

Why Are Multilayer (Deep) Neural Networks Needed?

Why Are Multilayer (Deep) Neural Networks Needed?

Multilayer (deep) neural networks are needed because they allow the model to learn hierarchical representations of data. While a single-layer neural network (or a shallow network) can only learn linear or simple decision boundaries, deep networks with multiple hidden layers can model complex patterns and non-linear relationships. Each layer in a deep network extracts more abstract features from the data, making it possible for deep neural networks to perform better on complex tasks such as image recognition, speech processing, and natural language understanding.

The deeper the network, the more layers of abstraction it can represent. For example, in image recognition:

- The first layer might learn to detect simple edges.
- The next layer might combine those edges into shapes.
- Deeper layers might combine shapes into objects or complex patterns.

This hierarchical learning ability is what makes deep networks superior in tasks that involve high-dimensional, structured data.

What is the Structure of Weight Matrix (How Many Rows and Columns)?

The weight matrix in a neural network is defined for each layer and connects the neurons of one layer to the neurons of the next layer. The structure of the weight matrix depends on the number of neurons in both the layers being connected:

- **Rows:** The number of rows in the weight matrix corresponds to the number of neurons in the next (output) layer.
- **Columns:** The number of columns corresponds to the number of neurons in the previous (input) layer.

For example, if layer (L) has (m) neurons and layer (L+1) has (n) neurons, the weight matrix between these two layers will have dimensions $(n \times m)$ (i.e., (n) rows and (m) columns).

Describe the Gradient Descent Method

Gradient descent is an optimization algorithm used to minimize a loss function (or cost function) by iteratively adjusting the model's parameters (weights and biases). The goal is to find the values of parameters that result in the lowest possible loss.

Steps of Gradient Descent:

1. **Initialize the Parameters:** Start with random values for weights and biases.
2. **Calculate the Gradient:** Compute the gradient of the loss function with respect to each parameter. This represents how much the loss changes with a small change in each parameter.
3. **Update Parameters:** Adjust the parameters by moving them in the opposite direction of the gradient. The learning rate controls the step size of the update.

$$\begin{aligned} & [\\ & \theta = \theta - \eta \cdot \nabla J(\theta) \\ &] \end{aligned}$$

where:

- θ represents the parameters (weights and biases).
 - η is the learning rate.
 - $\nabla J(\theta)$ is the gradient of the loss function with respect to θ .
4. **Repeat:** Continue updating the parameters until the loss function converges to a minimum (local or global).

Describe Forward Propagation and Backpropagation for Deep Neural Networks

Forward Propagation

Forward propagation is the process by which input data passes through the layers of the neural network to produce an output (prediction).

Steps in forward propagation:

1. **Input Layer:** The input data is fed into the input layer.
2. **Weighted Sum:** For each neuron in a layer, calculate the weighted sum of the inputs from the previous layer, plus a bias term.

$$\begin{aligned} & [\\ & z^{(l)} = W^{(l)} a^{(l-1)} + b^{(l)} \\ &] \end{aligned}$$

where $W^{(l)}$ is the weight matrix, $a^{(l-1)}$ are the activations from the previous layer, and $b^{(l)}$ is the bias term.

3. **Activation Function:** Apply an activation function to the weighted sum to produce the output (activation) for each neuron.

$$\begin{aligned} & [\\ & a^{(l)} = \sigma(z^{(l)}) \\ &] \end{aligned}$$

4. **Output Layer:** The final layer outputs a prediction, which may undergo another activation function (such as softmax for classification).

Backpropagation

Backpropagation is the process of calculating the gradients of the loss function with respect to the weights, using the chain rule of calculus, and then updating the weights using gradient descent.

Steps in backpropagation:

1. **Compute Loss:** Calculate the error between the predicted output and the actual target (e.g., using mean squared error or cross-entropy).
[
$$L = \frac{1}{m} \sum (y_{\text{pred}} - y_{\text{true}})^2$$

]
 2. **Calculate Gradients for Output Layer:** Compute the derivative of the loss with respect to the activations of the output layer and propagate it back to compute the gradients with respect to the weights.
 3. **Propagate Gradients Backwards:** Use the chain rule to calculate the gradients of the loss with respect to the weights and activations in each preceding layer.
[
$$\delta^{(l)} = \left(W^{(l+1)} \delta^{(l+1)} \right) \circ \sigma'(z^{(l)})$$

]
where $(\delta^{(l)})$ is the gradient of the loss with respect to the activations at layer (l) , and $(\sigma'(z^{(l)}))$ is the derivative of the activation function at layer (l) .
 4. **Update Weights:** Adjust the weights using the gradients and learning rate.
-

Activation Functions

1. Linear Activation

- **Formula:** $f(x) = x$
- **Purpose:** Linear activation is rarely used in hidden layers but can be applied in the output layer for regression tasks.
- **Use Case:** Regression problems.

2. ReLU (Rectified Linear Unit)

- **Formula:** $f(x) = \max(0, x)$
- **Purpose:** Introduces non-linearity while avoiding the vanishing gradient problem (compared to sigmoid/tanh). It allows for faster training and better performance in deep networks.
- **Use Case:** Commonly used in hidden layers of deep networks, especially in Convolutional Neural Networks (CNNs).

3. Sigmoid

- **Formula:** $f(x) = \frac{1}{1 + e^{-x}}$

- **Purpose:** Squashes the input to a range between 0 and 1, making it useful for binary classification.
- **Use Case:** Often used in the output layer of binary classification problems.

4. Tanh (Hyperbolic Tangent)

- **Formula:** $f(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$
- **Purpose:** Similar to sigmoid but squashes the input to a range between -1 and 1. It can center the data, which may lead to faster convergence.
- **Use Case:** Sometimes used in hidden layers, especially in Recurrent Neural Networks (RNNs).

5. Softmax

- **Formula:** $f(x_i) = \frac{e^{x_i}}{\sum_j e^{x_j}}$
- **Purpose:** Converts a vector of values into probabilities, where the sum of probabilities is 1. It is commonly used in multi-class classification tasks.
- **Use Case:** Output layer of multi-class classification models.

These activation functions, along with deep neural networks, allow models to learn complex patterns in data and apply them to tasks such as classification, regression, and image recognition. Each activation function has specific use cases based on the nature of the problem and the architecture of the network.