_2)$
- Loss: $L = \frac{1}{2} (y - \hat{y})^2$ (Mean Squared Error).

4. Step-by-Step Gradient Calculation:

For a network with ReLU activation ($\sigma(z) = \max(0, z)$):

- Forward Pass: Compute hidden layer: $h = \text{ReLU}(W_1 x + b_1)$
- Compute output: $\hat{y} = \text{ReLU}(W_2 h + b_2)$

Backward Pass:
 Output error: $\delta_{out} = (y - \hat{y}) \odot \text{ReLU}'(\hat{y})$
 Hidden error: $\delta_{hidden} = (W_2^T \cdot \delta_{out}) \odot \text{ReLU}'(h)$
 Update weights: $W_2 \leftarrow W_2 + \eta \cdot \delta_{out} \cdot h^T$
 $W_1 \leftarrow W_1 + \eta \cdot \delta_{hidden} \cdot x^T$

6. Why the Chain Rule Matters:

- Efficiently computes gradients for deep networks by breaking derivatives into local components.
- Enables optimization algorithms like SGD to update weights iteratively.

Memoization in MLP Training

1. What is Memoization?

Memoization is an optimization technique that stores the results of computationally expensive operations during the forward pass (e.g., layer activations, intermediate derivatives) to reuse them in the backward pass. This avoids redundant recalculations, significantly speeding up training.

2. Memoization in MLPs

During training, an MLP performs:

Forward Pass:
 Computes layer outputs (e.g., $z_1 = W_1 x + b_1$, $a_1 = \sigma(z_1)$).

Backward Pass:

Uses cached values to compute gradients via the chain rule.

Stores these values in memory (memoization).

Example: Reuses z_1 and a_1 to calculate $\frac{\partial L}{\partial W_1}$.

3. Code Example: Memoization in PyTorch

```
import torch

# Forward pass (memoizes z1, a1, z2)
x = torch.randn(10, 5) # Input
z1 = torch.matmul(W1, x) + b1 # Memoized
a1 = torch.relu(z1) # Memoized
z2 = torch.matmul(a1, W2) + b2 # Memoized

# Backward pass (uses cached values)
loss = ((z2 - y)**2).mean()
loss.backward() # Gradients reuse z1, a1, z2
```

4. Trade-offs

Speed: Reduces computation time by ~50%.

Memory: Requires storing intermediate values (challenge for large models).

5. Why It Matters

Memoization is foundational for efficient deep learning. Frameworks like PyTorch and TensorFlow automate this process through computation graphs, enabling scalable training of complex models like GPT-3 and ResNet.

Back Propagation

Backpropagation is a combination of memoization & chain rule. The procedure is as follows:

- Initialize w, b, j
- For each x_i in D
 - Forward Propagation: Pass x_i forward through the network
 - Compute Loss function $L(y_i, \hat{y}_i)$
 - Compute all derivatives using chain rule & memoization
 - Update the weights again from end to the start such that loss will be minimized
- Repeat step 2 till convergence

Note: Backpropagation works only if activation functions are differentiable. If these functions are easily or fast differentiable, it speeds up the training & Neural Network.

Activation Functions

Sigmoid (σ):

The sigmoid function is a non-linear activation function that squashes input values between 0 and 1. It is particularly useful for binary classification tasks.

Mathematical Definition:
 The sigmoid function is defined as:

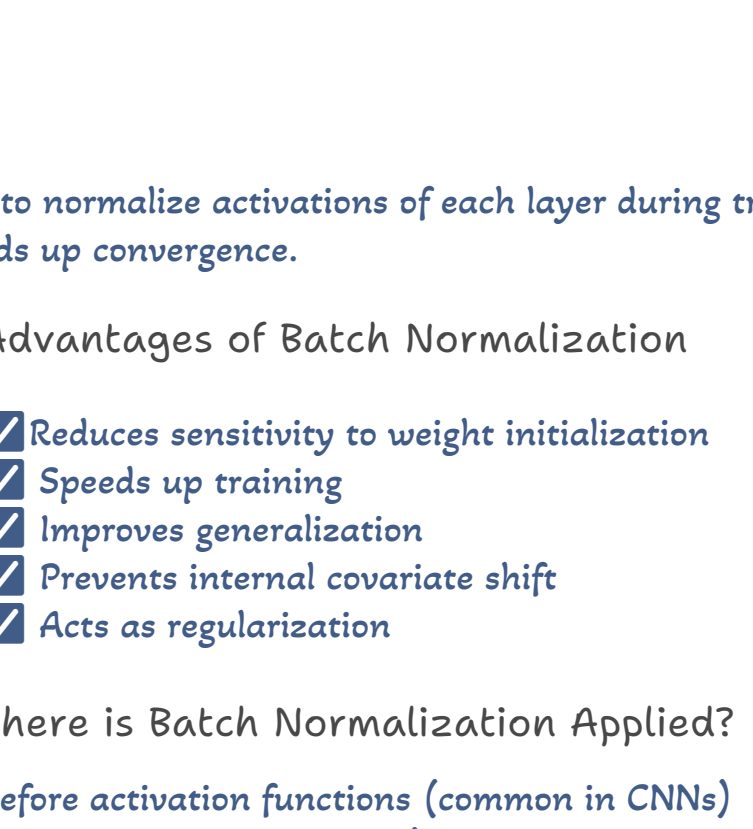
- $\sigma(z) = \frac{1}{1 + e^{-z}}$
- Where: $z = W \cdot x + b$ (weighted sum of inputs + bias)

Key Properties:

- Output range: (0, 1)
- Derivative range: (0, 0.25)
- Smooth, S-shaped curve

Derivation of the Derivative:

- The derivative of the sigmoid function is critical for backpropagation:
- Start with $\sigma(z) = \frac{1}{1 + e^{-z}}$
- Let $g(z) = (1 + e^{-z})^{-1}$, so $\sigma(z) = g(z)$
- Apply the chain rule: $\frac{d\sigma}{dz} = \left(\frac{dg}{dz}\right) \cdot \left(\frac{dz}{dz}\right)$
- Calculate individual terms: $\frac{dg}{dz} = -g \cdot z$; $\frac{dz}{dz} = -e^{-z}$
- Combine results: $\frac{d\sigma}{dz} = e^{-z} / (1 + e^{-z})^2 = \sigma(z) \cdot (1 - \sigma(z))$



Limitations
 Vanishing gradients for extreme inputs (saturates at 0/1)
 Not zero-centered (can slow down learning)

HyperbolicTangent (tanh):

The tanh function is a zero-centered activation function that squashes input values between -1 and 1. It is preferred over sigmoid for hidden layers in many architectures due to its symmetry.

Mathematical Definition:
 The tanh function is defined as:

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

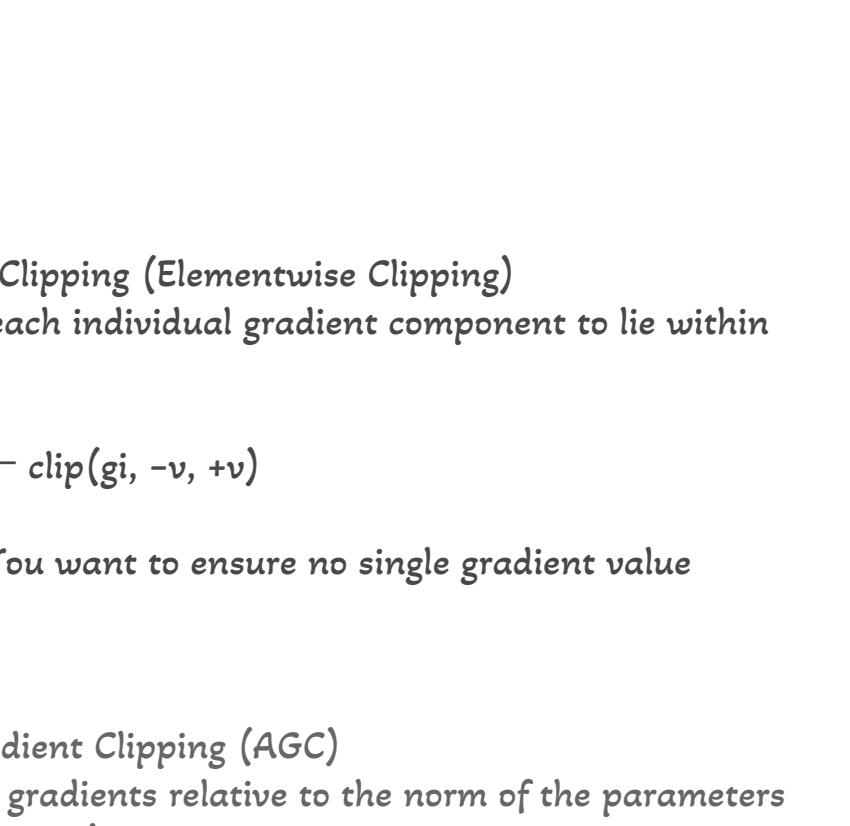
Key Properties:

- Output range: (-1, 1)
- Derivative range: (0, 1)
- Zero-centered (mitigates learning slowdown)

Derivation of the Derivative:

- The derivative of the tanh function is critical for backpropagation:
- Start with $\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$
- Differentiate using quotient rule:

$$\frac{d}{dz} \tanh(z) = 1 - \tanh^2(z)$$



Limitations
 Still suffers from vanishing gradients for extreme values
 Computationally more expensive than ReLU

Rectified Linear Unit (ReLU):

ReLU is the most widely used activation function for deep networks. It outputs the input directly if positive; otherwise, it outputs zero. Computationally efficient and mitigates vanishing gradients.

Mathematical Definition:
 The tanh function is defined as:

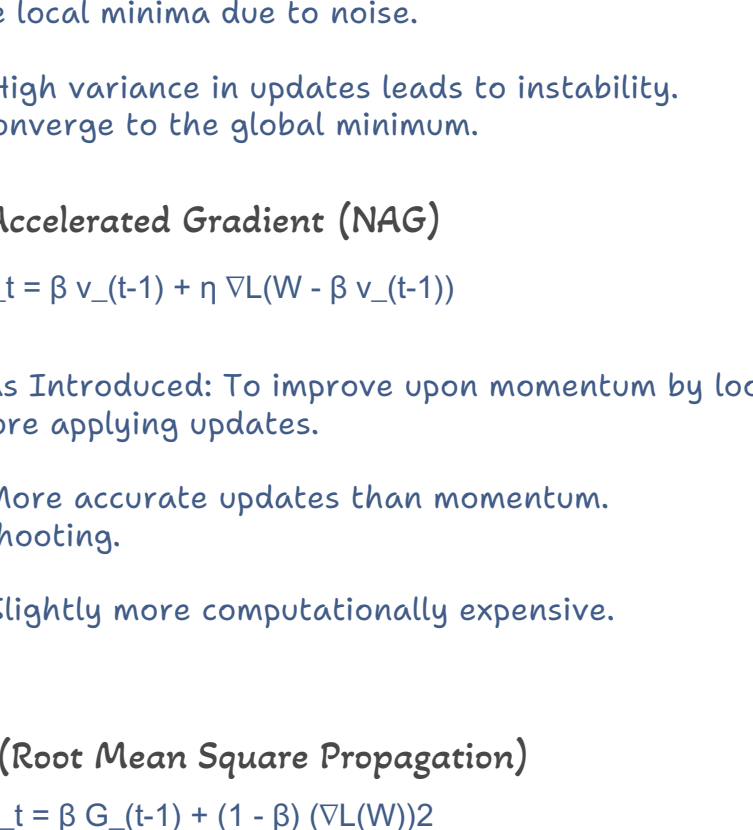
$$\text{ReLU}(z) = \max(0, z)$$

Key Properties:

- Output range: [0, ∞)
- Derivative range: (0, 1) if $z > 0$, 0 (otherwise)
- Sparsity-inducing (reduces overfitting)

Derivation of the Derivative: $\frac{d}{dz} \text{ReLU}(z) =$

$$\begin{cases} 1 & \text{if } z > 0 \\ 0 & \text{if } z \leq 0 \end{cases}$$



Limitations
 "Dying ReLU" problem: Neurons can get stuck in negative region
 Non-differentiable at $z = 0$

The Vanishing/Exploding Gradient Problem

Overview: In deep neural networks, gradients computed during backpropagation involve multiplying derivatives across layers. This can cause:

- Exploding Gradients:** Gradients grow exponentially if derivatives > 1 .
- Vanishing Gradients:** Gradients shrink to near-zero if derivatives < 1 .

Mathematical Explanation:

For a network with n layers, the gradient at the first layer is:

$$\frac{\partial L}{\partial W_1} = \frac{\partial L}{\partial y} \cdot \prod_{k=2}^n \left(\frac{\partial h_k}{\partial W_k} \right)$$

$\frac{\partial h_k}{\partial W_k} > 1 \rightarrow$ Gradient explodes
 e.g., $(1.5)^{10} \approx 58$

$\frac{\partial h_k}{\partial W_k} < 1 \rightarrow$ Gradient vanishes
 e.g., $(0.5)^{10} \approx 0.001$

Causes:

- Vanishing Gradients**
 Activation functions with small derivatives (e.g., sigmoid, tanh).
 Deep networks amplify the multiplicative effect.
- Exploding Gradients**
 Large initial weights or improper weight initialization.
 Unbounded activations (e.g., ReLU without normalization).

Solutions:

For Vanishing Gradients

- Use ReLU/Leaky ReLU activations.
- Weight initialization (He/Xavier).
- Skip connections (ResNet).

For Exploding Gradients

- Gradient clipping (torch.nn.utils.clip_grad_norm).
- Batch normalization.
- L2 regularization.

Conclusion:

The vanishing/exploding gradient problem is critical in deep learning but can be managed with techniques like ReLU, careful initialization, gradient clipping, and architectural innovations (e.g., ResNet). These strategies enable stable training of very deep networks.

Batch Normalization

What is Batch Normalization?

Batch Normalization (BatchNorm) is a technique used in deep neural networks to normalize activations of each layer during training. It improves training stability, prevents vanishing/exploding gradients, and speeds up convergence.

Why is Batch Normalization Needed?

- Internal Covariate Shift: Reduces shifting activation distributions.
- Faster Convergence: Reduces training time by keeping activations consistent.
- Prevents Vanishing/Exploding Gradients: Keeps values in a stable range.
- Regularization Effect: Acts as a mild form of regularization, reducing reliance on dropout.

Batch Normalization Formula

Compute Mean & Variance for Mini-Batch: $\mu_B = \frac{1}{m} \cdot \sum x_i$
 $\sigma_B^2 = \frac{1}{m} \cdot \sum (x_i - \mu_B)^2$

Normalize the Activations: $\hat{x}_i = \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$

Scale and Shift with Learnable Parameters: $y_i = \gamma \cdot \hat{x}_i + \beta$

Advantages of Batch Normalization

- Reduces sensitivity to weight initialization
- Speeds up training
- Improves generalization
- Prevents internal covariate shift
- Acts as regularization

Where is Batch Normalization Applied?

Before activation functions (common in CNNs)
 After activation functions (common in fully connected layers)
 Between layers in deep networks

Gradient Clipping

Overview: Gradient Clipping is a technique used to prevent the problem of exploding gradients, which can occur in deep networks, particularly in recurrent neural networks (RNNs) and long training processes.

Why Is Gradient Clipping Needed?

During backpropagation, gradients can become excessively large, leading to unstable updates in weight parameters. This can cause:

- Loss function divergence (failure to converge).
- Numerical instability (NaN values in computations).
- Poor generalization due to chaotic weight updates.

Types of Gradient Clipping:

Norm-Based Clipping (Global Norm Clipping)
 What: Clips the entire gradient vector if its L2 norm exceeds a threshold c .

Formula: If $\|g\| \geq c$, then: $g \leftarrow \frac{c}{\|g\|} \cdot g$

Used when: You want to control the overall gradient magnitude across the model.

Value-Based Clipping (Elementwise Clipping)
 What: Clips each individual gradient component to lie within $[-v, v]$.

Formula: $g_i \leftarrow \text{clip}(g_i, -v, v)$

Used when: You want to ensure no single gradient value explodes.

Per-Layer Gradient Clipping
 What: Apply norm-based or value-based clipping individually per layer.

Adaptive Gradient Clipping (AGC)
 What: Scales gradients relative to the norm of the parameters (not just gradients).

Why: Helps localize stability without affecting the entire model's gradient structure.

Formula (simplified): If $\|g\| > \epsilon \cdot \|\theta\| \Rightarrow$ scale $g \leftarrow \left(\frac{\epsilon}{\|g\|} \cdot \|\theta\|\right) \cdot g$
 Used in: Large models (e.g., Transformers, Stable Diffusion) for better stability.

Optimizers:

Optimizers are algorithms that adjust the weights of a neural network to minimize the loss function. The choice of optimizer affects convergence speed, accuracy, and overall model performance.

Gradient Descent (GD)

Formula: $W = W - \eta \nabla L(W)$

Why It Was Introduced: It is the fundamental optimization algorithm used in machine learning, forming the basis for all other optimizers.

- Pros: Simple and easy to understand.
- Guaranteed to find a minimum if convex loss function.
- Cons: Very slow for large datasets (Batch GD). May get stuck in local minima.

Momentum Optimizer

Formula: $v_t = \beta v_{t-1} + \eta \nabla L(W)$
 $W = W - v_t$

Why It Was Introduced: To reduce the oscillations in SGD and speed up convergence.

- Pros: Reduces oscillations in deep valleys. Faster convergence in certain cases.
- Cons: Can overshoot the optimal point.

Adagrad (Adaptive Gradient Algorithm)

Formula: $W = W - (\eta / \sqrt{G_t + \epsilon}) \nabla L(W)$

Why It Was Introduced: To adapt learning rates for each parameter, making it useful for sparse data.

- Pros: Good for sparse features (e.g., NLP