

Arjun Prakash

Student Id: 21239525

Course: MSc. AI (1MAI)

Assignment 2 - CT5165 Principles of Machine Learning

Description: Neural Network

The developed Neural network consists of two hidden layers with 6 neurons in each hidden layer and one output layer. The model takes input of n-dimension and classifies it into 2 classes, 0 or 1. Language used: Python. The model consists of 7 functions as mentioned below.

1. __init__(x,y)

This function takes in 2 parameters: (input and output labels) and initializes it to the object's input and output variable. It initializes random weights (for layer 1 to 2) of dimension (n,6) where n value depends on the input data dimension, random weights of dimension (6,6) for layer 2 to 3, random weights of dimension (6,1) for the layer 3 to 4. It also initializes loss variable for calculating the loss for each epoch.

2. feed_forward()

Feed forward is done by summing the products of each neuron value in the layer with its weights and passing it through an activation function, sigmoid in this case. The same process is done for every neuron of every layer in the network except the input layer.

Sigmoid function: $f(x) = \frac{1}{1+e^{-x}}$, where $(x) = \sum w \times a^{(L-1)} + bias$, (This is done for all neurons in the layer). Here bias is set to 0.

3. back_propagate()

Back propagation[3] is the algorithm where the model learns. This functions first calculates the loss between the actual 'Y' (label value) and the output layer (predicted value) using Mean squared error.

$$\text{Mean Squared Error (cost)} = \frac{1}{n} \sum_{i=0}^n (y_i - \hat{y})^2$$

After calculating the error, we backpropagate and change the weights with respect to the error. To do this, we need to compute the ratio of change in cost and change in weights $\left(\frac{\partial C_0}{\partial w^{(L)}}\right)$

$$\begin{aligned} C_0 &= (a^{(L)} - y)^2, & \text{where } a &= \text{activation of output layer and } y = \text{actual class label} \\ Z^{(L)} &= w^{(L)} a^{(L-1)} + b, & \text{where } w &= \text{weights of layer 'L' and 'a' is the activation (previous layer), } b = \text{bias} \\ a^{(L)} &= \sigma(Z^{(L)}), & \text{where } a &= \text{activation at layer 'L' and } \sigma = \text{sigmoid function as stated above.} \end{aligned}$$

To find change in cost w.r.t (With respect to) weights (or how sensitive the cost is with change in weights) in a Layer, we use chain rule which includes the change in all the intermediate variables like 'Z', activation, and weights. The change in weights $w^{(L)}$, causes some change in $Z^{(L)}$ which in-turn causes some changes in the activation $a^{(L)}$ which influences the cost C_0 . To find this [2][3][4]:

$$\frac{\partial C_0}{\partial w^{(L)}} = \frac{\partial Z^{(L)}}{\partial w^{(L)}} \frac{\partial a^{(L)}}{\partial Z^{(L)}} \frac{\partial C_0}{\partial a^{(L)}} \quad (\text{Using chain rule})$$

Where, in general:

Arjun Prakash

Student Id: 21239525

Course: MSc. AI (1MAI)

$$\frac{\partial z^{(L)}}{\partial w^{(L)}} = a^{(L-1)} \quad : \text{Activation of previous layer.}$$

$$\frac{\partial a^{(L)}}{\partial z^{(L)}} = \sigma'(z^{(L)}) \quad : \text{Derivative of sigmoid function, which is } \sigma'(x) = x(1 - x).$$

$$\frac{\partial C_0}{\partial a^{(L)}} = 2(a^{(L)} - y) \quad : \text{Derivative of cost function, where } a^{(L)} \text{ is the predicted output layer.}$$

The same is applied for all the layer and each neuron.

The Change in last layer weights w.r.t cost is calculated by (Applying chain rule as mentioned above):

$$\partial w_3^* = 2(a^{(L^3)} - y) \sigma'(L^3), \text{ where } L^3 \text{ is output layer, 'a' is activation, y is actual label}$$

$$\partial w_3 = a^{(L^3-1)} \cdot \partial w_3^*, \text{ These change in weights are added with the network's weight 3}$$

For all other weights, the loss will be calculated by:

$$\partial w_2^* = (\partial w_3^* \cdot w_3) \sigma'(L^2), \quad \text{where } L^2 \text{ is layer 2, 'w}_3 \text{' is weights of 3 – 4th layer}$$

$$\partial w_2 = a^{(L^2-1)} \cdot \partial w_2^*, \quad \text{These change in weights are added with the network's weight 2}$$

$$\partial w_1^* = (\partial w_2^* \cdot w_2) \sigma'(L^1), \quad \text{where } L^1 = \text{Layer 1, 'w}_1 \text{' is weights of 2 – 3th layer}$$

$$\partial w_1 = a^{(inp)} \cdot \partial w_1^*, \quad \text{These change in weights are added with the network's weight 1}$$

All the above changes in weights are then added with the actual weights of respectively layer. These changes in weights are done in every epoch. Updated loss is then saved in the loss (list), every epoch for further visualization. Loss function used for visualization: **Mean Squared Logarithmic Error Loss.**

4. fit_nn(epochs)

This function takes in epochs as parameter and calls the function feed_forward and back_propagate till the range of given epochs. The back_propagate function updates the weights whenever its called.

5. predict_nn(test_input)

The predict function takes in test input as a parameter and performs similar steps as feedforward function but using the test label as input and using the trained weights. Then it returns the predicted class label.

6. getweights()

Returns the latest trained weights of the model.

7. __str__()

The str() method prints the current network shape of each layer in below format:

Network is of shape: (136, 9), (6,), (6,), (1,)

Arjun Prakash

Student Id: 21239525

Course: MSc. AI (1MAI)

Preprocessing :

- Removed inconsistency from the dataset by removing spaces/tabs and converting the txt file to csv file.
- Replaced labels in dataframe from yes/no to 1/0.
- Dropped the label column from the data for training and testing set. Split the data into 2/3rd training and 1/3rd testing.
- Normalized the data using StandardScaler() [6], which scales the data and makes the mean of data as 0 and standard deviation as 1.

Training:

Own Model from Scratch:

Implemented a neural network with 2 hidden layers, with number of neurons as 6 in both the layers and activation function as Sigmoid.

Trained the classifier (using fit() function) using backpropagation with normalized training data for 1800 epochs. Stored the loss of each epoch and divided the number with the size of training data (completely optional) to make the loss value more normalized. Stored the F1 score for 10 different training samples. Also added an option to pickle the trained weights[5].

Sklearn Implementation:

Implemented MLPClassifier using same parameters. Hidden layers = 2 with same number of neurons as in scratch implementation and activation function as “tanh”(sigmoid).

Trained the classifier using the same normalized data that was used for the scratch implementation with the same epoch (1800) range. Stored the F1 score for 10 different training samples.

Testing:

Own Model from Scratch, F1 score:

Achieved F1 score in the range (testing data) = 0.8 – 0.94+

Average F1 score over 10 different training = 0.87

Sklearn Implementation F1 score:

Achieved F1 score in the range (testing data) = 0.78 – 0.94+

Average F1 score over 10 different training = 0.85

Observation:

Could see that Sklearn and the own implemented model had similar F1 Scores, for example, at any given point, if the sklearn model classified the data at 0.78 f1 score then the own implemented model also had the score around the same range (0.77-0.82) and when the score of sklearn was around 0.92, then the

Arjun Prakash

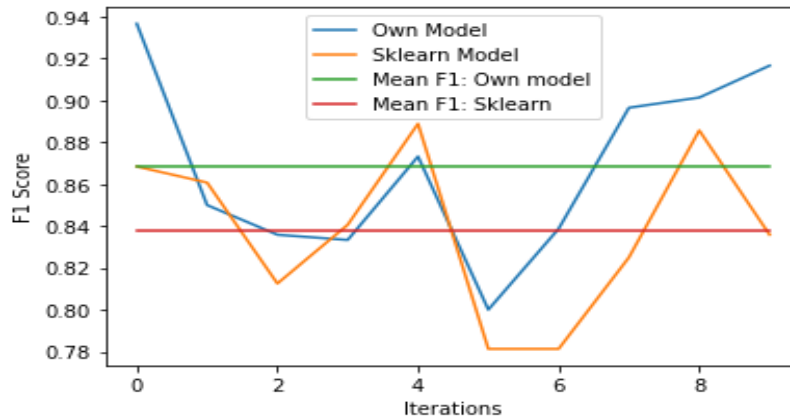
Student Id: 21239525

Course: MSc. AI (1MAI)

implemented model had the score in range (0.91-0.94). On average, most of the time, the sklearn model had an average f1 score of 2-3% below the f1 score of built from scratch model (For 10 distinct training samples).

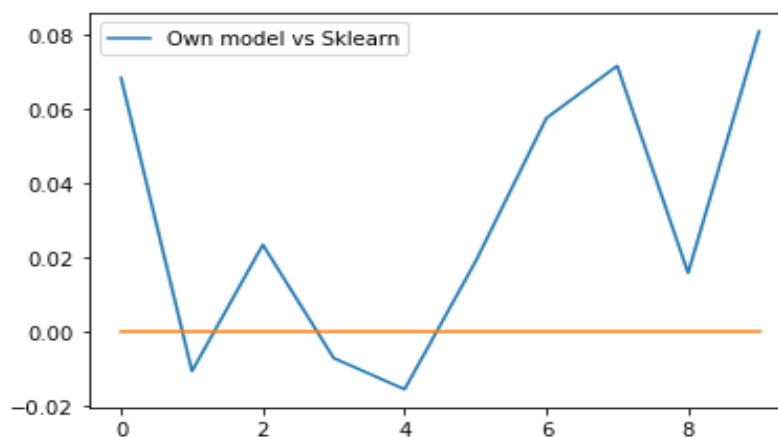
This can be seen in the below visualizations.

1. Model comparison



Here the Blue line represents the f1 score of the own built from scratch model and the orange line represents the f1 score of sklearn MLPClassifier over 10 different trainings. The Green line is the mean f1 score of the built from scratch model and the red line is the mean f1 score of sklearn MLPClassifier. Could see similar results every time the model was tested.

2. Difference in the accuracy/score



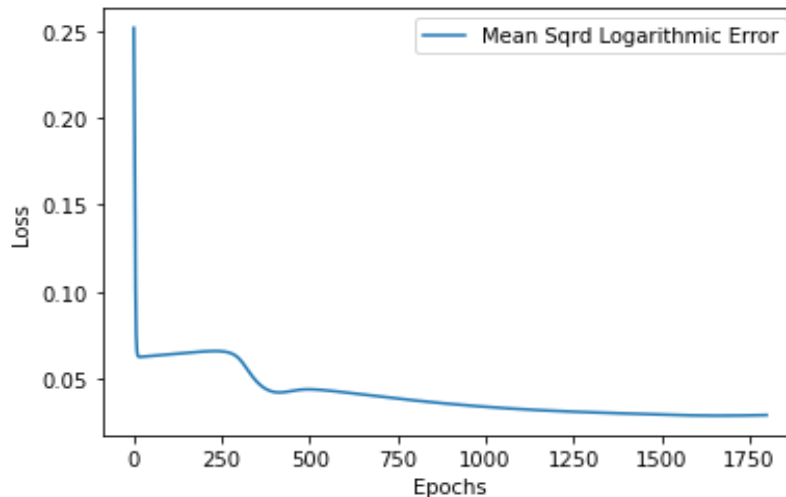
In this graph, the difference in both the models are calculated. When both the model classifies with the same accuracy/f1-score then their difference is 0, which means, if the datapoint touches the yellow line then both the models have same f1 score. For points below the yellow line, Sklearn classifier had better accuracy. Also, farther the data point is from the yellow line, higher the magnitude of difference in the f1 score of both the models.

Arjun Prakash

Student Id: 21239525

Course: MSc. AI (1MAI)

3. Loss function:



With the above loss function graph, we see that the loss is highest at first few epochs as the model has random weights assigned to it, thus the classification accuracy is not good initially. As the model gets trained, we see the loss function (or cost function) dropping till it saturates. (The loss is summed at axis=1 and multiplied with -1 to make the loss value positive)

Additional

- Exported classified data to file 'Predictions_x.csv' for all the 10 training iterations.
- To get the pickled weights data, uncomment the pickling line in training module, after the prediction file export. Pickle File name: 'model_x.pkl'.

References:

- [1] <https://www.youtube.com/watch?v=WUvTyaaNkzM> – 3blue1brown - Calculus
- [2] <https://www.youtube.com/watch?v=t1eHLnjs5U8> – 3blue1brown - Backpropagation calculus
- [3] <https://www.youtube.com/watch?v=llg3gGewQ5U> - 3blue1brown -Backpropagation
- [4] <https://www.youtube.com/watch?v=aircAruvnKk> - 3blue1brown – Neural network
- [5] <https://docs.python.org/3/library/pickle.html> - Pickling
- [6] <https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.StandardScaler.html>-scalar transformations

Appendix (code):