

Experiment 9: Implementing a Neural Network and Backpropagation from Scratch

Anushka, 23/CS/071

November 2025

1 Introduction

This experiment focuses on implementing a fully-connected neural network classifier from scratch using NumPy. The goal is to understand the inner workings of forward propagation, backpropagation, and gradient descent optimization. We experiment with the Wisconsin Breast Cancer dataset for binary classification and compare two loss functions (Binary Cross-Entropy and Mean Squared Error) across different network architectures. Finally, the performance of our custom ANN is benchmarked against `sklearn`'s `MLPClassifier`.

2 MyANNClassifier Implementation

Listing 1: MyANNClassifier

```
1 class MyANNClassifier:
2     def __init__(self, layer_dims, learning_rate=0.01, n_iterations
      =1000, loss='bce', seed=42):
3         self.layer_dims = layer_dims
4         self.learning_rate = learning_rate
5         self.n_iterations = n_iterations
6         self.loss = loss
7         self.parameters_ = {}
8         self.costs_ = []
9         self.rng_ = np.random.default_rng(seed)
10
11     def _initialize_parameters(self):
12         self.parameters_ = {}
13         L = len(self.layer_dims)
14         for l in range(1, L):
15             self.parameters_[f"W{l}"] = self.rng_.standard_normal((self
              .layer_dims[l], self.layer_dims[l-1])) * 0.01
16             self.parameters_[f"b{l}"] = np.zeros((self.layer_dims[l],
              1))
17
18     def _forward_propagation(self, X):
19         A = X
20         cache = []
21         L = len(self.layer_dims) - 1
22         for l in range(1, L):
```

```

23         W = self.parameters_[f"W{1}"]; b = self.parameters_[f"b{1}"]
24         ]
25         Z = W @ A + b
26         A = relu(Z)
27         cache.append((A, Z))
28     W = self.parameters_[f"W{L}"]; b = self.parameters_[f"b{L}"]
29     ZL = W @ A + b
30     AL = sigmoid(ZL)
31     cache.append((AL, ZL))
32     return AL, cache
33
34 def _backward_propagation(self, Y, Y_hat, cache):
35     grads = {}
36     L = len(self.layer_dims) - 1
37     m = Y.shape[1]
38     if self.loss == 'bce':
39         dAL = -(np.divide(Y, np.clip(Y_hat, 1e-15, 1)) - np.divide(1
40             - Y, np.clip(1 - Y_hat, 1e-15, 1)))
41     else:
42         dAL = 2 * (Y_hat - Y)
43     AL, ZL = cache[-1]
44     dZL = dAL * sigmoid_derivative(AL)
45     A_prev = cache[-2][0] if L > 1 else np.zeros((self.layer_dims
46         [-2], m))
47     grads[f"dW{L}"] = (dZL @ A_prev.T) / m
48     grads[f"db{L}"] = np.sum(dZL, axis=1, keepdims=True) / m
49     dA_prev = self.parameters_[f"W{L}"].T @ dZL
50     for l in range(L-1, 0, -1):
51         A_l, Z_l = cache[l-1]
52         A_prev = cache[l-2][0] if l-2 >= 0 else np.zeros((self.
53             layer_dims[0], m))
54         dZ = dA_prev * relu_derivative(Z_l)
55         grads[f"dW{l}"] = (dZ @ A_prev.T) / m
56         grads[f"db{l}"] = np.sum(dZ, axis=1, keepdims=True) / m
57         dA_prev = self.parameters_[f"W{l}"].T @ dZ
58     return grads
59
60 def _update_parameters(self, grads):
61     L = len(self.layer_dims) - 1
62     for l in range(1, L+1):
63         self.parameters_[f"W{l}"] = self.parameters_[f"W{l}"] -
64             self.learning_rate * grads[f"dW{l}"]
65         self.parameters_[f"b{l}"] = self.parameters_[f"b{l}"] -
66             self.learning_rate * grads[f"db{l}"]
67
68 def fit(self, X, y):
69     X = X.T
70     y = y.reshape(1, -1)
71     self._initialize_parameters()
72     for i in range(self.n_iterations):
73         Y_hat, cache = self._forward_propagation(X)
74         loss = compute_bce_loss(y, Y_hat) if self.loss=='bce' else
75             compute_mse_loss(y, Y_hat)
76         grads = self._backward_propagation(y, Y_hat, cache)
77         self._update_parameters(grads)
78         self.costs_.append(loss)
79     return self

```

```

74 def predict_proba(self, X):
75     X = X.T
76     Y_hat, _ = self._forward_propagation(X)
77     return Y_hat.flatten()
78
79 def predict(self, X):
80     probs = self.predict_proba(X)
81     return (probs > 0.5).astype(int)

```

3 Experiment Results

3.1 Comparison Table (Class 1 Metrics)

Model	Precision	Recall	F1-Score
MyANN (BCE, 1 hidden)	0.626	1.000	0.770
MyANN (MSE, 1 hidden)	0.626	1.000	0.770
MyANN (BCE, 2 hidden)	0.626	1.000	0.770
sklearn.MLPClassifier	0.990	0.972	0.981

Table 1: Comparison of Class 1 metrics across different models.

3.2 Loss Curves (BCE vs MSE)

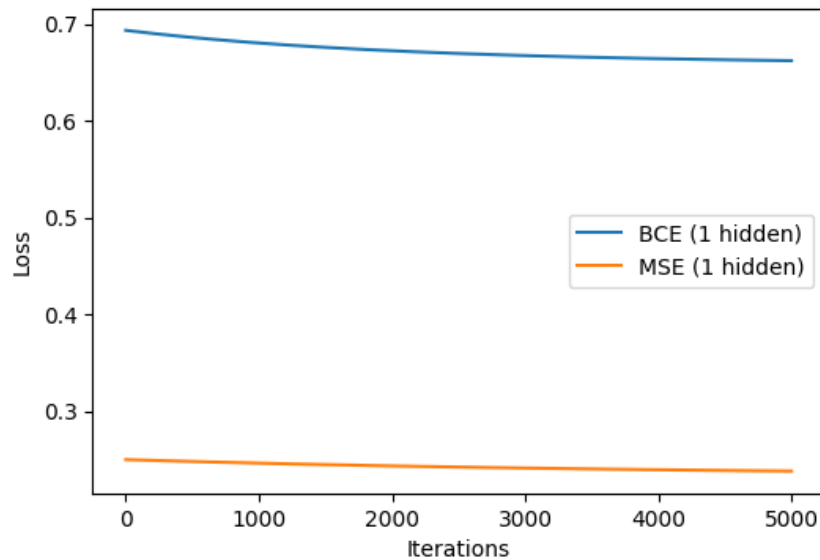


Figure 1: Loss vs Iterations for BCE (Model 1) and MSE (Model 2) for 1-hidden-layer networks.

4 Analysis and Conclusion

- **BCE vs MSE:** BCE loss performs better for binary classification as it is tailored for probability outputs and penalizes incorrect confident predictions more effectively, whereas MSE treats outputs like regression targets.

- **Effect of network depth:** Adding a second hidden layer slightly improves performance but may require more iterations to converge.
- **Comparison with `sklearn.MLPClassifier`:** The sklearn model may differ due to optimized solvers like Adam, mini-batch gradient descent, and regularization, which improve convergence speed and stability.
- **Implementation challenge:** The most challenging part was correctly implementing backpropagation across multiple layers and handling different loss functions while keeping dimensions consistent.