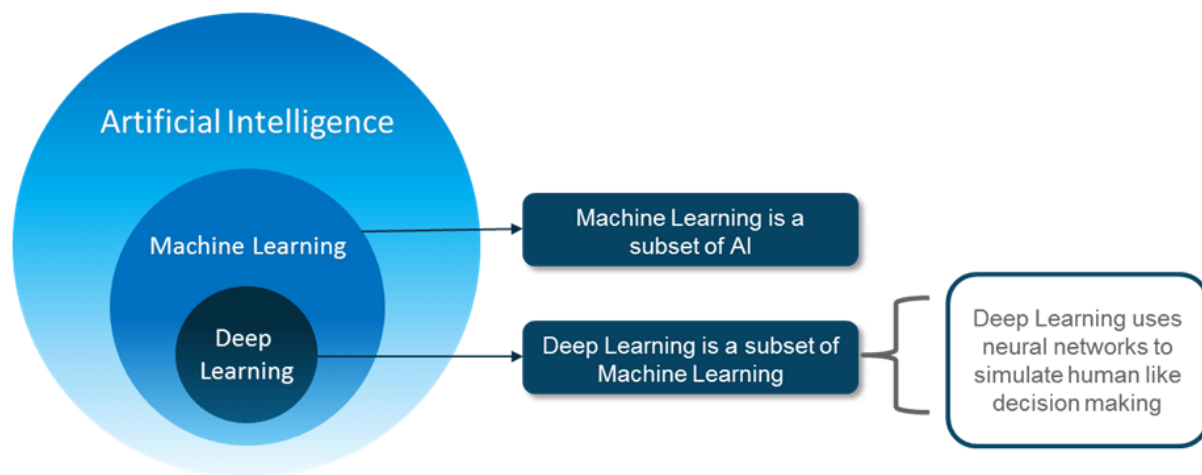


MODULE-I

MACHINE LEARNING WITH SHALLOW NEURAL NETWORKS

Deep Learning

Deep learning, a subset of machine learning, utilizes algorithms inspired by the human brain's structure and functioning. It employs artificial neural networks to build sophisticated models capable of addressing complex challenges. This technique is often applied to process and analyze unstructured data.



Artificial Intelligence (AI) refers to technology that enables machines to mimic human behavior and make independent decisions. This concept, both simple and intriguing, has long sparked debate. For years, many believed that computers could never match the complex capabilities of the human brain.

Historically, AI's development was hindered by limited data and insufficient computing power. However, with the rise of Big Data and the advancement of GPU technology, AI has evolved from a theoretical idea into a practical reality. **Machine Learning (ML)**, a subset of AI, uses statistical techniques to help machines improve performance by learning from past experiences.

Deep Learning (DL), a specialized branch of ML, allows the training of deep neural networks with multiple layers. By leveraging these neural networks, deep learning enables machines to replicate human decision-making processes.

The Need for Deep Learning

Traditional machine learning algorithms struggle with high-dimensional data—data with thousands of inputs and outputs across many dimensions. Handling such data is labor-intensive and computationally expensive, a challenge known as the **Curse of Dimensionality**. Deep learning, however, excels at managing high-dimensional data, efficiently identifying relevant features without human intervention.

What is Deep Learning?

Deep Learning is a branch of Machine Learning that utilizes similar algorithms to train deep neural networks, significantly improving accuracy. It mimics the cognitive processes of the human brain by learning through experience. Just as the brain's billions of neurons allow us to solve complex problems, deep learning uses artificial neurons (perceptrons) and artificial neural networks to replicate this cognitive power in machines.

What is a Neural Network?

A neural network, inspired by the biological neuron, is a collection of interconnected neurons that perform basic functions in our everyday activities. Examples include:

- A child learning to recognize an apple
- An animal distinguishing between its owner and others
- Sensing whether an object is hot or cold

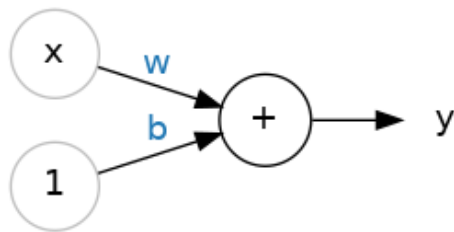
These tasks are handled by our brain's neural networks. Scientists have developed artificial neural networks (ANNs) that simulate these processes.

The Perceptron Model

An artificial neuron, or perceptron, is a supervised learning algorithm that classifies data into two categories, making it a binary classifier.

The Linear Unit:

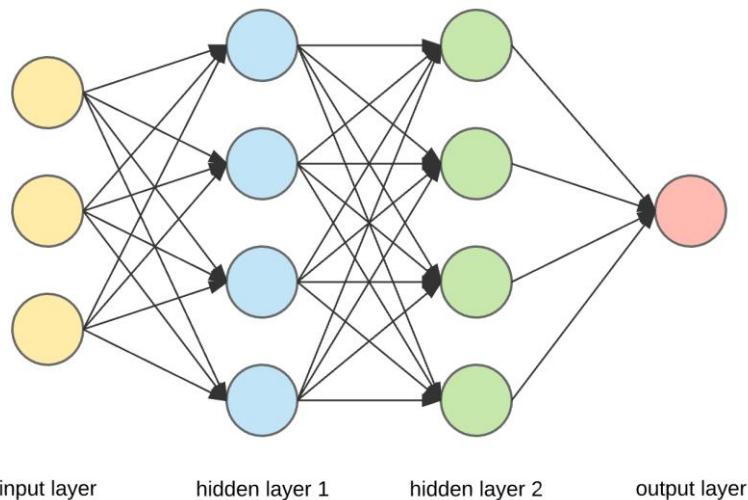
"The fundamental building block of a neural network is the individual neuron. Below is a basic illustration of a neuron (or unit) with a single input:"



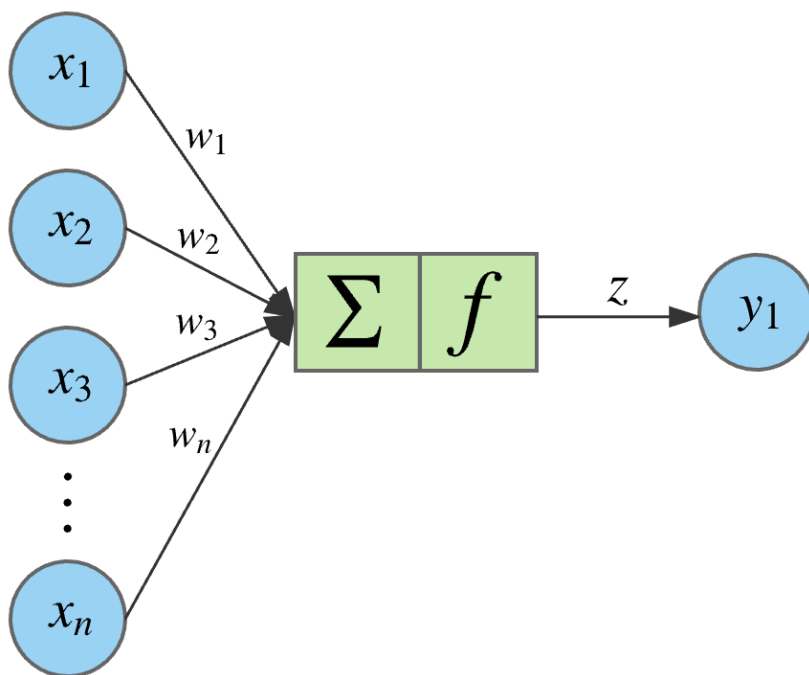
The Linear Unit: $y = wx + b$

- The neuron receives an input, denoted as x , connected by a weight w . When this input passes through the connection, it is multiplied by the weight, resulting in $w \times x$ reaching the neuron.
- The neural network improves its performance by adjusting these weights.
- Additionally, a special weight called the bias, denoted as b , is introduced. Unlike other weights, the bias does not correspond to any input. Instead, a fixed value of 1 is used, allowing the neuron to receive only the bias (since $1 \times b = b$).
- The bias enables the neuron to adjust its output independently of its inputs.
- The output of the neuron, represented as y , is calculated by summing all incoming values through its connections. The neuron's activation can be expressed using the formula $y = wx + b$, where w is the weight and x is the input.
- Are you familiar with the formula $y = wx + b$? This is the equation of a line, specifically the slope-intercept form, where w represents the slope and b denotes the y -intercept.

What is an Artificial Neural Network Model? An artificial neural network (ANN) is a type of neural network that consists of several layers: an input layer, one or more hidden layers, and an output layer.



Upon closer observation, it becomes evident that each node in a layer is intricately connected to every node in the adjacent layer. As more hidden layers are added, the network's depth increases. Let's take a deeper look at the characteristics of a single node within the output or hidden layers.



This node receives multiple inputs, computes the weighted sum, and then processes this sum through a non-linear activation function. The resulting output is passed as the input to the nodes in the next layer. It's important to note that the signal flows sequentially from left to right across the network. After all nodes process the data, the final output is produced.

Below is the equation that represents the operation within a node:

$$z = f(b + x \cdot w) = f \left(b + \sum_{i=1}^n x_i w_i \right)$$

$$x \in d_{1 \times n}, w \in d_{n \times 1}, b \in d_{1 \times 1}, z \in d_{1 \times 1}$$

The equation mentioned above indicates that b represents the bias, which serves as an input to all nodes and consistently takes a value of 1. The bias plays a key role in shifting the output of the activation function either to the left or the right.

Glossary of Artificial Neural Network Model

In this section, we'll introduce some fundamental terms that are essential for understanding an artificial neural network model.

Inputs:

The data initially fed into a neural network is called the "input." These inputs provide the network with the information it needs to make decisions or predictions. Typically, neural networks are designed to accept sets of real values as inputs, which are then passed to neurons in the input layer.

Training Set:

A training set consists of input data with known outputs. It is used to train the neural network, helping it learn and memorize the correct outputs for specific inputs.

Outputs:

A neural network generates outputs based on the data it processes. These outputs are usually real values or Boolean decisions produced by the neurons in the output layer.

Neuron:

The basic unit of a neural network is the neuron, also called a perceptron. Each neuron takes input values, processes them using a non-linear activation function, and generates an output. The output is then passed to the next layer of neurons. Common activation functions include TanH, Sigmoid, and ReLU, which are crucial for introducing non-linearity and improving the network's training efficiency.

Weight Space:

Each neuron in a neural network has a numerical weight. These weights are adjusted during the training process to produce the desired output. Training the network involves fine-tuning these weights and biases, often using a technique called backpropagation, to optimize performance.

Advantages of Artificial Neural Networks (ANN)

- **Parallel Processing Capability:**

Artificial neural networks can perform multiple tasks simultaneously, thanks to their ability to process data in parallel.

- **Distributed Data Storage:**

Unlike traditional programming, ANNs store data across the entire network, not in a centralized database. This distributed storage allows the network to function even if some data is lost in one location.

- **Ability to Operate with Limited Knowledge:**

After training, an ANN can still produce outputs even when incomplete data is provided. The performance impact depends on the importance of the missing data.

- **Memory Distribution:**

ANNs can adapt based on examples provided during training. The network learns from these examples and adjusts its behavior accordingly. The quality of the examples directly influences the performance, as incomplete or inaccurate data can lead to incorrect outputs.

- **Fault Tolerance:**

ANNs are resilient and can continue to function even if some of their components fail, thanks to their inherent fault tolerance.

Disadvantages of Artificial Neural Networks (ANN)

- **Determining the Correct Network Structure:**

There is no fixed guideline for determining the optimal structure of an ANN. Achieving the right structure often requires experimentation, practical knowledge, and learning from trial and error.

- **Unclear Network Behavior:**

A significant challenge with ANNs is their often opaque behavior. When they produce a

result, they typically do not provide an explanation for their decisions, which can reduce confidence in their outcomes.

- **Hardware Dependency:**

ANNs rely heavily on hardware with parallel processing capabilities to function effectively. This makes them dependent on specific equipment requirements.

- **Data Representation Challenges:**

ANNs process numerical data, so any problem must be translated into numbers before being input into the network. The way this data is represented directly affects the network's performance, and its success depends on the user's skill in converting the problem into an appropriate numerical format.

- **Uncertainty in Training Duration:**

The time required to train an ANN is often unpredictable and can be constrained by a predefined error threshold, which may not always yield the best results.

Neural Architectures for Binary Classification Models:

Least Square Regression Method:

Least squares is a popular technique in regression analysis used to estimate unknown parameters by constructing a model that minimizes the sum of the squared differences between observed and predicted values. One of the most commonly used methods for curve fitting involves reducing the sum of squared errors to minimize discrepancies. This approach helps generate the best-fitting line tailored to your specific data points.

Finding the Line of Best Fit Using Least Square Regression

When analyzing a set of paired numbers and its corresponding scatter plot, the line of best fit is the straight line that can be drawn through the scatter points to represent the relationship between them most accurately. The equation of the straight line is represented as:

$$Y = mX + c$$

where:

Y: Dependent Variable

m: Slope

X: Independent Variable

The goal is to find the values of the slope and y-intercept, which are then substituted into the equation along with the independent variable X to compute the dependent variable Y. For a dataset with 'n' data points, the slope can be determined using the following formula:.

$$m = \frac{\sum(x - \bar{x})(y - \bar{y})}{\sum(x - \bar{x})^2}$$

Then, the y-intercept is calculated using the formula:

$$c = \bar{y} - m * \bar{x}$$

At last, we substitute these values in the final equation

$$Y = mX + c.$$

Logistic regression:

Logistic Regression is a core machine learning algorithm primarily used for binary classification tasks. While it's simpler than more complex deep learning models, it remains relevant and frequently serves as a baseline or a foundational component in deep learning workflows. Below is an overview of logistic regression within the context of deep learning:

Logistic Regression Overview

Objective: The primary goal of logistic regression is to estimate the probability that a given input belongs to a specific class. It is typically used for binary classification problems, where there are two possible outcomes.

$$y' = p(y = 1|x)$$

Mathematical Model

Input: X is a matrix with dimensions $n \times m$, where n represents the number of features and m denotes the number of training examples.

Parameters: W represents a Weight Matrix with dimensions $n \times 1$, where n signifies the number of features within the dataset X . The Bias b serves the purpose of regulating the threshold value at which the activation function is activated.

Output:

$$y' = \sigma(W^T X + b)$$

Activation Function

Activation Functions are essential components in neural networks that introduce non-linearity, enabling the model to learn and capture complex patterns and relationships in the data. Without these functions, a neural network would essentially be limited to linear modeling, restricting its ability to solve more intricate problems.

Key Concepts of Activation Functions:

1. Purpose:
 - Non-linearity: Activation functions enable neural networks to model non-linear relationships between inputs and outputs.
 - Learning Complex Patterns: By applying non-linear transformations across multiple layers, activation functions help the network learn and represent intricate patterns in the data.

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

Loss Function

Loss Functions are critical elements in machine learning and deep learning that assess how well a model's predictions align with the actual target values. They quantify the discrepancy between predicted and true values, providing a measure of error that guides the optimization process during model training.

$$L(y', y) = -[y \log y' + (1 - y) \log(1 - y')]$$

Let's examine the effectiveness of this loss function for logistic regression:

1. When $y=1$, the loss function is represented as $L(y', y) = -\log y'$. To minimize this loss function, the value of $\log y'$ should be maximized. This can be achieved when y' approaches 1, resulting in predicted values that closely align with actual values.
2. When the loss function is calculated at $y=0$, it results in $L(y', y) = -\log(1-y')$. In order to minimize the loss function, the value of $\log(1-y')$ should be maximized. This can be achieved by reducing the value of y' , making it closer to 0. By doing so, the predicted value will be more aligned with the actual value.
3. The loss function described above exhibits convexity, ensuring the presence of a sole global minimum and preventing the network from getting trapped in local minima commonly found in non-convex loss functions.

Widrow–Hoff Rule or Delta Learning Rule:

The Least Mean Squares (LMS) method, introduced by Bernard Widrow and Marcian Hoff, is designed to minimize errors across all training patterns in a supervised learning setting. This algorithm employs a continuous activation function.

At its core, the LMS method is based on the gradient descent approach, which iteratively adjusts the synaptic weights to reduce the error between the predicted and desired output values. The delta rule is used to adjust these weights, progressively decreasing the overall input to the output unit to better align with the target value.

To adjust the synaptic weights, the delta rule can be expressed as:

$$\Delta w_i = \alpha \cdot x_i \cdot e_j$$

where:

Δw_i represents the change in weight for the i th pattern;

α is a positive and constant learning rate;

x_i is the input value from the pre-synaptic neuron;

e_j is the difference between the desired/target output and the actual output y_{in} ($e_j = t - y_{in}$).

It is important to note that the above delta rule is intended for a single output unit only.

The updating of weight can be done in the following two cases –

Closed - form solutions:

The closed-form solution is preferable for smaller datasets if the cost of computing a matrix inverse is not a concern. For larger datasets, or those where the inverse of $\mathbf{X}^T \mathbf{X}$ may not exist due to perfect multicollinearity, gradient descent (GD) or stochastic gradient descent (SGD) approaches are recommended. The linear regression model is defined as:

$$y = w_0 x_0 + w_1 x_1 + \dots + w_m x_m = \sum_{j=0}^m \mathbf{w}^T \mathbf{x}$$

The response variable is denoted as y , the sample vector with m dimensions is represented by \mathbf{x} , and the weight vector (a vector of coefficients) is symbolized as \mathbf{w} . It is important to acknowledge that w_0 signifies the intercept on the y -axis in the model, thus x_0 is equivalent to 1. By employing the closed-form solution, which is the normal equation method, we determine the weights of the model subsequently:

$$\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

Neural Architectures for Multiclass Models:

Multilayer Perceptron (MLP)

A Multilayer Perceptron (MLP) is a type of artificial neural network with multiple layers of neurons. It is one of the simplest deep learning models and is utilized for various predictive tasks, including classification and regression.

Structure of MLP

1. Layers:

- Input Layer: This is the first layer where the data is introduced into the network. Each neuron in this layer represents a feature from the input data.
- Hidden Layers: These are intermediate layers located between the input and output layers. Neurons in hidden layers process inputs through weighted connections and activation functions.
- Output Layer: This is the final layer that generates the model's predictions. For classification tasks, the output layer typically uses a softmax or sigmoid activation function. For regression tasks, a linear activation function is often employed.

2. Neurons:

- Each neuron in a layer receives inputs, calculates a weighted sum, adds a bias term, and then applies an activation function to produce an output.

3. Weights and Biases:

- Weights: These are parameters learned during training that control the strength of the connections between neurons.
- Biases: These are additional parameters added to the weighted sum before applying the activation function, which allows the model to fit the data more flexibly.

4. Activation Functions:

- Purpose: Activation functions introduce non-linearity into the model, enabling it to learn complex patterns.
- Common Functions: ReLU (Rectified Linear Unit), Sigmoid, Tanh.

How MLP Works

1. Forward Propagation:

- Input data is passed through the network, layer by layer.
- Each neuron computes a weighted sum of its inputs, adds a bias, and applies an activation function.
- The output of each layer serves as the input to the next layer, culminating in the final output.

2. Loss Calculation:

- The predicted output is compared to the true target using a loss function (e.g., cross-entropy loss for classification or mean squared error for regression).
- The loss function quantifies the error between the predicted values and the actual values.

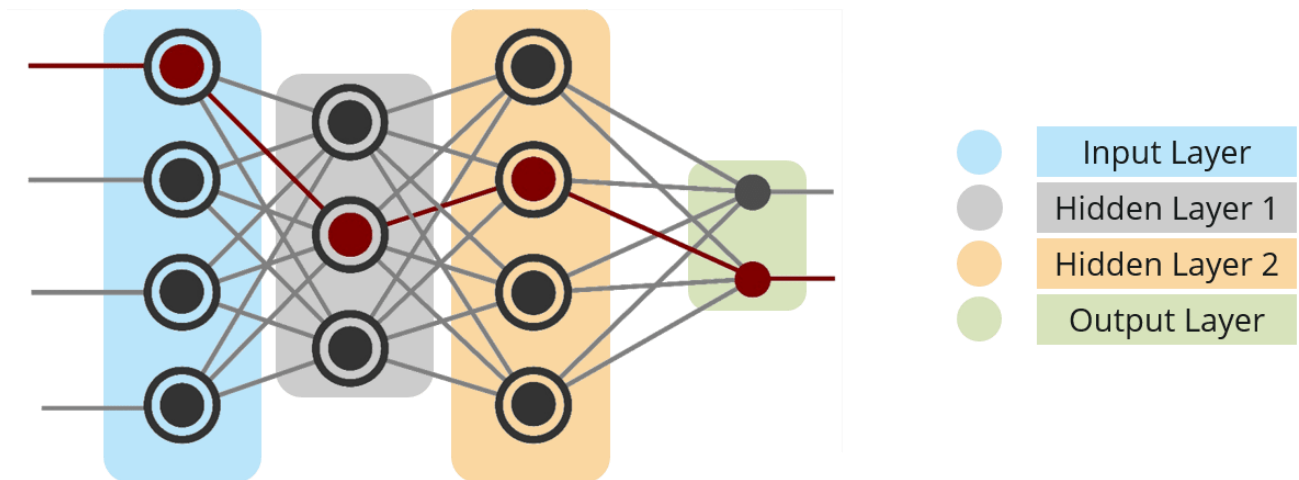
3. Backward Propagation:

- Gradient Computation: Gradients of the loss function with respect to the weights and biases are computed using the chain rule.
- Weight Update: Weights and biases are updated using an optimization algorithm (e.g., gradient descent or Adam) to minimize the loss function.

4. Training:

- The network is trained through multiple iterations (epochs) where forward and backward propagation are repeated.
- The optimization algorithm adjusts the weights and biases to enhance the model's performance.

Consider the diagram below:



Weston-Watkins SVM (or) Multi-class SVM Loss:

Weston-Watkins SVM (also known as Multi-class SVM or Weston-Watkins loss) is a variant of the Support Vector Machine (SVM) tailored for multi-class classification problems. Unlike binary SVMs, which are designed to separate data into two classes, multi-class SVMs extend this concept to handle scenarios involving more than two classes.

Weston-Watkins SVM (Multi-class SVM) Overview

Objective: The primary goal is to develop a model capable of classifying input data into one of several classes. This is achieved by minimizing a loss function that quantifies the discrepancy between predicted and actual class labels.

Weston-Watkins Loss Function

1. **Formulation:** The Weston-Watkins loss function extends the binary SVM loss function to accommodate multi-class classification. It introduces a structured loss that penalizes misclassifications by considering the relationships among multiple classes.
2. **Loss Function:** For a given data point, the Weston-Watkins loss function can be expressed as:

$$L(\mathbf{w}, \mathbf{x}, y) = \sum_{k \neq y} \max(0, 1 - (\mathbf{w}_y^T \mathbf{x} - \mathbf{w}_k^T \mathbf{x}))$$

where:

- \mathbf{x} is the input feature vector.
- y is the true class label.
- \mathbf{w}_y and \mathbf{w}_k are the weight vectors corresponding to the true class and the incorrect class k , respectively.

Softmax regression:

What is the Softmax Function?

The Softmax Function is a mathematical function commonly used in multi-class classification problems within machine learning and deep learning. It transforms raw scores (also known as logits) output by the model into probabilities, allowing for a probabilistic interpretation of the neural network's output. By applying the Softmax Function, each output value is converted into a probability that sums up to 1, representing the likelihood of each class label being the correct one.

Softmax Function Formula

$$\text{softmax}(z) = \frac{e^{z(i)}}{\sum_{j=0}^k e^{z(j)}}$$

Consider a scenario in which z is a vector of inputs with a length equal to the number of classes k . To gain a better understanding of the `softmax` function, let's explore an example by inputting a vector of numbers.

For instance, let $z=[1,3,4,7]$.

If we aim to compute the probability for the second entry, which is 3, we can do so by substituting our desired values into the formula.

$$\text{softmax}(z_2) = \frac{e^3}{e^1 + e^3 + e^4 + e^7} = 0.017$$

By using the `softmax` function on all elements in vector z , we obtain the resulting vector that adds up to 1:

$$\text{softmax}(z) = [0.002, 0.017, 0.047, 0.934]$$

Observe that the final element has a probability exceeding 90%. In a classification scenario, this means the observation would be assigned to the final class.

Multinomial Logistic Regression

To perform multinomial logistic regression, you construct a regression model in the following format:

$$z = \beta^t x$$

and applying the softmax function to it:

$$\hat{y} = \text{softmax}(\beta^t x)$$

Multinomial Logistic Regression Loss Function

$$\text{Cost}(\beta) = - \sum_{i=1}^k y_i \log(\hat{y}_i)$$

Let's demonstrate this with an example:

Consider a scenario where you need to categorize fruits into three categories, and the fruit in question is a banana. To simplify the discussion, we will focus on a single observation. The vector y can be represented as follows:

$$y = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} \textit{apple} \\ \textit{banana} \\ \textit{orange} \end{bmatrix}$$

Once the logistic regression model is trained on your dataset, it generates a prediction vector containing probabilities.

$$\hat{y} = \begin{bmatrix} 0.2 \\ 0.7 \\ 0.1 \end{bmatrix}$$

Now, we plug this into our cost function:

$$Cost(\beta) = -(0 \times \log(0.2) + 1 \times \log(0.7) + 0 \times \log(0.1))$$

One advantageous feature of the Softmax function is that when the entries in y are 0, all terms unrelated to the true class will effectively vanish. This property simplifies the calculation by focusing solely on the relevant class, which is especially useful for optimizing and evaluating classification models.

$$Cost(\beta) = -\log(0.7) = 0.36$$

In minimizing costs effectively, this function operates based on the probability y^{\wedge} . A higher probability y^{\wedge} associated with the true class leads to lower costs. The gradient of the cost function with respect to each parameter β_j is found by calculating the first derivative of the cost with respect to β_j , which is a straightforward process. It turns out to be:

$$\frac{\partial Cost(\beta)}{\partial \beta_j} = \hat{y} - y$$

To conclude, we can now apply gradient descent to iteratively minimize the cost. This is achieved by updating the parameters in the direction that reduces the cost, scaled by a learning rate α . This process gradually improves the model's performance by making incremental adjustments to the parameters based on the computed gradients.

$$Cost(\beta) = Cost(\beta) - \alpha \frac{\partial Cost(\beta)}{\partial \beta_j}$$

Hierarchical Softmax for Many Classes:

What is Hierarchical Softmax?

Hierarchical Softmax is an efficient approximation of the softmax function, particularly useful when dealing with a very large number of classes K . It is designed to speed up computations and reduce memory requirements for models with extensive output spaces, such as those encountered in natural language processing with large vocabularies.

Concept of Hierarchical Softmax

In standard softmax, calculating the probability distribution over a large number of classes involves computing exponentials and summing them up, which can be computationally intensive. Hierarchical Softmax mitigates this issue by organizing classes into a binary tree structure, thus reducing the number of required computations.

Structure

1. Binary Tree:

- Classes are arranged in a binary tree with K leaves, where each leaf represents a class.
- Internal nodes of the tree act as decision points that split the classes into two groups.

2. Path from Root to Leaf:

- Each class is linked to a unique path from the root of the tree to its corresponding leaf node.
- The path involves making binary decisions (e.g., left or right) at each internal node, leading to the final class prediction.

Backpropagated Saliency for Feature Selection:

Backpropagated Saliency is a technique used to interpret and understand the importance of features in a neural network by analyzing the gradients of the output with respect to the input features. It provides insights into which features have the most significant impact on the model's predictions, aiding in feature selection and model interpretation.

Concept of Backpropagated Saliency

The core concept involves using the gradients of the output (or loss) with respect to the input features to gauge feature importance. This method measures how changes in input features influence the model's output, thereby highlighting the most influential features.

Steps in Backpropagated Saliency

1. Forward Pass:

- Compute the network's output for a given input. This involves passing the input through the network to obtain predictions or scores.

2. Backward Pass:

- Calculate the gradients of the output (e.g., class score) with respect to each input feature. This is done through backpropagation, where gradients are computed layer by layer to determine how sensitive the output is to changes in each input feature.

3. Saliency Map Generation:

- Generate a saliency map by visualizing the absolute values of these gradients. The saliency map highlights which features (or parts of the input) have the most significant influence on the output.

Autoencoders:

Autoencoders are a type of artificial neural network used mainly for unsupervised learning. They aim to learn efficient data representations by encoding input data into a

lower-dimensional space and then reconstructing it as accurately as possible. They are commonly applied for dimensionality reduction, feature learning, and data denoising.

Components of Autoencoders

1. Encoder:

- **Purpose:** Compresses the input data into a lower-dimensional representation, known as the bottleneck or latent space.
- **Architecture:** Comprises one or more layers that transform the input data into a compact representation.

2. Latent Space (or Bottleneck):

- **Purpose:** Holds the compressed representation of the input data.
- **Characteristics:** A lower-dimensional space capturing the essential features of the data.

3. Decoder:

- **Purpose:** Reconstructs the original input data from the compressed representation.
- **Architecture:** Consists of one or more layers that transform the latent space representation back into the original data space.

4. Loss Function:

- **Purpose:** Measures the difference between the original input and the reconstructed output.
- **Common Choices:** Mean Squared Error (MSE) for continuous data, Binary Cross-Entropy for binary data.

How Autoencoders Work

1. Forward Pass:

- Input data is fed through the encoder to create a latent space representation.
- This representation is then passed through the decoder to reconstruct the original data.

2. Reconstruction Loss:

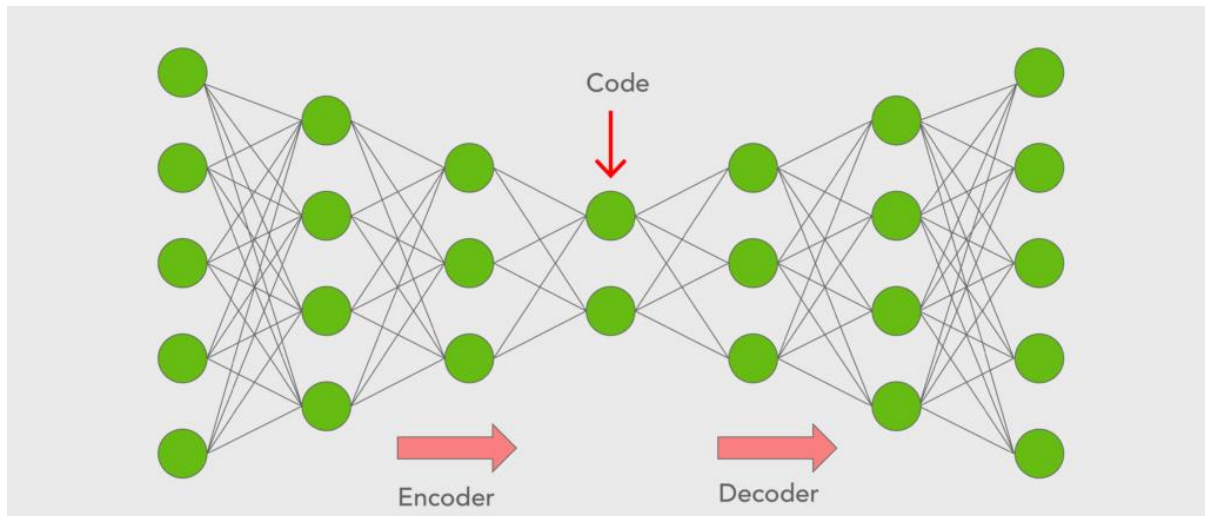
- The loss is computed by comparing the original input with the reconstructed output.
- The network is trained to minimize this loss, improving the quality of reconstruction.

3. Training:

- The network adjusts its weights using gradient descent or other optimization methods to minimize reconstruction loss.
- Both the encoder and decoder are updated simultaneously to enhance reconstruction accuracy.

Architecture of Autoencoders:

- **Encoder:** A fully connected neural network that compresses the input data into a latent space representation. It encodes the input into a reduced-dimensional form, which can be a distorted version of the original data.
- **Code:** This component stores the reduced representation of the input data, which is then passed to the decoder for reconstruction.
- **Decoder:** A feedforward neural network that mirrors the encoder's structure. It reconstructs the input data from the encoded representation, aiming to return it to its original dimensions.



Deep Autoencoders:

Deep Autoencoders extend the basic autoencoder architecture by adding more layers to both the encoder and decoder networks. This increased depth allows for learning more complex and hierarchical representations of the data, improving performance on tasks involving high-dimensional data or intricate patterns.

Architecture of Deep Autoencoders

1. Deep Encoder:

- Description: Multiple layers progressively transform the input data into a lower-dimensional representation. Each layer extracts increasingly abstract features from the input.
- Structure: Typically includes several fully connected (dense) layers. Convolutional layers may also be used for spatial data like images.

2. Latent Space:

- Description: The compressed representation of the input data after passing through the deep encoder. This space captures the essential features needed for accurate reconstruction.

3. Deep Decoder:

- Description: Multiple layers reconstruct the original input from the latent space representation. The decoder mirrors the encoder's architecture but in reverse, expanding the compressed representation back to the original data dimensions.
- Structure: Can include fully connected layers or convolutional layers, depending on the data type.

Training Deep Autoencoders

1. Forward Pass:

- Input data is passed through the deep encoder to obtain the latent space representation.
- The latent representation is then passed through the deep decoder to reconstruct the original input.

2. Reconstruction Loss:

- The difference between the original input and the reconstructed output is measured using a loss function, such as Mean Squared Error (MSE) or Binary Cross-Entropy.
- The network is trained to minimize this reconstruction loss by adjusting the weights of the encoder and decoder.

3. Optimization:

- Gradient-based optimization methods (e.g., stochastic gradient descent, Adam) are used to update the network weights and minimize the reconstruction loss.

Advantages of Deep Autoencoders

1. Hierarchical Feature Learning:

- The deep architecture enables learning hierarchical features, where each layer captures different levels of abstraction from the data.

2. Enhanced Representation Power:

- Deeper networks can capture more complex patterns and relationships within the data, leading to improved performance on tasks like dimensionality reduction and feature extraction.

3. Versatility:

- Deep autoencoders can be adapted for various data types, including images, text, and structured data, by using appropriate network architectures (e.g., convolutional layers for images, recurrent layers for text).

Application to Outlier Detection

Deep autoencoders can be particularly effective for detecting outliers or anomalies in high-dimensional spaces where traditional methods might struggle. Here's how they work for outlier detection:

1. Training the Autoencoder:

- Data Preparation: Train the autoencoder on a dataset where most data points are normal (non-anomalous).
- Encoding and Decoding: The autoencoder learns to encode normal data into a latent space and then decode it back to the original data space, minimizing reconstruction error.

2. Reconstruction Error:

- Definition: The difference between the original input and the reconstructed output from the autoencoder.
- Calculation: Measured using metrics such as Mean Squared Error (MSE) or Binary Cross-Entropy.

3. Detection of Outliers:

- Normal Data: For data similar to the training data, the autoencoder should reconstruct the input with low error.

- Anomalous Data: Outliers tend to have higher reconstruction errors because they are not well-represented in the learned latent space.
4. Threshold Setting:
- Determining Outliers: Set a threshold for the reconstruction error. Data points with errors exceeding this threshold are considered outliers.
 - Threshold Selection: Can be based on statistical methods (e.g., percentiles of the reconstruction error distribution) or empirical evaluation.

Advantages

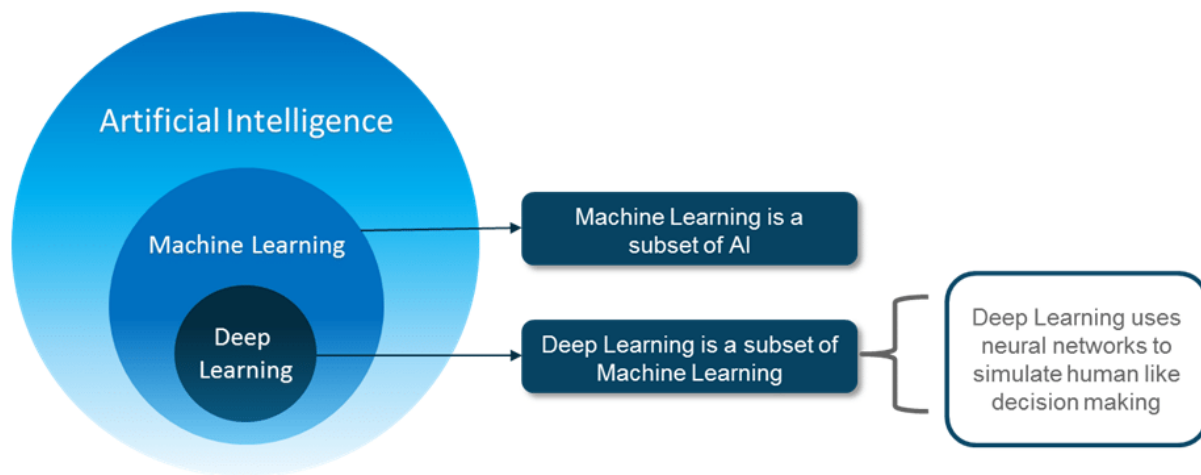
1. Handling High Dimensionality:
 - Deep autoencoders can manage high-dimensional data by learning compressed representations, making them suitable for complex datasets.
2. Learning Complex Patterns:
 - The deep architecture enables the model to capture intricate patterns and dependencies in the data, leading to more accurate anomaly detection.
3. Versatility:
 - Applicable to various data types, including images, time series, and structured data, by employing appropriate network architectures.

MODULE-I

MACHINE LEARNING WITH SHALLOW NEURAL NETWORKS

Deep Learning

Deep learning, a subset of machine learning, utilizes algorithms inspired by the human brain's structure and functioning. It employs artificial neural networks to build sophisticated models capable of addressing complex challenges. This technique is often applied to process and analyze unstructured data.



Artificial Intelligence (AI) refers to technology that enables machines to mimic human behavior and make independent decisions. This concept, both simple and intriguing, has long sparked debate. For years, many believed that computers could never match the complex capabilities of the human brain.

Historically, AI's development was hindered by limited data and insufficient computing power. However, with the rise of Big Data and the advancement of GPU technology, AI has evolved from a theoretical idea into a practical reality. **Machine Learning (ML)**, a subset of AI, uses statistical techniques to help machines improve performance by learning from past experiences.

Deep Learning (DL), a specialized branch of ML, allows the training of deep neural networks with multiple layers. By leveraging these neural networks, deep learning enables machines to replicate human decision-making processes.

The Need for Deep Learning

Traditional machine learning algorithms struggle with high-dimensional data—data with thousands of inputs and outputs across many dimensions. Handling such data is labor-intensive and computationally expensive, a challenge known as the **Curse of Dimensionality**. Deep learning, however, excels at managing high-dimensional data, efficiently identifying relevant features without human intervention.

What is Deep Learning?

Deep Learning is a branch of Machine Learning that utilizes similar algorithms to train deep neural networks, significantly improving accuracy. It mimics the cognitive processes of the human brain by learning through experience. Just as the brain's billions of neurons allow us to solve complex problems, deep learning uses artificial neurons (perceptrons) and artificial neural networks to replicate this cognitive power in machines.

What is a Neural Network?

A neural network, inspired by the biological neuron, is a collection of interconnected neurons that perform basic functions in our everyday activities. Examples include:

- A child learning to recognize an apple
- An animal distinguishing between its owner and others
- Sensing whether an object is hot or cold

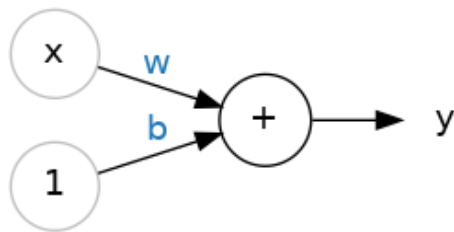
These tasks are handled by our brain's neural networks. Scientists have developed artificial neural networks (ANNs) that simulate these processes.

The Perceptron Model

An artificial neuron, or perceptron, is a supervised learning algorithm that classifies data into two categories, making it a binary classifier.

The Linear Unit:

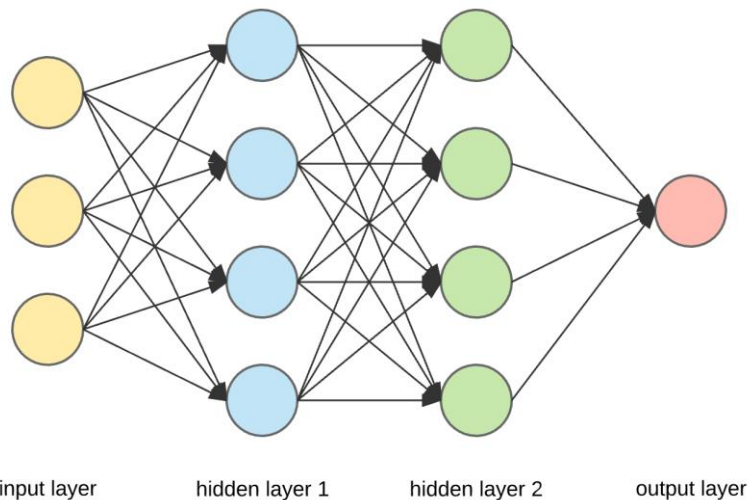
"The fundamental building block of a neural network is the individual neuron. Below is a basic illustration of a neuron (or unit) with a single input:"



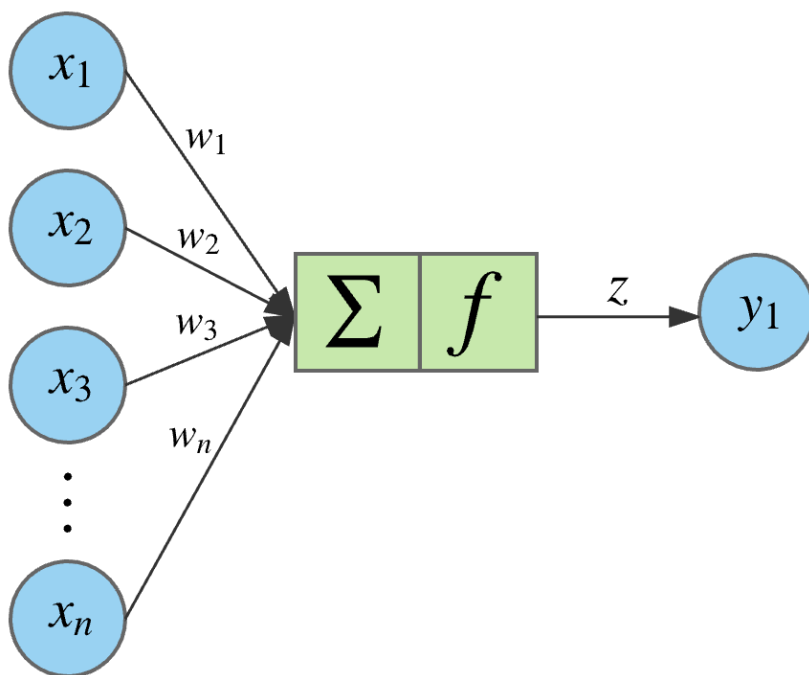
The Linear Unit: $y = wx + b$

- The neuron receives an input, denoted as x , connected by a weight w . When this input passes through the connection, it is multiplied by the weight, resulting in $w \times x$ reaching the neuron.
- The neural network improves its performance by adjusting these weights.
- Additionally, a special weight called the bias, denoted as b , is introduced. Unlike other weights, the bias does not correspond to any input. Instead, a fixed value of 1 is used, allowing the neuron to receive only the bias (since $1 \times b = b$).
- The bias enables the neuron to adjust its output independently of its inputs.
- The output of the neuron, represented as y , is calculated by summing all incoming values through its connections. The neuron's activation can be expressed using the formula $y = wx + b$, where w is the weight and x is the input.
- Are you familiar with the formula $y = wx + b$? This is the equation of a line, specifically the slope-intercept form, where w represents the slope and b denotes the y -intercept.

What is an Artificial Neural Network Model? An artificial neural network (ANN) is a type of neural network that consists of several layers: an input layer, one or more hidden layers, and an output layer.



Upon closer observation, it becomes evident that each node in a layer is intricately connected to every node in the adjacent layer. As more hidden layers are added, the network's depth increases. Let's take a deeper look at the characteristics of a single node within the output or hidden layers.



This node receives multiple inputs, computes the weighted sum, and then processes this sum through a non-linear activation function. The resulting output is passed as the input to the nodes in the next layer. It's important to note that the signal flows sequentially from left to right across the network. After all nodes process the data, the final output is produced.

Below is the equation that represents the operation within a node:

$$z = f(b + x \cdot w) = f \left(b + \sum_{i=1}^n x_i w_i \right)$$

$$x \in d_{1 \times n}, w \in d_{n \times 1}, b \in d_{1 \times 1}, z \in d_{1 \times 1}$$

The equation mentioned above indicates that b represents the bias, which serves as an input to all nodes and consistently takes a value of 1. The bias plays a key role in shifting the output of the activation function either to the left or the right.

Glossary of Artificial Neural Network Model

In this section, we'll introduce some fundamental terms that are essential for understanding an artificial neural network model.

Inputs:

The data initially fed into a neural network is called the "input." These inputs provide the network with the information it needs to make decisions or predictions. Typically, neural networks are designed to accept sets of real values as inputs, which are then passed to neurons in the input layer.

Training Set:

A training set consists of input data with known outputs. It is used to train the neural network, helping it learn and memorize the correct outputs for specific inputs.

Outputs:

A neural network generates outputs based on the data it processes. These outputs are usually real values or Boolean decisions produced by the neurons in the output layer.

Neuron:

The basic unit of a neural network is the neuron, also called a perceptron. Each neuron takes input values, processes them using a non-linear activation function, and generates an output. The output is then passed to the next layer of neurons. Common activation functions include TanH, Sigmoid, and ReLU, which are crucial for introducing non-linearity and improving the network's training efficiency.

Weight Space:

Each neuron in a neural network has a numerical weight. These weights are adjusted during the training process to produce the desired output. Training the network involves fine-tuning these weights and biases, often using a technique called backpropagation, to optimize performance.

Advantages of Artificial Neural Networks (ANN)

- **Parallel Processing Capability:**

Artificial neural networks can perform multiple tasks simultaneously, thanks to their ability to process data in parallel.

- **Distributed Data Storage:**

Unlike traditional programming, ANNs store data across the entire network, not in a centralized database. This distributed storage allows the network to function even if some data is lost in one location.

- **Ability to Operate with Limited Knowledge:**

After training, an ANN can still produce outputs even when incomplete data is provided. The performance impact depends on the importance of the missing data.

- **Memory Distribution:**

ANNs can adapt based on examples provided during training. The network learns from these examples and adjusts its behavior accordingly. The quality of the examples directly influences the performance, as incomplete or inaccurate data can lead to incorrect outputs.

- **Fault Tolerance:**

ANNs are resilient and can continue to function even if some of their components fail, thanks to their inherent fault tolerance.

Disadvantages of Artificial Neural Networks (ANN)

- **Determining the Correct Network Structure:**

There is no fixed guideline for determining the optimal structure of an ANN. Achieving the right structure often requires experimentation, practical knowledge, and learning from trial and error.

- **Unclear Network Behavior:**

A significant challenge with ANNs is their often opaque behavior. When they produce a

result, they typically do not provide an explanation for their decisions, which can reduce confidence in their outcomes.

- **Hardware Dependency:**

ANNs rely heavily on hardware with parallel processing capabilities to function effectively. This makes them dependent on specific equipment requirements.

- **Data Representation Challenges:**

ANNs process numerical data, so any problem must be translated into numbers before being input into the network. The way this data is represented directly affects the network's performance, and its success depends on the user's skill in converting the problem into an appropriate numerical format.

- **Uncertainty in Training Duration:**

The time required to train an ANN is often unpredictable and can be constrained by a predefined error threshold, which may not always yield the best results.

Neural Architectures for Binary Classification Models:

Least Square Regression Method:

Least squares is a popular technique in regression analysis used to estimate unknown parameters by constructing a model that minimizes the sum of the squared differences between observed and predicted values. One of the most commonly used methods for curve fitting involves reducing the sum of squared errors to minimize discrepancies. This approach helps generate the best-fitting line tailored to your specific data points.

Finding the Line of Best Fit Using Least Square Regression

When analyzing a set of paired numbers and its corresponding scatter plot, the line of best fit is the straight line that can be drawn through the scatter points to represent the relationship between them most accurately. The equation of the straight line is represented as:

$$Y = mX + c$$

where:

Y: Dependent Variable

m: Slope

X: Independent Variable

The goal is to find the values of the slope and y-intercept, which are then substituted into the equation along with the independent variable X to compute the dependent variable Y. For a dataset with 'n' data points, the slope can be determined using the following formula:.

$$m = \frac{\sum(x - \bar{x})(y - \bar{y})}{\sum(x - \bar{x})^2}$$

Then, the y-intercept is calculated using the formula:

$$c = \bar{y} - m * \bar{x}$$

At last, we substitute these values in the final equation

$$Y = mX + c.$$

Logistic regression:

Logistic Regression is a core machine learning algorithm primarily used for binary classification tasks. While it's simpler than more complex deep learning models, it remains relevant and frequently serves as a baseline or a foundational component in deep learning workflows. Below is an overview of logistic regression within the context of deep learning:

Logistic Regression Overview

Objective: The primary goal of logistic regression is to estimate the probability that a given input belongs to a specific class. It is typically used for binary classification problems, where there are two possible outcomes.

$$y' = p(y = 1|x)$$

Mathematical Model

Input: X is a matrix with dimensions $n \times m$, where n represents the number of features and m denotes the number of training examples.

Parameters: W represents a Weight Matrix with dimensions $n \times 1$, where n signifies the number of features within the dataset X . The Bias b serves the purpose of regulating the threshold value at which the activation function is activated.

Output:

$$y' = \sigma(W^T X + b)$$

Activation Function

Activation Functions are essential components in neural networks that introduce non-linearity, enabling the model to learn and capture complex patterns and relationships in the data. Without these functions, a neural network would essentially be limited to linear modeling, restricting its ability to solve more intricate problems.

Key Concepts of Activation Functions:

2. Purpose:

- Non-linearity: Activation functions enable neural networks to model non-linear relationships between inputs and outputs.
- Learning Complex Patterns: By applying non-linear transformations across multiple layers, activation functions help the network learn and represent intricate patterns in the data.

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

Loss Function

Loss Functions are critical elements in machine learning and deep learning that assess how well a model's predictions align with the actual target values. They quantify the discrepancy between predicted and true values, providing a measure of error that guides the optimization process during model training.

$$L(y', y) = -[y \log y' + (1 - y) \log(1 - y')]$$

Let's examine the effectiveness of this loss function for logistic regression:

1. When $y=1$, the loss function is represented as $L(y', y) = -\log y'$. To minimize this loss function, the value of $\log y'$ should be maximized. This can be achieved when y' approaches 1, resulting in predicted values that closely align with actual values.
2. When the loss function is calculated at $y=0$, it results in $L(y', y) = -\log(1-y')$. In order to minimize the loss function, the value of $\log(1-y')$ should be maximized. This can be achieved by reducing the value of y' , making it closer to 0. By doing so, the predicted value will be more aligned with the actual value.
3. The loss function described above exhibits convexity, ensuring the presence of a sole global minimum and preventing the network from getting trapped in local minima commonly found in non-convex loss functions.

Widrow–Hoff Rule or Delta Learning Rule:

The Least Mean Squares (LMS) method, introduced by Bernard Widrow and Marcian Hoff, is designed to minimize errors across all training patterns in a supervised learning setting. This algorithm employs a continuous activation function.

At its core, the LMS method is based on the gradient descent approach, which iteratively adjusts the synaptic weights to reduce the error between the predicted and desired output values. The delta rule is used to adjust these weights, progressively decreasing the overall input to the output unit to better align with the target value.

To adjust the synaptic weights, the delta rule can be expressed as:

$$\Delta w_i = \alpha \cdot x_i \cdot e_j$$

where:

Δw_i represents the change in weight for the i th pattern;

α is a positive and constant learning rate;

x_i is the input value from the pre-synaptic neuron;

e_j is the difference between the desired/target output and the actual output y_{in} ($e_j = t - y_{in}$).

It is important to note that the above delta rule is intended for a single output unit only.

The updating of weight can be done in the following two cases –

Closed - form solutions:

The closed-form solution is preferable for smaller datasets if the cost of computing a matrix inverse is not a concern. For larger datasets, or those where the inverse of $\mathbf{X}^T \mathbf{X}$ may not exist due to perfect multicollinearity, gradient descent (GD) or stochastic gradient descent (SGD) approaches are recommended. The linear regression model is defined as:

$$y = w_0 x_0 + w_1 x_1 + \dots + w_m x_m = \sum_{j=0}^m \mathbf{w}^T \mathbf{x}$$

The response variable is denoted as y , the sample vector with m dimensions is represented by \mathbf{x} , and the weight vector (a vector of coefficients) is symbolized as \mathbf{w} . It is important to acknowledge that w_0 signifies the intercept on the y -axis in the model, thus x_0 is equivalent to 1. By employing the closed-form solution, which is the normal equation method, we determine the weights of the model subsequently:

$$\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

Neural Architectures for Multiclass Models:

Multilayer Perceptron (MLP)

A Multilayer Perceptron (MLP) is a type of artificial neural network with multiple layers of neurons. It is one of the simplest deep learning models and is utilized for various predictive tasks, including classification and regression.

Structure of MLP

5. Layers:

- Input Layer: This is the first layer where the data is introduced into the network. Each neuron in this layer represents a feature from the input data.
- Hidden Layers: These are intermediate layers located between the input and output layers. Neurons in hidden layers process inputs through weighted connections and activation functions.
- Output Layer: This is the final layer that generates the model's predictions. For classification tasks, the output layer typically uses a softmax or sigmoid activation function. For regression tasks, a linear activation function is often employed.

6. Neurons:

- Each neuron in a layer receives inputs, calculates a weighted sum, adds a bias term, and then applies an activation function to produce an output.

7. Weights and Biases:

- Weights: These are parameters learned during training that control the strength of the connections between neurons.
- Biases: These are additional parameters added to the weighted sum before applying the activation function, which allows the model to fit the data more flexibly.

8. Activation Functions:

- Purpose: Activation functions introduce non-linearity into the model, enabling it to learn complex patterns.
- Common Functions: ReLU (Rectified Linear Unit), Sigmoid, Tanh.

How MLP Works

5. Forward Propagation:

- Input data is passed through the network, layer by layer.
- Each neuron computes a weighted sum of its inputs, adds a bias, and applies an activation function.
- The output of each layer serves as the input to the next layer, culminating in the final output.

6. Loss Calculation:

- The predicted output is compared to the true target using a loss function (e.g., cross-entropy loss for classification or mean squared error for regression).
- The loss function quantifies the error between the predicted values and the actual values.

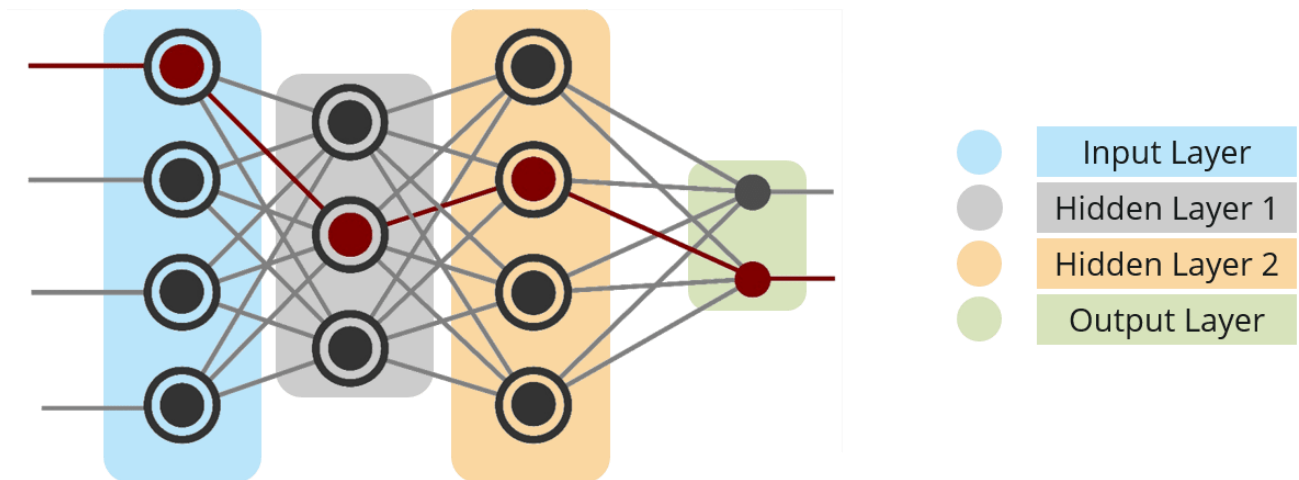
7. Backward Propagation:

- Gradient Computation: Gradients of the loss function with respect to the weights and biases are computed using the chain rule.
- Weight Update: Weights and biases are updated using an optimization algorithm (e.g., gradient descent or Adam) to minimize the loss function.

8. Training:

- The network is trained through multiple iterations (epochs) where forward and backward propagation are repeated.
- The optimization algorithm adjusts the weights and biases to enhance the model's performance.

Consider the diagram below:



Weston-Watkins SVM (or) Multi-class SVM Loss:

Weston-Watkins SVM (also known as Multi-class SVM or Weston-Watkins loss) is a variant of the Support Vector Machine (SVM) tailored for multi-class classification problems. Unlike binary SVMs, which are designed to separate data into two classes, multi-class SVMs extend this concept to handle scenarios involving more than two classes.

Weston-Watkins SVM (Multi-class SVM) Overview

Objective: The primary goal is to develop a model capable of classifying input data into one of several classes. This is achieved by minimizing a loss function that quantifies the discrepancy between predicted and actual class labels.

Weston-Watkins Loss Function

3. **Formulation:** The Weston-Watkins loss function extends the binary SVM loss function to accommodate multi-class classification. It introduces a structured loss that penalizes misclassifications by considering the relationships among multiple classes.
4. **Loss Function:** For a given data point, the Weston-Watkins loss function can be expressed as:

$$L(\mathbf{w}, \mathbf{x}, y) = \sum_{k \neq y} \max(0, 1 - (\mathbf{w}_y^T \mathbf{x} - \mathbf{w}_k^T \mathbf{x}))$$

where:

- \mathbf{x} is the input feature vector.
- y is the true class label.
- \mathbf{w}_y and \mathbf{w}_k are the weight vectors corresponding to the true class and the incorrect class k , respectively.

Softmax regression:

What is the Softmax Function?

The Softmax Function is a mathematical function commonly used in multi-class classification problems within machine learning and deep learning. It transforms raw scores (also known as logits) output by the model into probabilities, allowing for a probabilistic interpretation of the neural network's output. By applying the Softmax Function, each output value is converted into a probability that sums up to 1, representing the likelihood of each class label being the correct one.

Softmax Function Formula

$$\text{softmax}(z) = \frac{e^{z(i)}}{\sum_{j=0}^k e^{z(j)}}$$

Consider a scenario in which z is a vector of inputs with a length equal to the number of classes k . To gain a better understanding of the `softmax` function, let's explore an example by inputting a vector of numbers.

For instance, let $z=[1,3,4,7]$.

If we aim to compute the probability for the second entry, which is 3, we can do so by substituting our desired values into the formula.

$$\text{softmax}(z_2) = \frac{e^3}{e^1 + e^3 + e^4 + e^7} = 0.017$$

By using the `softmax` function on all elements in vector z , we obtain the resulting vector that adds up to 1:

$$\text{softmax}(z) = [0.002, 0.017, 0.047, 0.934]$$

Observe that the final element has a probability exceeding 90%. In a classification scenario, this means the observation would be assigned to the final class.

Multinomial Logistic Regression

To perform multinomial logistic regression, you construct a regression model in the following format:

$$z = \beta^t x$$

and applying the softmax function to it:

$$\hat{y} = \text{softmax}(\beta^t x)$$

Multinomial Logistic Regression Loss Function

$$\text{Cost}(\beta) = - \sum_{i=1}^k y_i \log(\hat{y}_i)$$

Let's demonstrate this with an example:

Consider a scenario where you need to categorize fruits into three categories, and the fruit in question is a banana. To simplify the discussion, we will focus on a single observation. The vector y can be represented as follows:

$$y = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} \textit{apple} \\ \textit{banana} \\ \textit{orange} \end{bmatrix}$$

Once the logistic regression model is trained on your dataset, it generates a prediction vector containing probabilities.

$$\hat{y} = \begin{bmatrix} 0.2 \\ 0.7 \\ 0.1 \end{bmatrix}$$

Now, we plug this into our cost function:

$$Cost(\beta) = -(0 \times \log(0.2) + 1 \times \log(0.7) + 0 \times \log(0.1))$$

One advantageous feature of the Softmax function is that when the entries in y are 0, all terms unrelated to the true class will effectively vanish. This property simplifies the calculation by focusing solely on the relevant class, which is especially useful for optimizing and evaluating classification models.

$$Cost(\beta) = -\log(0.7) = 0.36$$

In minimizing costs effectively, this function operates based on the probability y^{\wedge} . A higher probability y^{\wedge} associated with the true class leads to lower costs. The gradient of the cost function with respect to each parameter β_j is found by calculating the first derivative of the cost with respect to β_j , which is a straightforward process. It turns out to be:

$$\frac{\partial Cost(\beta)}{\partial \beta_j} = \hat{y} - y$$

To conclude, we can now apply gradient descent to iteratively minimize the cost. This is achieved by updating the parameters in the direction that reduces the cost, scaled by a learning rate α . This process gradually improves the model's performance by making incremental adjustments to the parameters based on the computed gradients.

$$Cost(\beta) = Cost(\beta) - \alpha \frac{\partial Cost(\beta)}{\partial \beta_j}$$

Hierarchical Softmax for Many Classes:

What is Hierarchical Softmax?

Hierarchical Softmax is an efficient approximation of the softmax function, particularly useful when dealing with a very large number of classes K . It is designed to speed up computations and reduce memory requirements for models with extensive output spaces, such as those encountered in natural language processing with large vocabularies.

Concept of Hierarchical Softmax

In standard softmax, calculating the probability distribution over a large number of classes involves computing exponentials and summing them up, which can be computationally intensive. Hierarchical Softmax mitigates this issue by organizing classes into a binary tree structure, thus reducing the number of required computations.

Structure

3. Binary Tree:

- Classes are arranged in a binary tree with K leaves, where each leaf represents a class.
- Internal nodes of the tree act as decision points that split the classes into two groups.

4. Path from Root to Leaf:

- Each class is linked to a unique path from the root of the tree to its corresponding leaf node.
- The path involves making binary decisions (e.g., left or right) at each internal node, leading to the final class prediction.

Backpropagated Saliency for Feature Selection:

Backpropagated Saliency is a technique used to interpret and understand the importance of features in a neural network by analyzing the gradients of the output with respect to the input features. It provides insights into which features have the most significant impact on the model's predictions, aiding in feature selection and model interpretation.

Concept of Backpropagated Saliency

The core concept involves using the gradients of the output (or loss) with respect to the input features to gauge feature importance. This method measures how changes in input features influence the model's output, thereby highlighting the most influential features.

Steps in Backpropagated Saliency

4. Forward Pass:

- Compute the network's output for a given input. This involves passing the input through the network to obtain predictions or scores.

5. Backward Pass:

- Calculate the gradients of the output (e.g., class score) with respect to each input feature. This is done through backpropagation, where gradients are computed layer by layer to determine how sensitive the output is to changes in each input feature.

6. Saliency Map Generation:

- Generate a saliency map by visualizing the absolute values of these gradients. The saliency map highlights which features (or parts of the input) have the most significant influence on the output.

Autoencoders:

Autoencoders are a type of artificial neural network used mainly for unsupervised learning. They aim to learn efficient data representations by encoding input data into a

lower-dimensional space and then reconstructing it as accurately as possible. They are commonly applied for dimensionality reduction, feature learning, and data denoising.

Components of Autoencoders

5. Encoder:

- **Purpose:** Compresses the input data into a lower-dimensional representation, known as the bottleneck or latent space.
- **Architecture:** Comprises one or more layers that transform the input data into a compact representation.

6. Latent Space (or Bottleneck):

- **Purpose:** Holds the compressed representation of the input data.
- **Characteristics:** A lower-dimensional space capturing the essential features of the data.

7. Decoder:

- **Purpose:** Reconstructs the original input data from the compressed representation.
- **Architecture:** Consists of one or more layers that transform the latent space representation back into the original data space.

8. Loss Function:

- **Purpose:** Measures the difference between the original input and the reconstructed output.
- **Common Choices:** Mean Squared Error (MSE) for continuous data, Binary Cross-Entropy for binary data.

How Autoencoders Work

4. Forward Pass:

- Input data is fed through the encoder to create a latent space representation.
- This representation is then passed through the decoder to reconstruct the original data.

5. Reconstruction Loss:

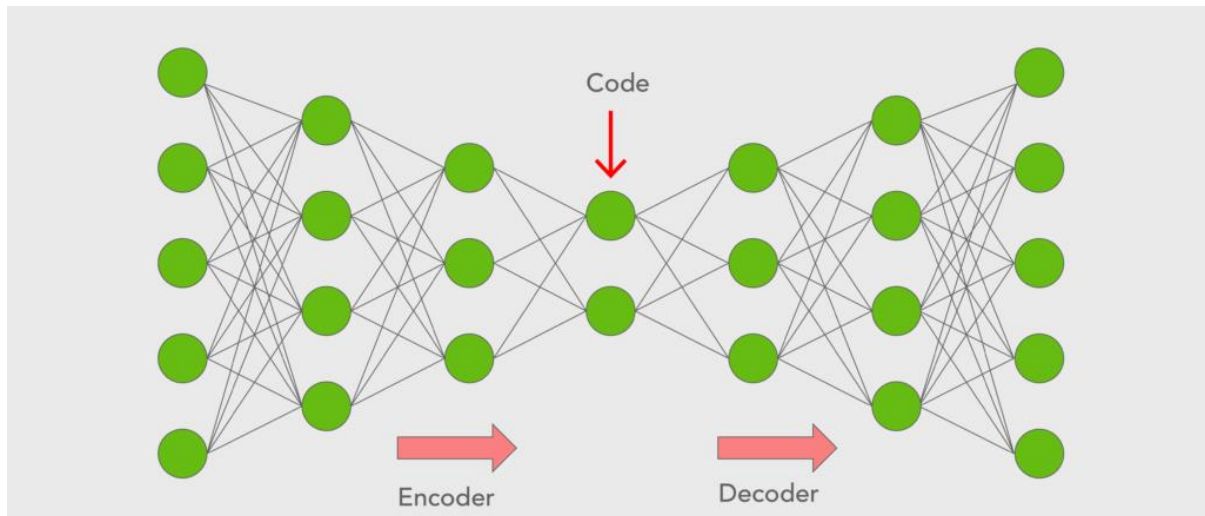
- The loss is computed by comparing the original input with the reconstructed output.
- The network is trained to minimize this loss, improving the quality of reconstruction.

6. Training:

- The network adjusts its weights using gradient descent or other optimization methods to minimize reconstruction loss.
- Both the encoder and decoder are updated simultaneously to enhance reconstruction accuracy.

Architecture of Autoencoders:

- **Encoder:** A fully connected neural network that compresses the input data into a latent space representation. It encodes the input into a reduced-dimensional form, which can be a distorted version of the original data.
- **Code:** This component stores the reduced representation of the input data, which is then passed to the decoder for reconstruction.
- **Decoder:** A feedforward neural network that mirrors the encoder's structure. It reconstructs the input data from the encoded representation, aiming to return it to its original dimensions.



Deep Autoencoders:

Deep Autoencoders extend the basic autoencoder architecture by adding more layers to both the encoder and decoder networks. This increased depth allows for learning more complex and hierarchical representations of the data, improving performance on tasks involving high-dimensional data or intricate patterns.

Architecture of Deep Autoencoders

4. Deep Encoder:

- Description: Multiple layers progressively transform the input data into a lower-dimensional representation. Each layer extracts increasingly abstract features from the input.
- Structure: Typically includes several fully connected (dense) layers. Convolutional layers may also be used for spatial data like images.

5. Latent Space:

- Description: The compressed representation of the input data after passing through the deep encoder. This space captures the essential features needed for accurate reconstruction.

6. Deep Decoder:

- Description: Multiple layers reconstruct the original input from the latent space representation. The decoder mirrors the encoder's architecture but in reverse, expanding the compressed representation back to the original data dimensions.
- Structure: Can include fully connected layers or convolutional layers, depending on the data type.

Training Deep Autoencoders

4. Forward Pass:

- Input data is passed through the deep encoder to obtain the latent space representation.
- The latent representation is then passed through the deep decoder to reconstruct the original input.

5. Reconstruction Loss:

- The difference between the original input and the reconstructed output is measured using a loss function, such as Mean Squared Error (MSE) or Binary Cross-Entropy.
- The network is trained to minimize this reconstruction loss by adjusting the weights of the encoder and decoder.

6. Optimization:

- Gradient-based optimization methods (e.g., stochastic gradient descent, Adam) are used to update the network weights and minimize the reconstruction loss.

Advantages of Deep Autoencoders

4. Hierarchical Feature Learning:

- The deep architecture enables learning hierarchical features, where each layer captures different levels of abstraction from the data.

5. Enhanced Representation Power:

- Deeper networks can capture more complex patterns and relationships within the data, leading to improved performance on tasks like dimensionality reduction and feature extraction.

6. Versatility:

- Deep autoencoders can be adapted for various data types, including images, text, and structured data, by using appropriate network architectures (e.g., convolutional layers for images, recurrent layers for text).

Application to Outlier Detection

Deep autoencoders can be particularly effective for detecting outliers or anomalies in high-dimensional spaces where traditional methods might struggle. Here's how they work for outlier detection:

5. Training the Autoencoder:

- Data Preparation: Train the autoencoder on a dataset where most data points are normal (non-anomalous).
- Encoding and Decoding: The autoencoder learns to encode normal data into a latent space and then decode it back to the original data space, minimizing reconstruction error.

6. Reconstruction Error:

- Definition: The difference between the original input and the reconstructed output from the autoencoder.
- Calculation: Measured using metrics such as Mean Squared Error (MSE) or Binary Cross-Entropy.

7. Detection of Outliers:

- Normal Data: For data similar to the training data, the autoencoder should reconstruct the input with low error.

- Anomalous Data: Outliers tend to have higher reconstruction errors because they are not well-represented in the learned latent space.

8. Threshold Setting:

- Determining Outliers: Set a threshold for the reconstruction error. Data points with errors exceeding this threshold are considered outliers.
- Threshold Selection: Can be based on statistical methods (e.g., percentiles of the reconstruction error distribution) or empirical evaluation.

Advantages

4. Handling High Dimensionality:

- Deep autoencoders can manage high-dimensional data by learning compressed representations, making them suitable for complex datasets.

5. Learning Complex Patterns:

- The deep architecture enables the model to capture intricate patterns and dependencies in the data, leading to more accurate anomaly detection.

6. Versatility:

- Applicable to various data types, including images, time series, and structured data, by employing appropriate network architectures.