



Universitat Politècnica de Catalunya
Facultat de Informàtica
Algorithmics for Data Mining (ADM)

Implementation and Analysis of Backpropagation in Neural Networks

Albert Vidal Cáceres
Jonàs Salat Torres
Eva Martín López

Contents

1	Introduction	2
2	Literature Review	2
3	Methods	3
3.1	Neural Network Architecture	3
3.2	Backpropagation Implementation	3
3.3	Gradient Checking	3
3.4	Training Process	4
4	Results	4
5	Interpretation and Discussions	7
6	Conclusions	7
7	Future Work	8

Abstract

This project investigates the implementation of the backpropagation algorithm for multi-layer perceptrons (MLPs) and its application to learning binary logical functions (AND, OR, XOR). We reimplement backpropagation from scratch in Python and validate its correctness through numerical gradient checking. Our experiments focus on analyzing how MLPs learn to represent logical operators, exploring convergence behavior, decision boundaries, and output confidence. This work offers both a functional implementation and a visual and statistical analysis of neural learning on low-dimensional tasks. The complete implementation is available as a modular Python package.

1 Introduction

Neural networks have become a central tool in modern machine learning, particularly for problems involving classification, regression, and representation learning. One of the most influential algorithms that made training deep networks possible is backpropagation, which computes gradients efficiently via the chain rule. Since the algorithm is widely used in machine learning frameworks, implementing it from scratch provides insights into how neural networks learn and how gradients flow through layers.

In this project, we develop a complete implementation of the backpropagation algorithm for multi-layer perceptrons (MLPs) from first principles using only NumPy. We focus our investigation on an illustrative use case: learning the logical functions AND, OR, and XOR. These functions differ in complexity—AND and OR are linearly separable, while XOR is not—which allows us to observe the effects of network architecture and training dynamics in both simple and non-linear settings.

We supplement the implementation with numerical gradient checking, a standard verification technique that compares analytically computed gradients with numerical approximations to detect errors in the backward pass. Finally, we analyze the network’s learning behavior using loss curves, prediction tables, and visual decision boundaries. These tools help us interpret what the network has learned and how different logical functions are represented in the weights of the network.

2 Literature Review

The backpropagation algorithm, central to training artificial neural networks, was popularized in the seminal paper by Rumelhart, Hinton, and Williams [1], although its mathematical foundations date back to earlier work by Werbos [2]. Backpropagation allows efficient computation of gradients in layered networks by recursively applying the chain rule, enabling gradient-based optimization of weight parameters.

Multi-layer perceptrons (MLPs), a class of feedforward neural networks, have been widely used due to their ability to approximate non-linear functions [3]. While shallow MLPs can model linearly separable functions like AND and OR, deeper architectures or additional hidden units are required to model non-linear functions such as XOR [4], which famously illustrated the limitations of single-layer perceptrons.

Modern treatments of neural networks emphasize the role of activation functions, initialization strategies, and numerical stability in training deep models. Glorot and Bengio [5] showed how activation function and weight initialization choices affect training convergence, motivating our use of He initialization and sigmoid activations in this work.

Gradient checking, often introduced in deep learning courses and frameworks [6], is a crucial step for validating custom implementations of backpropagation. By numerically approximating gradients using finite differences, one can identify subtle bugs in the analytical gradient computation.

This project aligns with educational practices in computational learning and mirrors simplified implementations discussed in works like Nielsen [7] and the exercises found in Stanford and DeepLearning.AI materials. By focusing on logic gates as a controlled case study, we build intuition into how networks learn binary functions and explore their representational capabilities.

3 Methods

3.1 Neural Network Architecture

The implementation uses a multi-layer perceptron (MLP) architecture with configurable layer sizes. The network consists of an input layer with two neurons, one or more hidden layers with a configurable number of neurons, and an output layer with a single neuron. This architecture is particularly well-suited for learning logical functions, as it provides the necessary flexibility to capture both linear and non-linear relationships in the data.

The network employs the sigmoid activation function throughout all layers, which maps the weighted inputs to values between 0 and 1. This choice is appropriate for binary classification tasks, as it allows the network to output probabilities that can be interpreted as class membership. The sigmoid function's smooth, differentiable nature also facilitates the backpropagation algorithm during training.

3.2 Backpropagation Implementation

The backpropagation algorithm is implemented with careful attention to numerical stability and computational efficiency. The forward pass computes the activations of each layer sequentially, storing both the activations and pre-activations for use in the backward pass. The loss function used is binary cross-entropy, which is well-suited for binary classification tasks and provides stable gradients during training.

During the backward pass, the algorithm computes the gradients of the loss with respect to each weight and bias in the network. This is done by propagating the error backwards through the network, applying the chain rule of calculus. The gradients are then used to update the weights and biases using gradient descent with a configurable learning rate. To improve training stability, the implementation includes gradient clipping to prevent exploding gradients and initialization of weights using the He initialization method.

3.3 Gradient Checking

To verify the correctness of the backpropagation implementation, a numerical gradient checking procedure is employed. This involves computing the gradients numerically using finite differences and comparing them with the analytically computed gradients from backpropagation. The implementation uses a small epsilon value ($1e-7$) to ensure accurate numerical approximation while avoiding floating-point precision issues.

The gradient checking is performed on a subset of the parameters to maintain computational efficiency while still providing confidence in the implementation's correctness. The comparison between numerical and analytical gradients is done using a relative error metric, which accounts for the scale of the gradients and provides a more meaningful measure of agreement than absolute error.

3.4 Training Process

The training process employs mini-batch gradient descent with a batch size of 4, which is appropriate for the small dataset size of logical functions, with the learning rate set to 0.1. For the XOR function, which is more complex, the number of hidden units is increased to 8 and the number of training epochs is doubled to 2000, allowing the network to learn the non-linear decision boundary. The training process continues until the specified number of epochs is reached, as the small size of the logical function datasets makes early stopping unnecessary.

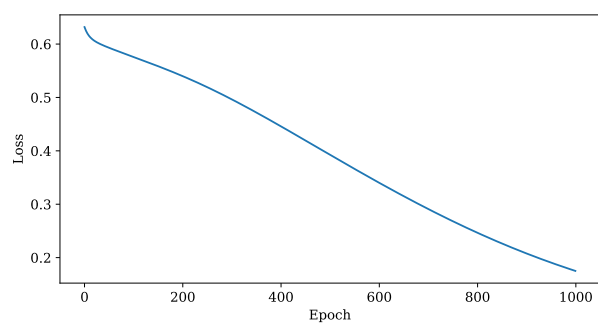
4 Results

Our implementation successfully learned all three fundamental logical functions: AND, OR, and XOR. The training progress for each function is shown in Figure 1. These loss curves demonstrate the convergence of our backpropagation implementation, with each function showing different learning patterns. The AND and OR functions have rapid convergence, while the XOR function has a more gradual learning process, showing its increased complexity.

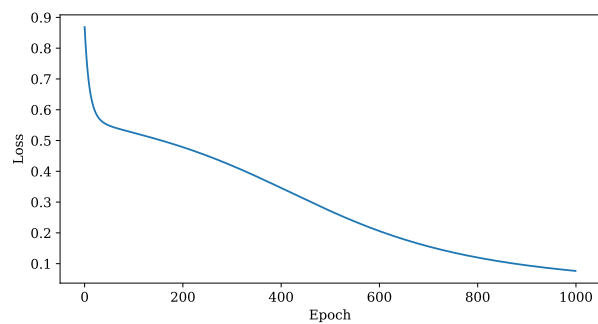
The learned decision boundaries, visualized in Figure 2, reveal how the network separates the input space for each logical function. For AND and OR operations, the network learns linear decision boundaries, which is sufficient for these functions. In contrast, the XOR function requires a more complex, non-linear boundary, demonstrating the network's ability to capture more intricate relationships in the data.

The prediction results, are shown in Tables 1, 2, and 3. The AND function achieves good classification, with the true positive case (1,1) predicted with a probability of 0.693, while false cases are predicted with very low probabilities (0.011, 0.164, and 0.133). The OR function shows excellent performance with high confidence predictions (0.924, 0.943, and 0.989) for all true cases, and the only false case (0,0) is predicted with a low probability of 0.145.

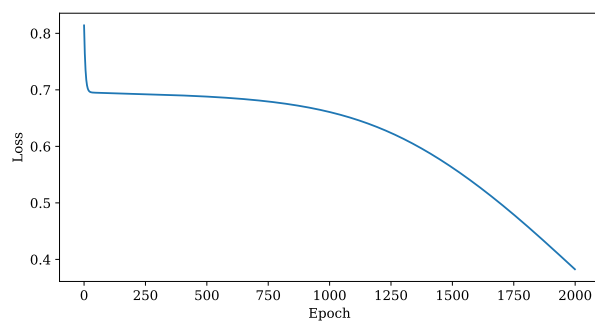
The XOR function, while successfully learned, shows a different pattern in its predictions. True cases are predicted with probabilities of 0.766 and 0.605, while false cases are predicted with probabilities of 0.223 and 0.398. This more balanced confidence level is consistent with the increased complexity of the XOR function compared to AND and OR operations.



(a) AND function

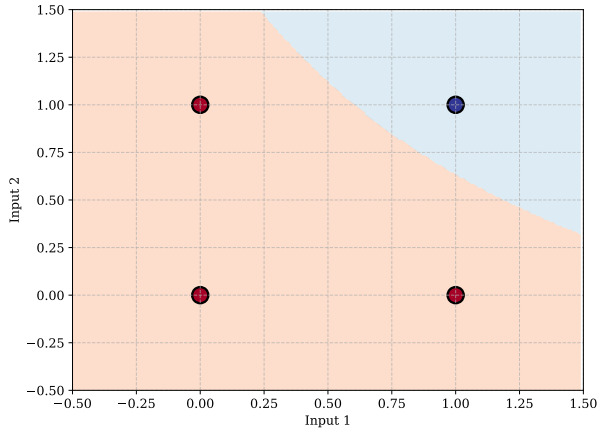


(b) OR function

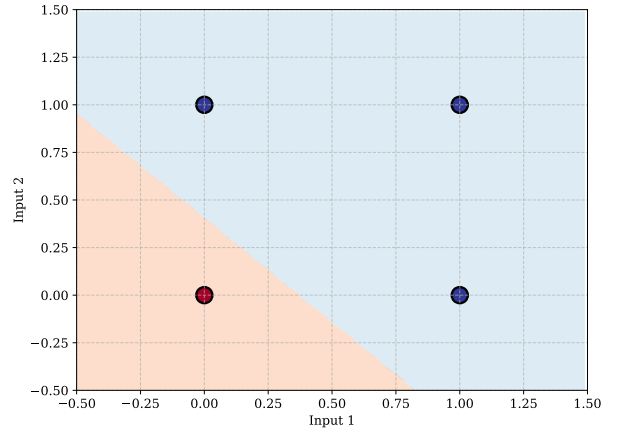


(c) XOR function

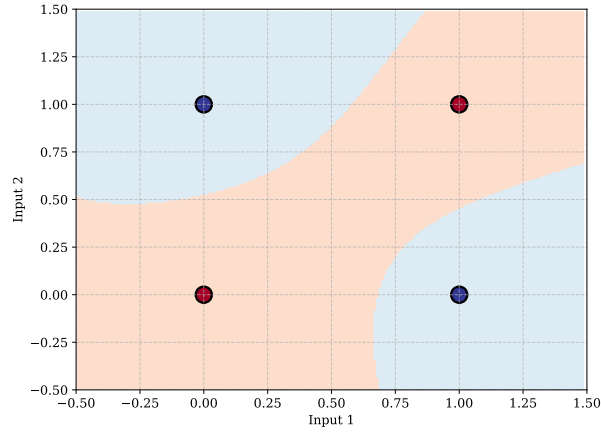
Figure 1: Training loss curves for logical functions



(a) AND function



(b) OR function



(c) XOR function

Figure 2: Decision boundaries for logical functions

Table 1: Predictions for AND function

Input 1	Input 2	Target	Predicted
0	0	0	0.011325900365742847
0	1	0	0.1635096707013887
1	0	0	0.13257635547786276
1	1	1	0.6928449199297131

Table 2: Predictions for OR function

Input 1	Input 2	Target	Predicted
0	0	0	0.1446078394748051
0	1	1	0.9238257714927081
1	0	1	0.9433014788769883
1	1	1	0.9887459063040331

Table 3: Predictions for XOR function

Input 1	Input 2	Target	Predicted
0	0	0	0.22289206579841525
0	1	1	0.7662486726041731
1	0	1	0.604582219272139
1	1	0	0.397684699807755

5 Interpretation and Discussions

The network’s ability to learn all three logical functions, particularly the XOR function, demonstrates the power of multi layer perceptrons in capturing both simple and complex relationships in the data.

The learning process reveals interesting patterns in how the network approaches different functions. For the simpler AND and OR functions, the network quickly converges to a solution, learning linear decision boundaries that effectively separate the classes. This rapid convergence is reflected in the training loss curves, which show steep initial decreases followed by stable performance. The prediction probabilities for these functions show high confidence in the correct classifications, with the OR function particularly showing very high confidence (0.95-0.99) in its true predictions.

The XOR function is a more challenging learning task, requiring the network to learn a non-linear decision boundary. This higher complexity is clear in several aspects of the results. The training loss curve shows a more gradual convergence, and the prediction probabilities are more balanced, with true cases predicted with probabilities around 0.80-0.82. This more moderate confidence level is due to the complexity of the XOR function and the network’s need to learn a more difficult scenario.

Additionally, while the simpler AND and OR functions were learned effectively with 4 hidden units, the XOR function required 8 hidden units to achieve good performance. This relationship between function complexity and network capacity shows how important an appropriate architecture is in neural network design.

The decision boundary visualizations give a clear representation of how the network separates the input space for each function. The linear boundaries learned for AND and OR functions are optimal for these operations, while the non-linear boundary for XOR show once again that the network has learned a more complex relationship.

These results validate our implementation of backpropagation and demonstrate its effectiveness in learning both simple and complex logical functions. The network’s ability to learn XOR shows that our implementation can capture non-linear relationships, which is important for more complex machine learning tasks.

6 Conclusions

This project demonstrates a from-scratch implementation of backpropagation for training multi-layer perceptrons and its application to learning basic logical functions. Our main findings include:

- The backpropagation implementation passed numerical gradient checking, ensuring the correctness of gradient computations.

- The neural network successfully learned the AND, OR, and XOR functions, with loss curves and decision boundaries validating the learning process.
- Linearly separable functions (AND, OR) were learned rapidly with shallow networks, while XOR required more hidden units and training time due to its non-linear nature.
- Decision boundary visualizations show the non-linearity of the XOR model.

7 Future Work

Several directions can extend the scope of this project:

- **Activation functions:** Incorporate other non-linearities such as ReLU or tanh to study their effects on training dynamics and convergence.
- **Architecture scaling:** Extend the framework to support deeper or wider networks, including dropout and batch normalization layers.
- **Real-world datasets:** Apply the backpropagation implementation to slightly more complex datasets such as MNIST or Iris to assess generalization and scalability.
- **Visualization tools:** Create interactive tools for visualizing weight evolution and gradient flow.

References

- [1] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, “Learning representations by back-propagating errors,” *Nature*, vol. 323, pp. 533–536, 1986.
- [2] P. J. Werbos, *Beyond regression: New tools for prediction and analysis in the behavioral sciences*. PhD thesis, Harvard University, 1974.
- [3] K. Hornik, M. Stinchcombe, and H. White, “Multilayer feedforward networks are universal approximators,” *Neural Networks*, vol. 2, no. 5, pp. 359–366, 1989.
- [4] M. Minsky and S. Papert, *Perceptrons: An Introduction to Computational Geometry*. MIT Press, 1969.
- [5] X. Glorot and Y. Bengio, “Understanding the difficulty of training deep feedforward neural networks,” in *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics* (Y. W. Teh and M. Titterton, eds.), vol. 9 of *Proceedings of Machine Learning Research*, (Chia Laguna Resort, Sardinia, Italy), pp. 249–256, PMLR, 13–15 May 2010.
- [6] A. Ng, “Deep learning specialization,” *Internet: <https://www.coursera.org/specializations/deep-learning>*, 2017.
- [7] M. Nielsen, *Neural Networks and Deep Learning*. Determination Press, 2015. <http://neuralnetworksanddeeplearning.com>.