



Neural Networks - Fundamentals and Multi-Layer Perceptrons (MLPs)

This guide will cover:

1. Fundamentals of Neural Networks

- Neurons, weights, and activation functions
- Forward propagation
- Loss functions
- Backpropagation
- Gradient descent

2. Multi-Layer Perceptrons (MLPs)

- Structure and layers
- Fully connected networks
- Training an MLP
- Manual implementation (fully commented)

1. Fundamentals of Neural Networks

What is a Neural Network?

A neural network is a function that **maps inputs to outputs** by simulating the behavior of biological neurons. Each neuron **takes inputs, applies a weight, adds bias, passes it through an activation function, and sends the result to the next layer**.

Basic Structure:

- **Input Layer** → Takes feature inputs.
- **Hidden Layers** → Perform computations.
- **Output Layer** → Provides final prediction.

Example:

For an image classifier:

- Input: **Pixels of an image**
- Hidden layers: **Learn features like edges, shapes, patterns**
- Output: **Predicted class (dog, cat, etc.)**

2. Neurons and Activation Functions

Neuron Computation (Perceptron)

Each neuron performs:

1. Weighted sum of inputs:

$$z = w_1x_1 + w_2x_2 + \cdots + w_nx_n + b$$

2. Apply activation function:

$$a = f(z)$$

Common Activation Functions

Function	Formula	When to Use
Sigmoid	$\sigma(z) = \frac{1}{1+e^{-z}}$	Output layer of binary classification
ReLU	$\max(0, z)$	Hidden layers (fast, avoids vanishing gradient)
Tanh	$\frac{e^z - e^{-z}}{e^z + e^{-z}}$	Hidden layers, when values should be centered around zero
Softmax	$\frac{e^{z_i}}{\sum e^{z_j}}$	Output layer for multi-class classification

3. Forward Propagation

Forward propagation moves data **through the network**:

1. Compute **weighted sum z** .
2. Apply **activation function**.
3. Pass output to the next layer.

For an MLP:

$$a^{(l+1)} = f(W^{(l)}a^{(l)} + b^{(l)})$$

4. Loss Functions

The loss function measures how far the predictions are from the actual values.

Loss Function	Formula	Used For
Mean Squared Error (MSE)	$L = \frac{1}{m} \sum (y - \hat{y})^2$	Regression tasks
Cross-Entropy Loss	$L = - \sum y \log \hat{y}$	Classification tasks

5. Backpropagation & Gradient Descent

To train the network, we use backpropagation and gradient descent:

1. Compute **loss** using forward propagation.
2. Compute **gradients** using chain rule.
3. Update weights:

$$w := w - \alpha \frac{\partial L}{\partial w}$$

Backpropagation **propagates errors backward** to update weights.

1. Universal Approximation Theorem

What It States

- A feedforward neural network with at least one hidden layer can approximate **any continuous function** given enough neurons and a **non-linear activation function**.

Why It Matters

- This explains why **deep learning works**—with enough capacity, an NN can model **any complex relationship**.
- However, just because a neural network **can** approximate a function doesn't mean it will **generalize well**—this leads us to **overfitting and regularization**.

2. Overfitting and Generalization in Neural Networks

Why Neural Networks Overfit

- **Too many parameters** → Can memorize training data instead of learning patterns.
- **Too few training samples** → High variance in learned weights.
- **Lack of regularization** → The model relies too much on the dataset and doesn't generalize.

How to Prevent Overfitting

Method	What It Does
Dropout	Randomly removes neurons during training to prevent reliance on specific pathways.
L1/L2 Regularization	Penalizes large weights to reduce model complexity (L2 = Ridge, L1 = Lasso).
Early Stopping	Stops training once validation loss stops improving.
Data Augmentation	Increases dataset diversity to improve generalization.
Batch Normalization	Helps control activation values, reducing sensitivity to small changes in inputs.

3. Optimization Methods and Training Stability

✓ Key Optimization Algorithms

Optimizer	How It Works	Pros	Cons
Gradient Descent (GD)	Uses all data to compute gradients	Stable, good for convex functions	Slow for large datasets
Stochastic Gradient Descent (SGD)	Uses a single random sample per step	Fast, less memory usage	High variance, noisy updates
Mini-Batch Gradient Descent	Uses small batches instead of full dataset	Balanced speed and accuracy	Requires tuning batch size
Adam (Adaptive Moment Estimation)	Combines momentum & adaptive learning rates	Fast, works well on noisy data	More hyperparameters, may not generalize well

🚀 Key Concepts You Should Know

- **Momentum:** Helps accelerate gradient descent in the right direction.
- **Learning Rate Scheduling:** Reducing the learning rate over time improves convergence.
- **Exploding and Vanishing Gradients:** Happens when gradients become too large (explode) or too small (vanish), affecting deep networks.

4. Activation Functions – Why They Matter

✓ Key Activation Functions

Function	Formula	When to Use	Pros	Cons
Sigmoid	$\sigma(z) = \frac{1}{1+e^{-z}}$	Binary classification (output layer)	Squashes values between (0,1)	Vanishing gradients
Tanh	$\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$	Hidden layers	Zero-centered	Still has vanishing gradients
ReLU	$f(z) = \max(0, z)$	Most hidden layers	Computationally efficient	Dead neurons (outputs stuck at 0)
Leaky ReLU	$f(z) = \max(0.01z, z)$	When ReLU suffers from dead neurons	Fixes dead neurons	Still unbounded
Softmax	$\frac{e^{z_i}}{\sum e^{z_j}}$	Multi-class classification (output layer)	Probabilities sum to 1	Computationally expensive

🔥 Why Does Activation Function Choice Matter?

- ReLU is **preferred** for hidden layers because it avoids vanishing gradients.
- Sigmoid and Tanh **work well for shallow networks** but are poor for deep ones.
- Softmax is the best choice for **multi-class classification**.

5. Neural Network Weight Initialization

✓ Why Initialization Matters

- Poor initialization can cause slow training or dead neurons.

🚀 Weight Initialization Strategies

Method	Formula	When to Use
Random Initialization	Small random values	Small networks
Xavier Initialization	$W \sim \mathcal{N}(0, \frac{1}{n})$	Sigmoid/Tanh activations
He Initialization	$W \sim \mathcal{N}(0, \frac{2}{n})$	ReLU-based networks

6. Batch Normalization (BN)

✓ What It Does

- Normalizes activations in each layer to stabilize training.
- Reduces internal covariate shift (shifting distributions during training).
- Improves gradient flow, making deep networks train faster.

🚀 When to Use It

- If training deep networks (many layers).
- If using high learning rates (BN allows higher rates without instability).
- If activation values become too large or too small.

7. Training Best Practices & Hyperparameter Tuning

✓ Hyperparameters to Tune

Hyperparameter	Effect
Learning Rate	Too high → won't converge, Too low → slow training
Batch Size	Small → noisy, Large → smoother but needs more memory
Number of Layers	More layers → better representation, but harder to train
Number of Neurons per Layer	Too few → underfitting, Too many → overfitting
Dropout Rate	Higher → more regularization, Lower → risk of overfitting

🚀 Tuning Strategies

1. **Grid Search** – Try all possible hyperparameter combinations (computationally expensive).
2. **Random Search** – Randomly pick hyperparameters, faster than Grid Search.
3. **Bayesian Optimization** – Smartly searches best hyperparameters based on past results.
4. **Gradient-Based Optimization** – AutoML techniques like **Hyperband**, **Optuna**.

8. Transfer Learning & Pretraining

✓ Why Transfer Learning Works

- Instead of training from scratch, start with a **pretrained model**.
- Fine-tune **only the last layers** for a new task.
- Saves **time and computation**.

🚀 When to Use It

- When you **don't have much data** (e.g., medical images, specialized datasets).
- When **pretrained models already exist** (e.g., ResNet, BERT for NLP).

9. Neural Network Interpretability

✓ Why It's Hard

- NNs are **black boxes**, hard to understand compared to linear models.
- We need ways to explain **why a model made a certain prediction**.

🚀 How to Interpret Neural Networks

Method	How It Helps
SHAP Values	Measures each feature's impact on predictions.
Feature Visualization	See which parts of an image activate neurons.
Saliency Maps	Highlights important pixels for a CNN.



6. Multi-Layer Perceptrons (MLPs)

An MLP is a type of neural network with:

- **Multiple hidden layers**
- **Fully connected neurons**
- **Non-linear activation functions**

It's the foundation for **deep learning models**.

Revised Pseudo Code with Gradient Descent Types

plaintext

Copy

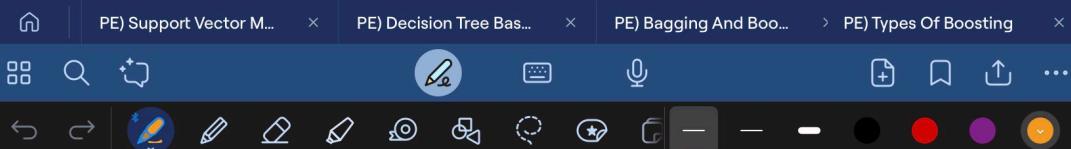
Edit

```

1 # Initialize Neural Network
2 1. Set hyperparameters (learning rate, epochs, hidden layer size, batch size).
3 2. Initialize weights (random small values) and biases (zeros).
4
5 # Training Process
6 FOR each epoch:
7   FOR each batch (if using Mini-Batch GD or SGD):
8     1. Forward Propagation:
9       a. Compute weighted sum for hidden layer:  $Z_1 = W_1 * X + b_1$ 
10      b. Apply activation function:  $A_1 = \text{activation}(Z_1)$ 
11      c. Compute weighted sum for output layer:  $Z_2 = W_2 * A_1 + b_2$ 
12      d. Apply activation function:  $A_2 = \text{activation}(Z_2)$  (final prediction)
13
14     2. Compute Loss:
15       a. Use Cross-Entropy Loss for classification
16       b. Use Mean Squared Error (MSE) for regression
17
18     3. Backpropagation:
19       a. Compute gradient of loss w.r.t.  $A_2$  (output layer).
20       b. Compute gradient of loss w.r.t.  $W_2$  and  $b_2$ .
21       c. Compute error signal backpropagated to hidden layer.
22       d. Compute gradient of loss w.r.t.  $W_1$  and  $b_1$ .
23
24     4. Update Weights (Gradient Descent Step):
25       a. **Batch Gradient Descent (BGD):** Update weights using all training samples.
26          $W_1 := W_1 - \text{learning\_rate} * dw_1$ 
27          $b_1 := b_1 - \text{learning\_rate} * db_1$ 
28          $W_2 := W_2 - \text{learning\_rate} * dw_2$ 
29          $b_2 := b_2 - \text{learning\_rate} * db_2$ 
30
31       b. **Stochastic Gradient Descent (SGD):** Update weights using one random sample at a time.
32         Select a random  $(X_i, y_i)$ 
33         Compute gradients using  $(X_i, y_i)$ 
34          $W_1 := W_1 - \text{learning\_rate} * dw_1$  (for sample i)
35          $b_1 := b_1 - \text{learning\_rate} * db_1$ 
36          $W_2 := W_2 - \text{learning\_rate} * dw_2$ 
37          $b_2 := b_2 - \text{learning\_rate} * db_2$ 
38

```

```
38
39      c. **Mini-Batch Gradient Descent:** Update weights using a small batch of samples.
40          Select a batch of `m` random samples.
41          Compute average gradients over batch.
42          W1 := W1 - learning_rate * dW1 (batch average)
43          b1 := b1 - learning_rate * db1
44          W2 := W2 - learning_rate * dW2
45          b2 := b2 - learning_rate * db2
46
47      5. Print loss every 100 epochs for monitoring.
48
49 # Making Predictions
50 1. Perform forward propagation on new data.
51 2. If using classification, convert probability to class label.
52 #K for Command, #L for Cascade
```



3. [25 points] A Simple Neural Network

Let $X = \{x^{(1)}, \dots, x^{(n)}\}$ be a dataset of n samples with 2 features, i.e. $x^{(i)} \in \mathbb{R}^2$. The samples are classified into 2 categories with labels $y^{(i)} \in \{0, 1\}$. A scatter plot of the dataset is shown in Figure 1.

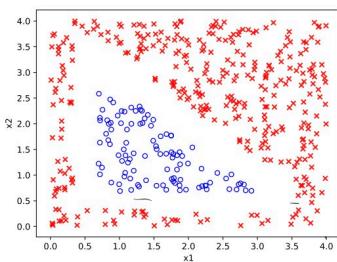


Figure 1: Plot of dataset X .

The examples in class 1 are marked as as “ \times ” and examples in class 0 are marked as “ \circ ”. We want to perform binary classification using a simple neural network with the architecture shown in Figure 2.

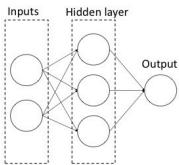


Figure 2: Architecture for our simple neural network.

Denote the two features x_1 and x_2 , the three neurons in the hidden layer h_1, h_2 , and h_3 , and the output neuron as o . Let the weight from x_i to h_j be $w_{i,j}^{[1]}$ for $i \in \{1, 2\}, j \in \{1, 2, 3\}$, and the weight from h_j to o be $w_j^{[2]}$. Finally, denote the intercept weight for h_j as $w_{0,j}^{[1]}$, and the intercept weight for o as $w_0^{[2]}$. For the loss function, we'll use average squared loss instead of the usual negative log-likelihood:

$$l = \frac{1}{n} \sum_{i=1}^n \left(o^{(i)} - y^{(i)} \right)^2,$$

where $o^{(i)}$ is the result of the output neuron for example i .

- (a) [5 points] Suppose we use the sigmoid function as the activation function for h_1, h_2, h_3 and o . What is the gradient descent update to $w_{j,2}^{[1]}$, assuming we use a learning rate of α ? Your answer should be written in terms of $x^{(i)}$, $o^{(i)}$, $y^{(i)}$, and the weights.

For hidden layer:

$$\begin{aligned} z_j^{(i)} &\approx w_{0,j} + \sum_{j=1}^3 w_{j,j}^{[1]} x^{(i)} \\ \Rightarrow h_j^{(i)} &= \sigma(z_j^{(i)}) = \frac{1}{1 + e^{-z_j^{(i)}}} \end{aligned}$$

For Output layer:

$$\begin{aligned} z^{(i)} &\approx w_{0,0}^{[0]} + \sum_{j=1}^3 w_{j,0}^{[0]} h_j^{(i)} \\ \Rightarrow o &= \sigma(z^{(i)}) = \frac{1}{1 + e^{-z^{(i)}}} \\ \ell &\leq \frac{1}{n} \sum_{i=1}^n (o^{(i)} - y^{(i)})^2 = \ell^{(i)} = (o^{(i)} - y^{(i)})^2 \end{aligned}$$

To get $\frac{\partial \ell^{(i)}}{\partial w_{j,0}^{[0]}}$, backprop:

$$\begin{aligned} \frac{\partial \ell^{(i)}}{\partial w_{j,0}^{[0]}} &= \frac{\partial \ell^{(i)}}{\partial o^{(i)}} \frac{\partial o^{(i)}}{\partial z^{(i)}} \frac{\partial z^{(i)}}{\partial h_j^{(i)}} \frac{\partial h_j^{(i)}}{\partial z^{(i)}} \frac{\partial z^{(i)}}{\partial w_{j,0}^{[0]}} \\ \frac{\partial \ell^{(i)}}{\partial o^{(i)}} &= 2(o^{(i)} - y^{(i)}) \\ \frac{\partial o^{(i)}}{\partial z^{(i)}} &= \left\{ \text{numerically} \right\} \frac{e^{-z^{(i)}}}{(1 + e^{-z^{(i)}})^2} = \frac{1}{(1 + e^{-z^{(i)}})} \left(\frac{e^{-z^{(i)}}}{(1 + e^{-z^{(i)}})} \right) = \sigma^{(i)}(1 - \sigma^{(i)}) \\ \frac{\partial z^{(i)}}{\partial h_j^{(i)}} &= w_j^{[0]} \\ \frac{\partial h_j^{(i)}}{\partial z^{(i)}} &= \left\{ \text{numerically} \right\} \frac{\partial o^{(i)}}{\partial z^{(i)}} = h_j^{(i)}(1 - h_j^{(i)}) \\ \frac{\partial z^{(i)}}{\partial w_{j,0}^{[0]}} &= x_j^{(i)} \end{aligned} \quad \left. \begin{array}{l} \frac{\partial \ell^{(i)}}{\partial w_{j,0}^{[0]}} = 2(o^{(i)} - y^{(i)}) \sigma^{(i)}(1 - \sigma^{(i)}) w_j^{[0]} h_j^{(i)}(1 - h_j^{(i)}) x_j^{(i)} \end{array} \right\}$$

$$\Rightarrow w_{j,0}^{[0](t+1)} = w_{j,0}^{[0](t)} - \lambda \frac{\ell^{(i)}}{n w_{j,0}^{[0]}}$$

\Rightarrow Case where $j=1, j=2$:

$$\underbrace{w_{j,1}^{[1](t+1)} = w_{j,1}^{[1](t+1)} - \lambda \left[2(o^{(i)} - y^{(i)}) \sigma^{(i)}(1 - \sigma^{(i)}) w_1^{[2]} h_1^{(i)}(1 - h_1^{(i)}) x_1^{(i)} \right]}_{\ell^{(i)}}$$

```
1 import numpy as np
2
3 class MLP:
4     """
5         A simple Multi-Layer Perceptron (MLP) with one hidden layer.
6         Uses sigmoid activation in output layer (binary classification).
7     """
8     def __init__(self, input_size, hidden_size, output_size, learning_rate=0.01):
9         """
10            Initialize network parameters.
11            :param input_size: Number of input features.
12            :param hidden_size: Number of neurons in hidden layer.
13            :param output_size: Number of output neurons.
14            :param learning_rate: Step size for weight updates.
15        """
16        self.input_size = input_size
17        self.hidden_size = hidden_size
18        self.output_size = output_size
19        self.learning_rate = learning_rate
20
21        # Initialize weights and biases with small random values
22        self.W1 = np.random.randn(hidden_size, input_size) * 0.01 # Weights: Input → Hidden
23        self.b1 = np.zeros((hidden_size, 1)) # Bias: Hidden Layer
24        self.W2 = np.random.randn(output_size, hidden_size) * 0.01 # Weights: Hidden → Output
25        self.b2 = np.zeros((output_size, 1)) # Bias: Output Layer
26
27    def sigmoid(self, z):
28        """
29            Compute sigmoid activation function.
30        """
31        return 1 / (1 + np.exp(-z))
32
33    def sigmoid_derivative(self, a):
34        """
35            Compute derivative of sigmoid function.
36        """
37        return a * (1 - a)
```

```
34
35     def forward(self, X):
36         """
37             Perform forward propagation.
38             :param X: Input feature matrix.
39             :return: Output probabilities.
40         """
41         self.Z1 = np.dot(self.W1, X) + self.b1 # Compute hidden layer input
42         self.A1 = self.sigmoid(self.Z1) # Apply activation function
43
44         self.Z2 = np.dot(self.W2, self.A1) + self.b2 # Compute output layer input
45         self.A2 = self.sigmoid(self.Z2) # Output probabilities
46
47         return self.A2
48
49     def backward(self, X, y):
50         """
51             Perform backpropagation to compute gradients and update weights.
52             :param X: Input features.
53             :param y: True labels.
54         """
55         m = X.shape[1] # Number of samples
56
57         # Compute loss gradient at output layer
58         dZ2 = self.A2 - y # Derivative of binary cross-entropy loss
59         dW2 = (1 / m) * np.dot(dZ2, self.A1.T) # Gradient w.r.t. W2
60         db2 = (1 / m) * np.sum(dZ2, axis=1, keepdims=True) # Gradient w.r.t. b2
61
62         # Backpropagate to hidden layer
63         dA1 = np.dot(self.W2.T, dZ2) # Gradient w.r.t. hidden layer activations
64         dZ1 = dA1 * self.sigmoid_derivative(self.A1) # Gradient of activation function
65         dW1 = (1 / m) * np.dot(dZ1, X.T) # Gradient w.r.t. W1
66         db1 = (1 / m) * np.sum(dZ1, axis=1, keepdims=True) # Gradient w.r.t. b1
67
68         # Update weights using gradient descent
69         self.W1 -= self.learning_rate * dW1
70         self.b1 -= self.learning_rate * db1
71         self.W2 -= self.learning_rate * dW2
72         self.b2 -= self.learning_rate * db2
```

```
73
74     def train(self, X, y, epochs=1000):
75         """
76             Train the MLP using gradient descent.
77             :param X: Input features.
78             :param y: Target labels.
79             :param epochs: Number of iterations.
80         """
81         for epoch in range(epochs):
82             self.forward(X) # Forward propagation
83             self.backward(X, y) # Backpropagation
84
85             # Print loss every 100 iterations
86             if epoch % 100 == 0:
87                 loss = -np.mean(y * np.log(self.A2) + (1 - y) * np.log(1 - self.A2))
88                 print(f"Epoch {epoch}: Loss = {loss:.4f}")
89
90     def predict(self, X):
91         """
92             Make predictions.
93             :param X: Feature matrix.
94             :return: Predicted labels.
95         """
96         A2 = self.forward(X)
97         return (A2 > 0.5).astype(int) # Convert probabilities to binary predictions
98
%K for Command, %L for Cascade
```

2. Scikit-Learn Implementation of Multi-Layer Perceptron (MLP)

python

Copy

Edit

```
from sklearn.neural_network import MLPClassifier
from sklearn.datasets import make_classification
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler

# Generate synthetic classification dataset
X, y = make_classification(n_samples=1000, n_features=10, random_state=42)

# Split dataset into training and test sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_st

# Standardize features (important for neural networks)
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)

# Initialize MLP Classifier
mlp = MLPClassifier(hidden_layer_sizes=(10,), activation='relu', solver='adam',
                     learning_rate_init=0.01, max_iter=1000, random_state=42)

# Train the model
mlp.fit(X_train, y_train)

# Evaluate the model
accuracy = mlp.score(X_test, y_test)
print(f"MLP Accuracy: {accuracy:.4f}")
```

Final Takeaways

- ✓ Neural Networks can approximate any function, but generalization requires regularization.
 - ✓ Optimization and activation functions determine training efficiency.
 - ✓ Proper initialization, batch normalization, and tuning improve model stability.
 - ✓ Interpretability remains a challenge, but tools like SHAP and saliency maps help.
-