

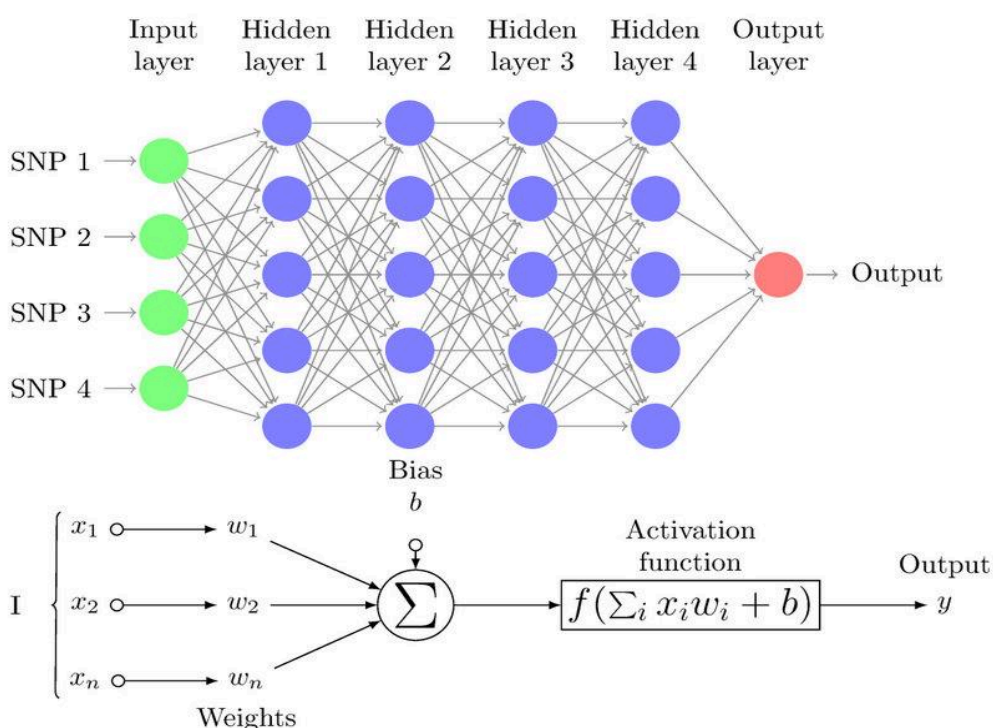
Perceptron in Deep Learning

The perceptron is one of the fundamental building blocks of neural networks and deep learning. It is a type of artificial neuron or the simplest form of a neural network, designed to mimic the function of a biological neuron. Introduced by Frank Rosenblatt in 1957, the perceptron laid the groundwork for modern neural networks.

Structure of a Perceptron

A perceptron consists of several key components:

- Inputs (Features):**
 - The perceptron takes multiple inputs, each representing a feature of the input data.
- Weights:**
 - Each input is associated with a weight, which determines the input's importance. These weights are adjusted during the learning process to minimize error.
- Bias:**
 - An additional parameter, the bias, is added to the weighted sum of the inputs. It allows the activation function to be shifted to the left or right, improving the model's flexibility.
- Activation Function:**
 - The perceptron uses an activation function to produce an output. The most common activation function for a basic perceptron is the step function, which outputs a binary value (0 or 1) based on a threshold.



Learning in a Perceptron

The learning process of a perceptron involves adjusting the weights and bias to minimize the difference between the predicted output and the actual output. This is typically done using the Perceptron Learning Rule:

1. **Initialize Weights:**
 - Start with small random values for the weights and bias.
2. **Compute Output:**
 - Calculate the output using the current weights and bias.
3. **Update Weights:**
 - Adjust the weights and bias based on the error (difference between predicted and actual output).

$$w_i^{n+1} = w_i^n + \eta(y_i - \hat{y}_i)x_i$$

learning rate ("eta") ↓
↑ *new weight* *current weight*

$$\eta = 0.1$$

Activation Functions in Neural Networks

Activation functions play a crucial role in neural networks by introducing non-linearity into the model, allowing it to learn and represent complex patterns. Without activation functions, a neural network would simply perform linear transformations, limiting its capability to solve complex tasks. Here's a detailed overview of activation functions in neural networks:

Purpose of Activation Functions

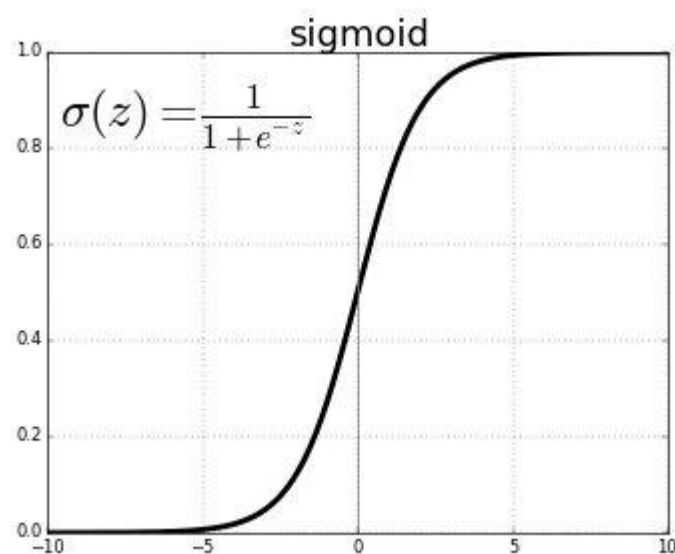
- **Non-Linearity:** Activation functions introduce non-linearities, enabling neural networks to model complex relationships and solve non-linear problems.
- **Output Transformation:** They transform the weighted sum of inputs into an output signal, which can be passed to the next layer or used as the final output.
- **Gradient Flow:** They impact the gradients used in backpropagation, influencing how the network learns.

Sigmoid Activation Function

The sigmoid activation function is one of the most commonly used activation functions in neural networks, especially in the early days of deep learning. It is a type of logistic function that maps any real-valued number into a value between 0 and 1. This squashing behavior makes it particularly useful for binary classification tasks, where the output can be interpreted as a probability.

Mathematical Definition

The sigmoid function is defined by the following equation:



Use Cases

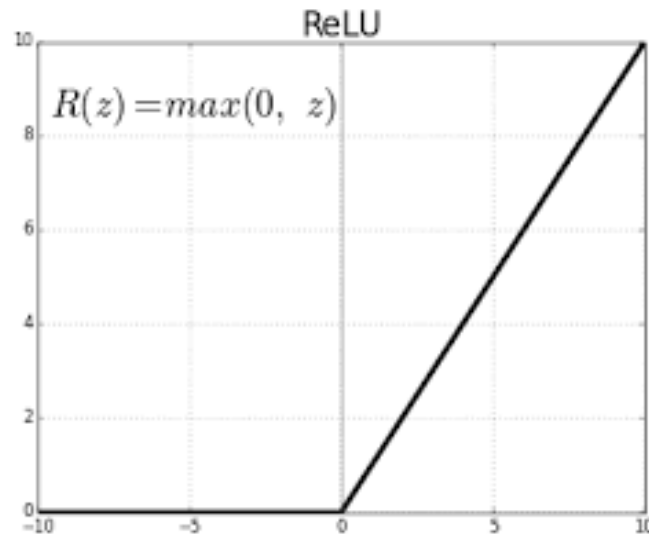
- **Binary Classification:** Often used in the output layer of a binary classification model, where the output needs to be a probability indicating the likelihood of one class over the other.
- **Logistic Regression:** The sigmoid function is the core of logistic regression models.
- **Simple Neural Networks:** Historically, it was used in early neural network architectures.

ReLU Activation Function

The Rectified Linear Unit (ReLU) is one of the most popular activation functions used in neural networks today. It has gained widespread adoption due to its simplicity and effectiveness in addressing some of the limitations of earlier activation functions like sigmoid and tanh.

Definition

The ReLU activation function is defined as:



Properties

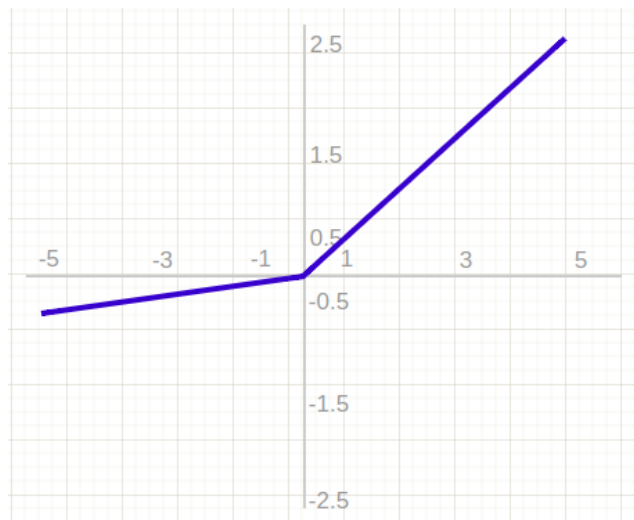
1. **Range:** The output of the ReLU function ranges from 0 to ∞ .
2. **Non-Linearity:** ReLU introduces non-linearity into the model, which is essential for learning complex patterns.
3. **Computational Efficiency:** ReLU is computationally efficient because it involves simple operations (max and comparison).

Disadvantages

1. **Dying ReLU Problem:**
 - If a neuron's input is always negative, it can "die," meaning it will output zero for all inputs and stop contributing to the model's learning. This issue is known as the dying ReLU problem.

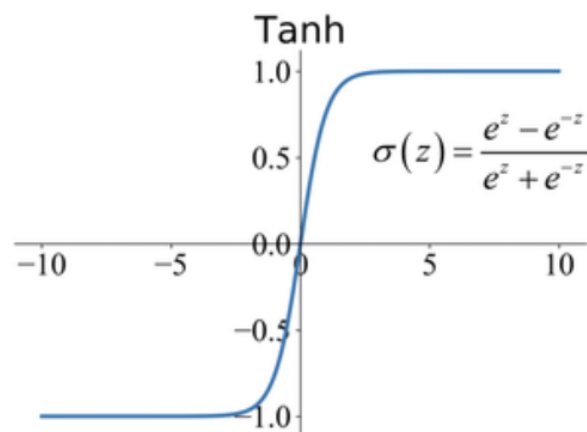
Leaky ReLU Activation Function

Leaky ReLU (Leaky Rectified Linear Unit) is a variant of the ReLU activation function designed to address the "dying ReLU" problem, where neurons can become inactive and always output zero for any input. This issue can occur in standard ReLU when the input to a neuron is always negative, causing the gradient to be zero and preventing the neuron from updating its weights during training.



Tanh Activation Function

The tanh (hyperbolic tangent) activation function is another widely used activation function in neural networks. It is similar to the sigmoid function but has some advantages that make it preferable in certain contexts.



Properties

1. **Range:** The output of the tanh function ranges from -1 to 1 .
2. **Non-Linearity:** Like other activation functions, tanh introduces non-linearity, allowing neural networks to model complex patterns.
3. **Zero-Centered Output:** The tanh function is zero-centered, meaning its outputs are evenly distributed around zero, which can help in faster convergence during training.

Softmax Activation Function

The softmax activation function is commonly used in the output layer of a neural network for multi-class classification problems. It converts raw score outputs (logits) from the network into probabilities, which sum to 1, making it easy to interpret the output as a probability distribution over multiple classes.

$$\text{softmax}(z_j) = \frac{e^{z_j}}{\sum_{k=1}^K e^{z_k}} \text{ for } j = 1, \dots, K$$

Use in Neural Networks

The softmax function is typically used in the final layer of a neural network for multi-class classification tasks. It converts the logits (raw class scores) into probabilities, allowing the model to output a probability distribution over the classes. This is particularly useful for:

- **Multi-Class Classification:** When there are more than two classes, softmax is used to provide a probability for each class.

Forward Propagation

Definition: Forward propagation refers to the process where input data is fed through the neural network, layer by layer, to generate an output.

Steps Involved:

1. **Input Layer:** The process begins with the input layer, where the initial data (features) are fed into the network.
2. **Hidden Layers:** Each hidden layer receives inputs from the previous layer. The neurons in each hidden layer calculate a weighted sum of their inputs (from the previous layer) and apply an activation function to produce an output. This output becomes the input for the next layer.
3. **Output Layer:** The final hidden layer's output is passed through the output layer, where it undergoes a final activation function. The output of the output layer depends on the nature of the problem being solved (e.g., regression, binary classification, multi-class classification).
4. **Prediction:** The output layer's activation generates predictions or classifications based on the input data.

Purpose: Forward propagation serves to compute predictions or classifications based on the current weights and biases of the neural network. It is a deterministic process where the inputs are transformed through the network's layers to produce an output.

Backward Propagation (Backpropagation)

Definition: Backward propagation, or backpropagation, refers to the process where the network computes gradients of the loss function with respect to each weight and bias in the network.

Steps Involved:

1. **Loss Calculation:** First, the loss function is calculated using the predicted output and the actual target values. The loss function quantifies how well or poorly the model performed on the training data.
2. **Gradient Calculation:** Starting from the output layer, the gradients of the loss function with respect to the weights and biases of each layer are computed using the chain rule of calculus. These gradients indicate how much each weight and bias contributed to the error in the prediction.
3. **Weight Update:** The gradients computed in the previous step are used to update the weights and biases of the neural network. This step involves using an optimization algorithm (e.g., gradient descent) to adjust the weights in the direction that reduces the error and improves the model's performance.
4. **Iterative Process:** Steps 1-3 are repeated iteratively for each batch of data (in mini-batch gradient descent) or for each data point (in stochastic gradient descent) until the model converges to optimal weights that minimize the loss function.

Purpose: Backward propagation enables the neural network to learn from the errors it makes during forward propagation. By computing gradients and updating weights accordingly, the network adjusts its parameters to improve its predictions over time.

Backward Propagation

$$W_{\text{new}} = W_{\text{old}} - \eta \frac{\partial E}{\partial W}$$

→ using sigmoid function $G = \frac{1}{1 + e^{-x}}$ gradient descent

$$\begin{array}{c}
 z_1 = w_1 a_0 \quad z_2 = w_2 a_1 \\
 x_1 \rightarrow \text{circle} \rightarrow \text{circle} \rightarrow \text{circle} \rightarrow \hat{y} / \text{predicted} = a_2 / G(z_2) \\
 a_0 = x_1 \quad a_1 = G(z_1) \quad a_2 = G(z_2) \\
 \text{I/P} \quad \quad \quad \text{H/L} \quad \quad \quad \text{O/P}
 \end{array}$$

$a_0 = \text{output of I/P layer}$
 $a_1 = \text{output of H/L}$
 $a_2 = \text{final output}$

$\hat{y} = a_2$

1) $E = (y - \hat{y})^2 \rightarrow (y - a_2)^2 \rightarrow F(a_2) \quad [2(y - a_2)]$

2) $a_2 = G(z_2) \rightarrow F(z_2) \quad [G(z_2)(1 - G(z_2))]$

3) $z_2 = w_2 a_1 \rightarrow F(w_2) \quad a_1$

BY APPLYING chain rule

$$\begin{aligned}
 \frac{\partial E}{\partial w_2} &= \frac{\partial E}{\partial a_2} \times \frac{\partial a_2}{\partial z_2} \times \frac{\partial z_2}{\partial w_2} \\
 \frac{\partial E}{\partial w_2} &= -2(y - a_2) [G(z_2)(1 - G(z_2))] \cdot a_1
 \end{aligned}$$

$$W_{\text{new}} = W_{\text{old}} - \eta \frac{\partial E}{\partial w_2}$$