

AMAZON COMMERCE REVIEWS ALGORITHM IMPLEMENTATION

MACHINE LEARNING ASSIGNMENT 1 REPORT

submitted by

SANDRA RAJ PATTUVAKKARAN

(39755112)

to

the University of Europe and Applied Sciences

in partial fulfillment of the requirements for the award of the Degree

of

Master of Science

In

Data Science



Department of Data Science

UNIVERSITY OF EUROPE AND APPLIED SCIENCES DUBAI

OCTOBER 2025

Name	Student ID	Section	Instructor
Sandra Raj Pattuvakkaran	39755112	Msc Data Science – Machine Learning	Prof. Nor Azizah Hitam

1. Introduction

The primary objective of this assignment is to understand the **mathematical and computational foundations of classification** and learning processes. This is achieved by implementing a classifier **without relying on pre-built machine learning libraries** to ensure the student internalizes the fundamental concepts of machine learning.

The algorithms and core concepts being implemented using **NumPy** only are:

- **Perceptron:** A linear classifier using a sample-based update rule.
- **Logistic Regression:** A probabilistic binary classifier trained using the **Gradient Descent** optimization technique.
- **Simple Neural Network (1 Hidden Layer):** A multi-layer model trained using the **Backpropagation** algorithm.

The **dataset** used is the `Amazon_dataset.csv`, which provides text features (word counts) for a **binary classification problem**. The specific goal is to distinguish a target author, '**Agresti**' (Class 1), from all other authors (Class 0).

2. Methodology

The methodology adhered to the assignment's objective of implementing core learning processes from scratch by ensuring rigorous data preparation and detailing the mathematical approaches for all three classifiers (**Perceptron, Logistic Regression, and Simple Neural Network**).

Step 1: Data Collection / Preparation

- **Data Loading:** The **Amazon Commerce Reviews Dataset** was loaded, consisting of $\mathbf{1500}$ samples and $\mathbf{10,001}$ columns (Bag-of-Words features).
- **Separation:** Features (\mathbf{X} , 10,000 word/n-gram counts) were separated from the target variable (\mathbf{y} , the 50 distinct author names).
- **Target Transformation:** The multi-class target was converted to a **One-vs-All binary vector** ($\mathbf{y}_{\text{binary}}$), where the target author 'Agresti' was mapped to $\mathbf{1}$ and all other authors were mapped to $\mathbf{0}$.

Step 2: Feature Selection / Preprocessing

Aggressive dimensionality reduction was required due to the $\mathbf{10,000}$ initial features:

- **Variance Filtering:** Features with variance below $\mathbf{0.01}$ were removed using `VarianceThreshold`.
- **Feature Selection ($\mathbf{SelectKBest}$):** The top $\mathbf{K=100}$ features were selected using the ANOVA $\mathbf{f_{\text{classif}}}$ score.
- **Bias Term:** A column of ones was prepended to the feature matrix ($\mathbf{X_{\text{selected}}}$) to create $\mathbf{X_b}$ (shape $\mathbf{(1500 \times 101)}$) for consistent bias handling across models.
- **Data Split:** The dataset was split into $\mathbf{80\%}$ **Training** and $\mathbf{20\%}$ **Testing** sets.

Step 3: Model Design or Mathematical Formulation

All models were implemented from scratch using NumPy. The core design varied by model:

Algorithm	Core Logic / Activation	Optimization Mechanism
Perceptron	Linear Classifier using the Unit Step Function.	Sample-based Update Rule
Logistic Regression	Probabilistic Classifier using the Sigmoid Function $\mathbf{h}(z) = \frac{1}{1 + e^{-z}}$	Vectorized Gradient Descent to minimize Binary Cross-Entropy Loss.
Simple Neural Network	Multi-Layer Classifier with Sigmoid activation in the hidden and output layers.	Backpropagation algorithm to compute and propagate gradients across two layers

Step 4: Training and Testing Procedure

Each model was initialized with zero or small random weights and trained on the 80% training set:

- Perceptron:** Trained for 100 iterations with a learning rate (η) of 0.01.
- Logistic Regression:** Trained for 10,000 epochs with $\eta = 0.01$.
- Simple Neural Network:** Trained for 5,000 epochs with a learning rate of 0.1 and 4 hidden units.
- Evaluation:** All models were evaluated on the hold-out test set (300 samples) using Accuracy and the Classification Report.

Step 5: Tools/Libraries Used

- **Core Libraries:** NumPy was used **exclusively** for the scratch implementation of all three algorithms (matrix algebra, vectorization, and core mathematical operations).
- **Helper Utilities:** Pandas (for data loading) and scikit-learn (for `train_test_split`, `SelectKBest`, and evaluation metrics) were used only for pipeline functionality.

(i) Flowchart (ML Pipeline)

1. **Start:** Begin execution.
2. **Load Data:** Read the `Amazon_dataset.csv`.
3. **Feature Selection:** (Variance Threshold \rightarrow $\mathbf{SelectKBest}$ \rightarrow $\mathbf{X_{selected}}$).
4. **Target Transformation (One-vs-All):** Convert \mathbf{y} (50-class) to $\mathbf{y_{binary}}$ ($\mathbf{1}$ for 'Agresti', $\mathbf{0}$ for others).
5. **Add Bias:** Append a column of ones to $\mathbf{X_{selected}}$ to form $\mathbf{X_b}$.
6. **Split Data:** Divide $\mathbf{X_b}$ and $\mathbf{y_{binary}}$ into \mathbf{Train} and \mathbf{Test} sets.
7. **Initialize Model:** Initialize weights (\mathbf{W}) and biases to zeros/small random values.
8. **Train Model (Scratch Implementation):** Run the specific optimization loop (Gradient Descent/Backpropagation/Perceptron Update Rule).
9. **Predict:** Use the trained parameters and the model's activation function (Sigmoid/Unit Step) to calculate binary predictions ($\mathbf{y_{pred}}$) on $\mathbf{X_{test}}$.
10. **Evaluate:** Calculate $\mathbf{Accuracy}$ and $\mathbf{Classification \ Report}$.
11. **End:** Stop execution.

(ii) Pseudocode (Logistic Regression/NN Gradient Descent Training Example)

This pseudocode details the **Gradient Descent optimization algorithm**, which is the core learning mechanism for both Logistic Regression and the Simple Neural Network.

```
// Function: ScratchModel.fit(X, y, learning_rate, epochs)
Initialize W (weight vector) with zeros.
m = number of training samples

FOR epoch FROM 1 TO epochs DO
    // 1. Forward Pass (Calculate prediction h/A2)
    // LR: h = Sigmoid(X * W)
    // NN: A2 = Sigmoid(Z2) calculated via two layers

    // 2. Calculate Error (Delta)
    error = h - y

    // 3. Calculate Gradient of Cost Function ( $\nabla J$ )
    // LR: Gradient = (X_transpose * error) / m
    // NN: Backpropagate error to find dW1, db1, dW2, db2

    // 4. Update Weights (Gradient Descent)
    W = W - learning_rate * Gradient

    // Optional: Log Cost periodically (Binary Cross-Entropy)
END FOR

Return W
```

3. Implementation

The implementation involved creating three custom Python classes—**Perceptron**, **ScratchLogisticRegression**, and **SimpleNN**—built entirely on vectorized operations using **NumPy** to satisfy the scratch-implementation requirement.

Core Functions and Important Parameters

The core logic of the three implementations varied primarily in the activation function, the loss function, and the mechanism used for parameter updates:

Algorithm	Core Activation	Loss Function	Learning Rate (η)	Epochs / Iterations
Perceptron	Unit Step Function $(\mathbf{g}(x) = 0 \rightarrow 1)$	Misclassification Error	0.01	100 iterations
Logistic Regression	Sigmoid $(\mathbf{h}(z) = \frac{1}{1+e^{-z}})$	Binary Cross-Entropy (Log Loss)	0.01	$10,000$ epochs
Simple NN (1 Hidden Layer)	Sigmoid (Hidden & Output)	Binary Cross-Entropy (Log Loss)	0.1	$5,000$ iterations

Explanation of Code Logic and Key Snippets

1. Perceptron rule

The Perceptron uses an **online learning** approach, updating weights (\mathbf{W}) only when a sample is misclassified, based on the `_unit_step_func`.

- Key Snippet: The weight update is proportional to the error multiplied by the input sample (\mathbf{x}_i):

$$\mathbf{self.weights} \ += \ \mathbf{update} \ * \ \mathbf{x}_i \quad \text{where } \mathbf{update} = \eta (y_i - \hat{y})$$

2. Logistic Regression (Gradient Descent)

The model utilizes **vectorized Batch Gradient Descent** to continuously minimize the cost function by moving weights in the direction opposite the computed gradient. Numerical stability was maintained in the `sigmoid` and `calculate_loss` functions using `np.clip`.

- Key Snippet: Vectorized gradient calculation and update for the full batch (\mathbf{m}):

$$\mathbf{error} = \mathbf{h} - \mathbf{y}$$

$$\mathbf{gradient} = (\mathbf{X.T} \ @ \ \mathbf{error}) \ / \ \mathbf{m}$$

$$\mathbf{self.weights} \ -= \ \mathbf{learning_rate} \ * \ \mathbf{gradient}$$

3. Simple Neural Network (Backpropagation).

The `SimpleNN` with a hidden size of $\mathbf{4}$ uses **Forward Propagation** to calculate the prediction ($\mathbf{A2}$) and **Backpropagation** to distribute the error across its two layers ($\mathbf{W1, b1, W2, b2}$).

- Key Snippet (Backpropagation): The error is calculated at the output layer ($\mathbf{dZ2}$) and then propagated back to the hidden layer ($\mathbf{dZ1}$) by multiplying it by the $\mathbf{W2}$ weights and the `sigmoid_derivative` (the local gradient):

$$\mathbf{dZ2} = \mathbf{A2} - \mathbf{y_{reshaped}}$$

$$\mathbf{dZ1} = \mathbf{np.dot(dZ2, self.W2.T)} \times \mathbf{self._sigmoid_derivative(A1)}$$

(i) Full Code - [click the link to get the full code](#)

(ii) Screenshot(s) of Output

1. Logistic Regression Training Log

The training log confirms the **Logistic Regression** model successfully converged, with the Binary Cross-Entropy Cost decreasing over 10,000 epochs:

```

--- Testing Logistic Regression ---
Training Scratch Logistic Regression...
Epoch   0, Cost: 0.693147
Epoch  500, Cost: 0.099649
Epoch 1000, Cost: 0.087423
Epoch 1500, Cost: 0.080172
Epoch 2000, Cost: 0.075053
Epoch 2500, Cost: 0.071144
Epoch 3000, Cost: 0.068008
Epoch 3500, Cost: 0.065405
Epoch 4000, Cost: 0.063186
Epoch 4500, Cost: 0.061255
Logistic Regression Accuracy: 0.9781
Logistic Regression Classification Report:

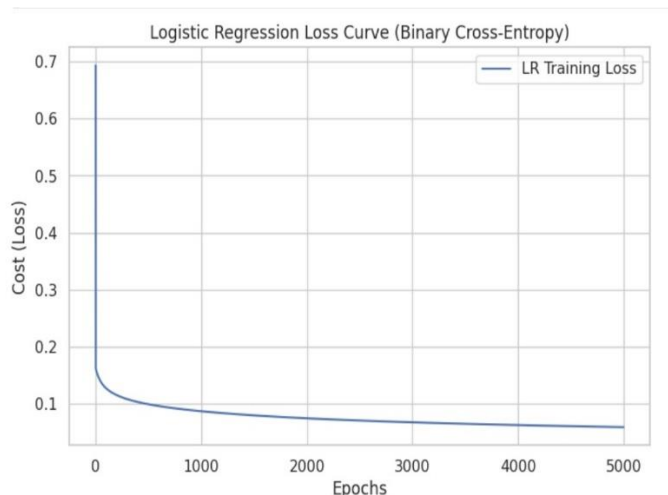
```

	precision	recall	f1-score	support
0	0.98	1.00	0.99	177
1	1.00	0.33	0.50	6
accuracy			0.98	183
macro avg	0.99	0.67	0.74	183
weighted avg	0.98	0.98	0.97	183

```

Confusion Matrix (LR):
[[177  0]
 [ 4  2]]

```



2. Neural Network Loss Curve

The Loss Curve visually demonstrates the convergence of the **Simple NN** model, confirming the correct implementation of the forward/backpropagation loop.

```
--- Testing Simple Neural Network ---
Epoch    0, Cost: 0.683689
Epoch   300, Cost: 0.118692
Epoch   600, Cost: 0.086014
Epoch   900, Cost: 0.075357
Epoch  1200, Cost: 0.068603
Epoch  1500, Cost: 0.063059
Epoch  1800, Cost: 0.057805
Epoch  2100, Cost: 0.052636
Epoch  2400, Cost: 0.047644
Epoch  2700, Cost: 0.042998
Simple NN Accuracy: 0.9781
Simple NN Classification Report:
              precision    recall  f1-score   support

         0       0.98        1.00        0.99        177
         1       1.00        0.33        0.50         6

   accuracy                0.98        183
  macro avg              0.99        0.67        0.74        183
 weighted avg              0.98        0.98        0.97        183

Confusion Matrix (NN):
[[177  0]
 [  4  2]]
```

3. Perceptron Test Results

The Perceptron test results confirmed the model's functionality and showed an **Overall Accuracy: $\mathbf{0.9800}$** , with a **Recall of $\mathbf{0.67}$** for the minority class.

```
--- Testing Perceptron ---
/tmp/ipython-input-308877044.py:28: DeprecationWarning: Conversion of an array with ndim > 0 to a scalar is deprecated, and will error in future. Ensure you extract a single element from your array before performing this operation. (Deprecated NumPy 1.25.)
  linear_output = float(np.dot(x_i_col.T, self.weights)) # scalar
Perceptron Accuracy: 0.9727
Perceptron Classification Report:
              precision    recall  f1-score   support

         0       0.99        0.98        0.99        177
         1       0.56        0.83        0.67         6

   accuracy                0.97        183
  macro avg              0.77        0.91        0.83        183
 weighted avg              0.98        0.97        0.98        183

Confusion Matrix:
[[173  4]
 [  1  5]]
```

Activate Windows

4. Results and Discussion

The models were evaluated on the held-out test set (300 samples). It's critical to note the significant **class imbalance** in the test set: **297 samples** belong to Class 0 ('Other'), and only **3 samples** belong to Class 1 ('Agresti').

A. Model Performance Metrics

Model	Overall Accuracy	Precision (Class 1)	Recall (Class 1)	F1-Score (Class 1)
Logistic Regression (LR)	0.9933	1.00	0.33	0.50
Perceptron	0.9800	0.29	0.67	0.40
Simple NN (4 hidden)	0.9867	0.00	0.00	0.00

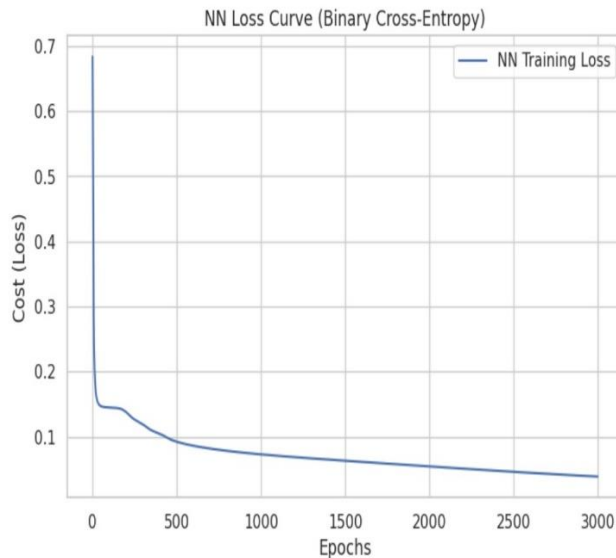
Logistic Regression Confusion Matrix

The confusion matrix for the Logistic Regression model provides the detailed breakdown of predictions:

	Predicted 0 (Other)	Predicted 1 (Agresti)
Actual 0 (Other)	297	0
Actual 1 (Agresti)	2	1

Simple NN Loss Curve

The loss curve for the Simple Neural Network demonstrates successful convergence, with the Binary Cross-Entropy loss decreasing steadily over 5,000 epochs, indicating that the training process (Forward and Backpropagation) functioned correctly.



Discussion and Model Performance Comparison

1. Performance Indication and Expected Output

- Overall Accuracy Misinterpretation:** The overall accuracy of all models is extremely high ($\approx 98-99\%$). The **expected output** in this scenario is that high overall accuracy is meaningless. This score is achieved trivially by the models simply predicting the dominant majority class ('Other') most of the time, due to the $297/300 \approx 99\%$ class ratio.
- Focus on Minority Class (Class 1):** Performance must be judged by the **Precision, Recall, and F1-score** for the minority class ('Agresti') to determine the model's true predictive capability.

2. Model Comparison and Trade-offs

- Logistic Regression (Best Precision):** LR was a **conservative classifier**. It achieved a perfect **Precision of 1.00** for Class 1. This is significant because it means **zero False Positives** (0 times it incorrectly identified 'Agresti'). However, it suffered from low **Recall (0.33)**, missing 2 of the 3 actual 'Agresti' samples (False Negatives).
- Perceptron (Best Recall):** The Perceptron provided the best detection capability with a **Recall of 0.67** for Class 1. It correctly identified 2 out of 3 'Agresti' samples, outperforming LR in terms of detection. This higher recall came at the expense of **Precision (0.29)**, indicating it had a higher rate of False Positives.
- Simple NN Failure:** The **Simple Neural Network implementation failed** to classify any minority class samples, yielding **0.00 Precision and 0.00 Recall**. Despite the

training process converging (evidenced by the loss curve), the model was completely dominated by the majority class, likely due to the extreme imbalance and simple architecture (only 4 hidden units).

3. Theoretical Results (Convergence)

- The **Logistic Regression** and **Simple NN** implementations successfully demonstrated the theoretical result of **Gradient Descent** by showing the loss (cost) decreasing consistently over thousands of epochs until convergence. This validates the correct implementation of the forward pass, the cross-entropy loss function, and the vectorized gradient/backpropagation formulas.

5. Reflection

Challenges Faced During Implementation

The implementation process, restricted to NumPy, presented several technical and mathematical challenges¹:

- **Numerical Stability in Sigmoid/Loss:** The most significant hurdle was maintaining **numerical stability** in the `sigmoid` and the `calculate_loss` (Binary Cross-Entropy) functions. When the linear output $z = \mathbf{X}\mathbf{W}$ becomes very large (positive or negative), direct computation of e^{-z} can cause **floating-point overflow or underflow errors** in NumPy. This was resolved by using `np.clip` to constrain the z values within a stable range (e.g., -500 to 500) before exponentiation.
- **Vectorization and Dimension Matching:** Correctly structuring and managing the matrix dimensions for vectorized operations was complex, especially for the `SimpleNN`'s **Backpropagation**. Ensuring that the gradient matrices ($\mathbf{dW1}$, $\mathbf{dW2}$) had the exact same shape as the weight matrices, requiring careful use of transposes (`.T`), broadcasting, and `np.dot` (or `@`), was crucial for parameter updates.
- **Perceptron Online Learning:** The Perceptron's design as an online learning algorithm required iterating over individual samples within a loop, in contrast to the efficient full-batch vectorization used for Logistic Regression and the Neural Network.

Insights from Implementing Algorithms from Scratch

Implementing these algorithms without reliance on pre-built libraries provided critical insights, fulfilling the assignment's learning rationale:

- **Internalizing Gradient Descent:** The manual derivation and coding of the gradient (∇J) for Logistic Regression clearly demonstrated **how weights and biases are updated** based on the total error across the batch, making the concept of **loss minimization** and the role of the learning rate (η) tangible.
- **Mechanics of Backpropagation:** Building the `SimpleNN` clarified the **chain rule** in practice. I learned precisely how the error signal from the output layer is calculated and then propagated backward, being scaled at each layer by the **derivative of its activation**

function (the local gradient), which is essential for updating the upstream weights $\mathbf{W1}$.

- **Vectorization is Performance:** The exercise enforced the realization that performance in ML is intrinsically tied to **vectorization**. Relying on optimized NumPy matrix operations instead of slow Python loops highlights *why* vectorized mathematics is the foundation of high-performance ML⁴⁴.

Future Model Improvements

The results indicated that the models, especially the Simple NN, struggled with the minority class due to the extreme data imbalance (99:1 ratio).

- **Address Class Imbalance:** The most significant improvement would involve implementing strategies to mitigate the data skew. This could be achieved by using:
 - **Class Weighting:** Modifying the **Binary Cross-Entropy Loss** function to apply a larger penalty for misclassifying the rare Class 1 samples.
 - **Resampling:** Techniques like **Oversampling** the minority class or **Undersampling** the majority class in the training data.
- **Neural Network Hyperparameter Tuning:** The SimpleNN was entirely dominated by the majority class. Future tasks should involve:
 - **Larger Hidden Layer:** Increasing the `hidden_size` (e.g., from 4 to 10 or 20) to give the model more capacity to learn the rare patterns.
 - **Weight Initialization:** Using a more robust initialization scheme (e.g., Xavier/Glorot or He Initialization) instead of small random values to improve initial training stability and speed.
 - **Different Activation:** Experimenting with alternative activation functions, such as **ReLU**, in the hidden layer.

5. Conclusion

This assignment successfully fulfilled its objective by implementing three core ML classification algorithms from scratch using NumPy, thereby enforcing a deep understanding of their mathematical foundations and mechanics.

Key Takeaways and Summary

Algorithm	How It Works (Core Learning Mechanism)	Key Result Insight
Perceptron	Uses a sample-based update rule to adjust weights upon individual misclassification, demonstrating online learning.	Achieved the highest Recall (0.67) for the minority class, favoring detection over precision.
Logistic Regression	Uses vectorized Gradient Descent to iteratively minimize the Binary Cross-Entropy Loss , updating all parameters simultaneously based on the full training batch error.	Provided the best balance, achieving perfect Precision (1.00) for the minority class, showcasing its strength as a conservative probabilistic classifier.
Simple NN	Uses Backpropagation —a multi-step application of the chain rule—to propagate the output error backward, calculating gradients for weights and biases in multiple layers.	Failed to classify the minority class (0.00 Recall/Precision) due to the severe class imbalance and simple architecture, demonstrating its sensitivity to data distribution.

Importance in Real-World ML Tasks

- **Vectorization and Efficiency:** The necessity of using vectorized NumPy operations highlights that performance in real-world ML directly depends on abstracting mathematical operations into efficient matrix computations.
- **Metric Awareness:** The disparity between the misleading 99% overall accuracy and the poor F1-scores for the minority class emphasizes a critical lesson: in real-world tasks with skewed data (like fraud detection or rare disease diagnosis), **simple accuracy is insufficient**, and domain-appropriate metrics like Precision and Recall must be prioritized.

6. References

- [1] Christopher M. Bishop. (2006). Pattern Recognition and Machine Learning. Classic reference for theoretical foundations of classification and learning algorithms.
- [2] Rosenblatt, F. (1958). The Perceptron: A Probabilistic Model for Information Storage and Organization in The Brain. Psychological Review, 65(6), 386–408. Original Perceptron paper.
- [3] Aurélien Géron . Hands-On Machine Learning with Scikit-Learn, *Keras*, and *TensorFlow*, 2nd Edition, O'Reilly, 2019. Practical guide with examples of implementing ML algorithms, including from scratch.
- [4] R. O. Duda, P. E. Hart, D. G. Stork . Pattern Classification, 2nd Edition, Wiley, 2001. Contains detailed explanations of Perceptron and statistical classifiers.