

CS 6375

ASSIGNMENT _____2_____

Names of students in your group:

Ching-Yi Chang - CXC190002

Xiaokai Rong - XXR230000

Number of free late days used: _____0_____

Note: You are allowed a **total** of 4 free late days for the **entire semester**. You can use at most 2 for each assignment. After that, there will be a penalty of 10% for each late day.

Please list clearly all the sources/references that you have used in this assignment.

UCI ML Repository: <https://archive.ics.uci.edu/dataset/9/auto+mpg> Auto MPG dataset

1. Theoretical Part (40 points)

1.1. Revisiting Backpropagation Algorithm

In class we had derived the backpropagation algorithm for the case where each of the hidden and output layer neurons used the sigmoid activation function:

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

Revise the backpropagation algorithm for the case where each hidden and output layer neuron uses the

a. tanh activation function

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

b. ReLu activation function:

$$\text{ReLu}(x) = \max(0, x)$$

Show all steps of your derivation and the final equation for output layer and hidden layers.

A:

To derive the gradient descent update rule for the alternative error function $E(\vec{w})$, we will calculate the partial derivatives of $E(\vec{w})$ with respect to the weights.

By chain rule:

$$\frac{\partial E}{\partial(w_{ji})} = \frac{\partial E}{\partial(o_j)} \frac{\partial(o_j)}{\partial(net_j)} \frac{\partial(net_j)}{\partial(w_{ji})}$$

First, we need to calculate the derivative of the error function with respect to the output of the network, denoted as o_j .

$$\frac{dE}{d(o_j)} = \frac{\partial}{\partial o_j} \left[\frac{1}{2} \sum_{k \in \text{outputs}} (t_k - o_k)^2 \right] = \frac{d \left[\frac{1}{2} (t_j - o_j)^2 \right]}{d(o_j)} = -(t_j - o_j)$$

Next, we calculate the derivative of the output with respect to the input of the output layer, denoted as net_j . Here, f' represents the derivative of the activation function used in the output layer.

$$\frac{\partial(o_j)}{\partial(net_j)} = \frac{\partial f(net_j)}{\partial(net_j)} = f'(net_j)$$

Now, let's calculate the derivative of the input of the output layer, net_j , with respect to the weights connecting the hidden layer to the output layer, denoted as w_{ji} .

$$\frac{\partial(net_j)}{\partial(w_{ji})} = \frac{\partial(\sum_i w_{ji}x_{ji})}{\partial(w_{ji})} = x_{ji}$$

Applying the chain rule, we can calculate the derivative of the error function with respect to the weights connecting the hidden layer to the output layer.

$$\frac{\partial E}{\partial(w_{ji})} = \frac{\partial E}{\partial(o_j)} \frac{\partial(o_j)}{\partial(net_j)} \frac{\partial(net_j)}{\partial(w_{ji})} = -(t_j - o_j) f'(net_j) x_{ji}$$

$$\Delta \vec{w} = -\eta \frac{\partial E}{\partial w}$$

$$\Delta w_{ji} = \eta(t_j - o_j) f'(net_j) x_{ji} = \eta \delta_j x_{ji}$$

where δ_j represents the error term for the j-th output neuron, o_j is the output of the j-th output neuron, t_j is the target output, and net_j is the weighted sum of inputs to the j-th output neuron, Δw_{ji} represents the change in weight from the i-th hidden neuron to the j-th output neuron, η is the learning rate.

Case1: When j is an output unit

$$\frac{\partial E}{\partial(net_j)} = -\delta_j = -(t_j - o_j) f'(net_j)$$

$$\delta_j = (t_j - o_j) f'(net_j)$$

Case2: When j is a hidden unit

$$\frac{\partial E}{\partial(net_j)} = \sum_{k \in \text{Downstream}(j)} -\delta_k \frac{\partial net_k}{\partial o_j} \frac{\partial o_j}{\partial net_j} = \sum_{k \in \text{Downstream}(j)} -\delta_k w_{kj} f'(net_j)$$

$$\delta_j = f'(net_j) \sum_{k \in \text{Downstream}(j)} \delta_k w_{kj}$$

Unit j sends its output to Downstream(j) units. δ_k is the delta for the unit k which is part of Downstream(j) and represents the error term for the i-th hidden neuron

a. tanh activation function:

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

The backpropagation algorithm involves updating the weights and biases based on the gradients of the loss function with respect to these parameters. We'll derive the equations for the output layer and hidden layers using the tanh activation function.

For the output layer:

Calculate the error term for each output neuron:

$$\begin{aligned} f'(\tanh(z)) &= \frac{\frac{\partial}{\partial z} \sinh(z)}{\frac{\partial}{\partial z} \cosh(z)} \\ &= \frac{\frac{\partial}{\partial z} \sinh(z) \cosh(z) - \frac{\partial}{\partial z} \cosh(z) \sinh(z)}{\cosh^2(z)} \\ &= \frac{\cosh^2(z) - \sinh^2(z)}{\cosh^2(z)} \\ &= 1 - \frac{\sinh^2(z)}{\cosh^2(z)} \\ &= 1 - \tanh^2(z) \end{aligned}$$

$$\delta_j = (t_j - o_j) f'(\tanh(net_j)) = (t_j - o_j) (1 - \tanh^2(net_j))$$

For the hidden layers:

Calculate the error term for each hidden neuron:

$$\delta_j = (1 - \tanh^2(net_j)) \sum_{k \in \text{Downstream}(j)} \delta_k w_{kj}$$

b. ReLU activation function:

$$\text{ReLU}(x) = \max(0, x)$$

$$f'(\text{ReLu}(z)) = 1 \text{ when } z \geq 0$$

$$f'(\text{ReLu}(z)) = 0 \text{ when } z < 0$$

The derivation for the ReLU activation function follows a similar process. We can use a subgradient instead.

For the output layer:

Calculate the error term for each output neuron:

$$\delta_j = (t_j - o_j)f'(net_j) = (t_j - o_j) \text{ when } net_j \geq 0$$

$$\delta_j = (t_j - o_j)f'(net_j) = 0 \text{ when } net_j < 0$$

For the hidden layers:

Calculate the error term for each hidden neuron:

$$\delta_j = f'(net_j) \sum_{k \in \text{Downstream}(j)} \delta_k w_{kj} = \sum_{k \in \text{Downstream}(j)} \delta_k w_{kj} \text{ when } net_j \geq 0$$

$$\delta_j = 0 \text{ when } net_j < 0$$

1.2. Gradient Descent

Derive a gradient descent training rule for a single unit neuron with output o , defined as:

$$o = w_0 + w_1(x_1 + x_1^2) + \dots + w_n(x_n + x_n^2)$$

where x_1, x_2, \dots, x_n are the inputs, w_1, w_2, \dots, w_n are the corresponding weights, and w_0 is the bias weight.

Note that the above is the final output of the neuron, you don't need to consider any activation function. Show all steps of your derivation and the final result for weight update. You can assume a learning rate of η .

A: To derive the gradient descent training rule for the given single-unit neuron with output o and the target output $t^{(i)}$, we'll need to calculate the partial derivatives of the output with respect to each weight. Let's go through the derivation step by step:

$$o = w_0 + w_1(x_1 + x_1^2) + \dots + w_n(x_n + x_n^2)$$

$$\text{Generalized cost function } E(w) = \frac{1}{2m} \sum_{i=1}^m (o_w(x^{(i)}) - t^{(i)})^2$$

$$\begin{aligned} \frac{\partial}{\partial w_j} E(w) &= \frac{\partial}{\partial w_j} \frac{1}{2m} \sum_{i=1}^m (o_w(x^{(i)}) - t^{(i)})^2 \\ &= \frac{\partial}{\partial w_j} \frac{1}{2m} \sum_{i=1}^m \left(\left(w_0 + w_1(x_1^{(i)} + x_1^{2(i)}) + \dots + w_n(x_n^{(i)} + x_n^{2(i)}) \right) - t^{(i)} \right)^2 \\ &= \frac{1}{m} \sum_{i=1}^m (o_w(x^{(i)}) - t^{(i)}) (x_j^{(i)} + x_j^{2(i)}) \end{aligned}$$

Initialize the weight update rule using gradient descent:

$$\Delta w_j = -\eta \frac{\partial E(w)}{\partial w_j} = -\eta \frac{1}{m} \sum_{i=1}^m (o_w(x^{(i)}) - t^{(i)}) (x_j^{(i)} + x_j^{2(i)})$$

Calculate the derivative of the output with respect to each weight

The final result for weight update for each weight w_j

$$w_j = w_j + \Delta w_j = w_j - \eta \frac{1}{m} \sum_{i=1}^m (o_w(x^{(i)}) - t^{(i)}) (x_j^{(i)} + x_j^{2(i)})$$

This rule allows you to update each weight iteratively during the training process using gradient descent.

1.3. Comparing Activation Function

Consider a neural net with 2 input layer neurons, one hidden layer with 2 neurons, and 1 output layer neuron as shown in Figure 1. Assume that the input layer uses the identity activation function i.e. $f(x) = x$, and each of the hidden layers and output layer use an activation function $h(x)$. The weights of each of the connections are marked in the figure.

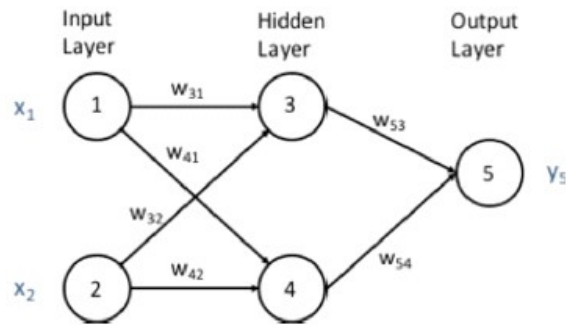


Figure 1: A neural net with 1 hidden layer having 2 neurons

a. Write down the output of the neural net y_5 in terms of weights, inputs, and a general activation function $h(x)$.

The input layer uses the identity activation function i.e. $f(x) = x$

The output of the neural net y_5

$$X_3 = W_{31}X_1 + W_{32}X_2$$

$$X_4 = W_{41}X_1 + W_{42}X_2$$

$$X_5 = W_{53}h(X_3) + W_{54}h(X_4)$$

$$= W_{53}h(W_{31}X_1 + W_{32}X_2) + W_{54}h(W_{41}X_1 + W_{42}X_2)$$

$$y_5 = h(X_5) = h(W_{53}h(W_{31}X_1 + W_{32}X_2) + W_{54}h(W_{41}X_1 + W_{42}X_2))$$

b. Now suppose we use vector notation, with symbols defined as below:

$$X = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$$

$$W^{(1)} = \begin{pmatrix} w_{3,1} & w_{3,2} \\ w_{4,1} & w_{4,2} \end{pmatrix}$$

$$W^{(2)} = (w_{5,3} \ w_{5,4})$$

Write down the output of the neural net in vector format using above vectors.

Calculate the hidden layer and output layer.

$$X_3 = W^{(1)}[0][0]X[0] + W^{(1)}[0][1]X[1]$$

$$X_4 = W^{(1)}[1][0]X[0] + W^{(1)}[1][1]X[1]$$

$$X_5 = (W^{(2)}[0]h(W^{(1)}[0][0]X[0] + W^{(1)}[0][1]X[1]) \\ + W^{(2)}[1]h(W^{(1)}[1][0]X[0] + W^{(1)}[1][1]X[1]))$$

$$y_5 = h(X_5) = h((W^{(2)}[0]h(W^{(1)}[0][0]X[0] + W^{(1)}[0][1]X[1]) \\ + W^{(2)}[1]h(W^{(1)}[1][0]X[0] + W^{(1)}[1][1]X[1])))$$

c. Now suppose that you have two choices for activation function $h(x)$, as shown below:

Sigmoid:

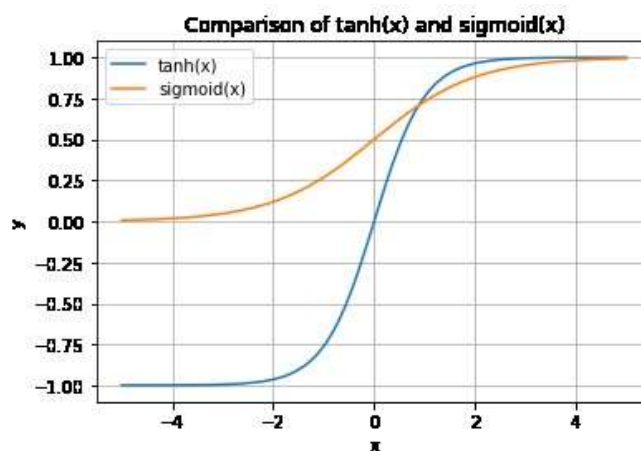
$$h_s(x) = \frac{1}{1 + e^{-x}}$$

Tanh:

$$h_t(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

Show that neural nets created using the above two activation functions can generate the same function.

Hint: First compute the relationship between $h_s(x)$ and $h_t(x)$ and then show that the output functions are same, with the parameters differing only by linear transformations and constants.



Sigmoid and Tanh:

Both $\tanh(x)$ and $\text{sigmoid}(x)$ have an S-shaped curve, similar shape.

Range of Output: $\tanh(x)$ has a range of $[-1, 1]$, $\text{sigmoid}(x)$ has a range of $(0, 1)$. In the range $(0,1)$, the curves of $\tanh(x)$ and $\text{sigmoid}(x)$ will be very close to each other, exhibiting similar behavior.

The logistic sigmoid function is symmetric around the origin.

$$1 - \sigma(x) = \sigma(-x)$$
$$1 - \frac{1}{1 + e^{-x}} = \frac{1}{1 + e^x}$$

The tanh function into a similar form by:

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} = \frac{e^x + e^{-x} - 2e^{-x}}{e^x + e^{-x}} = 1 + \frac{-2e^{-x}}{e^x + e^{-x}} = 1 - \frac{2}{e^{2x} + 1}$$

From the logistic sigmoid's,

$$\tanh(x) = 1 - \frac{2}{e^{2x} + 1} = 1 - 2\sigma(-2x) = 1 - 2(1 - \sigma(2x)) = 2\sigma(2x) - 1$$

Programming Part (60 points)

In this part, you will code a neural network (NN) having **at least one hidden layers, besides the input and output layers**. You are required to pre-process the data and then run the processed data through your neural net. Below are the requirements and suggested steps of the program

- The programming language for this assignment will be Python 3.x
- **You cannot use any libraries for neural net creation.** You are free to use any other libraries for data loading, pre-processing, splitting, model evaluation, plotting, etc.
- Your code should be in the form of a Python class with methods like pre-process, train, test within the class. I leave the other details up to you.
- As the first step, pre-process and clean your dataset. There should be a method that does this.
- Split the pre-processed dataset into training and testing parts. You are free to choose any reasonable value for the train/test ratio, but be sure to mention it in the README file.
- Code a neural net having **at least one hidden layer**. You are free to select **the number of neurons in each layer**. Each neuron in the hidden and output layers should have a bias connection.
- You are required to add an optimizer on top of the basic backpropagation algorithm. This could be the one you selected in the previous assignment or a new one. Some good resources for gradient descent optimizers are:

<https://arxiv.org/pdf/1609.04747.pdf>

<https://ruder.io/optimizing-gradient-descent/>

<https://towardsdatascience.com/10-gradient-descent-optimisation-algorithms-86989510b5e9>

- You are **required to code three different activation functions**:
 1. Sigmoid
 2. Tanh
 3. ReLu

The earlier part of this assignment may prove useful for this stage. The activation function should be a parameter in your code.

- Code a method for creating a neural net model from the training part of the dataset. Report the training accuracy.
- Apply the trained model on the test part of the dataset. Report the test accuracy.
- You have to tune model parameters like learning rate, activation functions, etc. Report your results in a tabular format, with a column indicating the parameters used, a column for training accuracy, and one for test accuracy.

Dataset

You can use any one dataset from the UCI ML repository, or any other standard repository such as Kaggle.

What to submit

You need to submit the following for the programming part:

- Link to the dataset used. Please do not include the data as part of your submission.
- Your source code and a README file indicating how to run your code. Do not hardcode any paths to your local computer. It is fine to code any public paths, such as AWS S3.
- Output for your dataset summarized in a tabular format for different combination of parameters
- A brief report summarizing your results. For example, which activation function performed the best and why do you think so.
- Any assumptions that you made.

Report:

Use Python 3.8

Step 1: Load the dataset

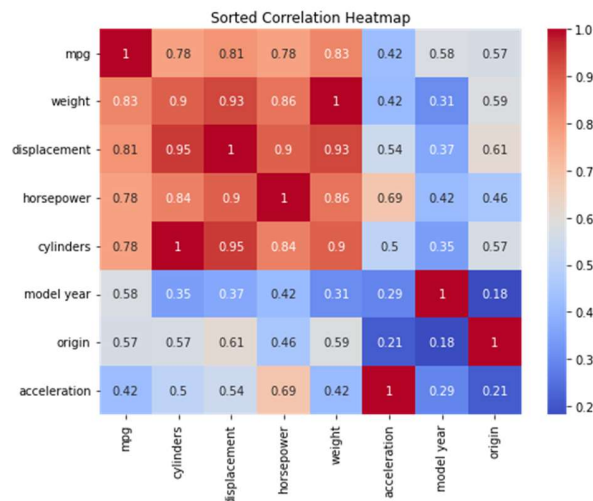
```
url = 'https://raw.githubusercontent.com/SparrowChang/CS6375_assignment1/main/auto%20mpg/auto-mpg.data'
```

Read the CSV file from the URL into a DataFrame, in the original dataset, total with 8 features ('cylinders', 'displacement', 'horsepower', 'weight', 'acceleration', 'model year', 'origin', 'car name'), 1 target('mpg'). MPG means miles per gallon, the number of miles your vehicle can travel down the streets of Deltona on one gallon of gas.

Step 2: Pre-processing

- Convert categorical variables to numerical variables
- Drop the columns 'car name' that are not relevant for the regression analysis
- Create a StandardScaler object and normalized the dataframe
- Remove null or NA values, dropna()
- Remove redundant rows
- Calculate the correlation matrix. Sort absolute correlation with 'mpg'. We would know 'weight', 'displacement', 'horsepower', 'cylinders' are high correlation with target.

Finally, the dataset is with 7 features ('cylinders', 'displacement', 'horsepower', 'weight', 'acceleration', 'model year', 'origin') and 392 instances.



Step 3: Split the dataset into training and test sets

We split the training and test sets to 80/20 (test_size = 0.2, random_state =42)

Step 4: To define Activation Functions

- sigmoid(x): $1 / (1 + \exp(-x))$ and its derivative: $\text{sigmoid}(x) * (1 - \text{sigmoid}(x))$
- tanh(x): $\tanh(x)$ and its derivative: $1 - \tanh(x)^2$
- relu(x): $\max(0, x)$ and its derivative: $1 (x > 0)$

Step 5: To define Neural Network Class

```
class NeuralNetwork:
    def __init__(self, input_size, hidden_size, output_size, activation):
        self.input_size = input_size
        self.hidden_size = hidden_size
        self.output_size = output_size
        self.activation = activation

        # Initialize weights and biases
        self.W1 = np.random.randn(self.input_size, self.hidden_size)
        self.b1 = np.zeros((1, self.hidden_size))
        self.W2 = np.random.randn(self.hidden_size, self.output_size)
        self.b2 = np.zeros((1, self.output_size))

    def forward(self, X):
        # Input to hidden layer
        self.z1 = np.dot(X, self.W1) + self.b1

        # Apply activation function to the hidden layer
        if self.activation == 'sigmoid':
            self.a1 = sigmoid(self.z1)
        elif self.activation == 'tanh':
            self.a1 = tanh(self.z1)
        elif self.activation == 'relu':
            self.a1 = relu(self.z1)

        # Hidden to output layer
        self.z2 = np.dot(self.a1, self.W2) + self.b2
        self.output = self.z2 # Since it's a regression problem, no activation on output

        return self.output
```

```
def backward(self, X, y, output):
    # Compute the loss and the derivative of the loss with respect to output
    loss = np.mean((output - y)**2)
    d_loss_output = 2 * (output - y) / X.shape[0]

    # Backpropagate the error
    if self.activation == 'sigmoid':
        d_output_z2 = sigmoid_derivative(output)
    elif self.activation == 'tanh':
        d_output_z2 = tanh_derivative(output)
    elif self.activation == 'relu':
        d_output_z2 = relu_derivative(output)

    d_loss_z2 = d_loss_output * d_output_z2
    d_loss_a1 = np.dot(d_loss_z2, self.W2.T)
    d_loss_z1 = d_loss_a1 * sigmoid_derivative(self.a1) # Use sigmoid derivative for sigmoid activ
```

```

        # Compute the gradients
        d_loss_W2 = np.dot(self.a1.T, d_loss_z2)
        d_loss_b2 = np.sum(d_loss_z2, axis=0, keepdims=True)
        d_loss_W1 = np.dot(X.T, d_loss_z1)
        d_loss_b1 = np.sum(d_loss_z1, axis=0, keepdims=True)

        return loss, d_loss_W1, d_loss_b1, d_loss_W2, d_loss_b2

    def update_weights(self, d_loss_W1, d_loss_b1, d_loss_W2, d_loss_b2, learning_rate):
        # Update weights and biases using gradient descent
        self.W1 -= learning_rate * d_loss_W1
        self.b1 -= learning_rate * d_loss_b1
        self.W2 -= learning_rate * d_loss_W2
        self.b2 -= learning_rate * d_loss_b2

```

```

def train(self, X, y, learning_rate, num_epochs, tolerance):
    for epoch in range(num_epochs):
        # Forward pass
        output = self.forward(X)

        # Backward pass
        loss, d_loss_W1, d_loss_b1, d_loss_W2, d_loss_b2 = self.backward(X, y, output)

        # Update weights and biases
        self.update_weights(d_loss_W1, d_loss_b1, d_loss_W2, d_loss_b2, learning_rate)

        if (epoch + 1) % 100 == 0:
            print(f'Epoch {epoch+1}/{num_epochs}, Loss: {loss:.2f}')

        # check if the loss is below the tolerance level
        if loss < tolerance:
            print('Training converged!')
            break

    print('Training completed!')

def predict(self, X):
    output = self.forward(X)
    return output

```

Step 6: To set up Neural Network model

We create 1 hidden layer (64 hidden units) and 1 output layer (1 output units). In each layer, with its own weights and bias. 3 activations (sigmoid, tanh, relu)

For hyperparameters:

epochs: 100, 500, 1000.

learning rate: 0.001, 0.01, 0.1, 0.5

```

iteration_range = [100, 500, 1000]
learning_rate_range = [0.001, 0.01, 0.1, 0.5]
activ = ['sigmoid', 'tanh', 'relu']

```

Step 7: train and evaluate test result

```

model = NeuralNetwork(input_size, hidden_size, output_size, activation=activ[1])

# Train the model
model.train(X_train, y_train, learning_rate, num_epochs, tolerance)

```

```

# Evaluate the model
train_predictions = model.predict(X_train).flatten()
train_loss = np.mean((train_predictions - y_train.flatten())**2)
train_accuracy = np.mean(np.abs(train_predictions - y_train.flatten()))

test_predictions = model.predict(X_test).flatten()
test_loss = np.mean((test_predictions - y_test.flatten())**2)
test_accuracy = np.mean(np.abs(test_predictions - y_test.flatten()))

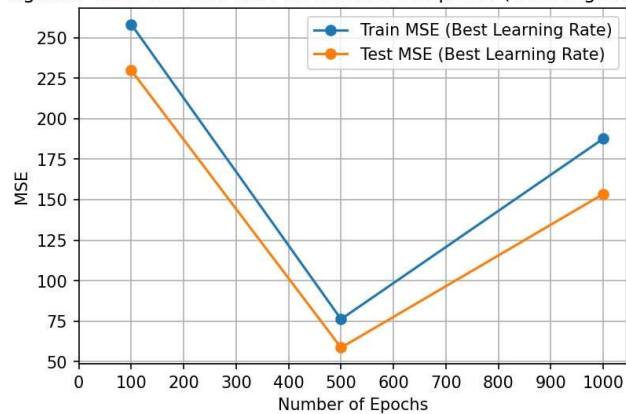
```

```

sigmoid Learning Rate: 0.001, Num Epochs: 100, Train Loss: 433.08, Train Accuracy: 19.11, Test Loss: 413.88, Test Accuracy: 18.88
-----
sigmoid Learning Rate: 0.001, Num Epochs: 500, Train Loss: 347.36, Train Accuracy: 16.72, Test Loss: 299.86, Test Accuracy: 15.70
-----
sigmoid Learning Rate: 0.001, Num Epochs: 1000, Train Loss: 254.82, Train Accuracy: 14.28, Test Loss: 231.58, Test Accuracy: 13.85
-----
sigmoid Learning Rate: 0.01, Num Epochs: 100, Train Loss: 325.55, Train Accuracy: 16.45, Test Loss: 291.38, Test Accuracy: 15.67
-----
sigmoid Learning Rate: 0.01, Num Epochs: 500, Train Loss: 269.90, Train Accuracy: 14.73, Test Loss: 227.94, Test Accuracy: 13.64
-----
sigmoid Learning Rate: 0.01, Num Epochs: 1000, Train Loss: 242.60, Train Accuracy: 13.48, Test Loss: 223.47, Test Accuracy: 13.20
-----
sigmoid Learning Rate: 0.1, Num Epochs: 100, Train Loss: 258.28, Train Accuracy: 14.21, Test Loss: 229.88, Test Accuracy: 13.60
-----
sigmoid Learning Rate: 0.1, Num Epochs: 500, Train Loss: 76.20, Train Accuracy: 6.32, Test Loss: 58.62, Test Accuracy: 5.55
-----
sigmoid Learning Rate: 0.1, Num Epochs: 1000, Train Loss: 187.66, Train Accuracy: 11.22, Test Loss: 153.52, Test Accuracy: 10.12
-----
sigmoid Learning Rate: 0.5, Num Epochs: 100, Train Loss: 2327.21, Train Accuracy: 47.28, Test Loss: 2336.21, Test Accuracy: 47.56
-----
sigmoid Learning Rate: 0.5, Num Epochs: 500, Train Loss: 478.53, Train Accuracy: 17.80, Test Loss: 452.05, Test Accuracy: 16.56
-----
sigmoid Learning Rate: 0.5, Num Epochs: 1000, Train Loss: 206.68, Train Accuracy: 12.98, Test Loss: 226.90, Test Accuracy: 13.96
-----
Best Mean squared error (MSE):58.62
Best parameters: {'learning_rate': 0.1, 'num_epochs': 500}
-----

```

sigmoid Train and Test MSE vs. Number of Epochs (Learning Rate: 0.1)

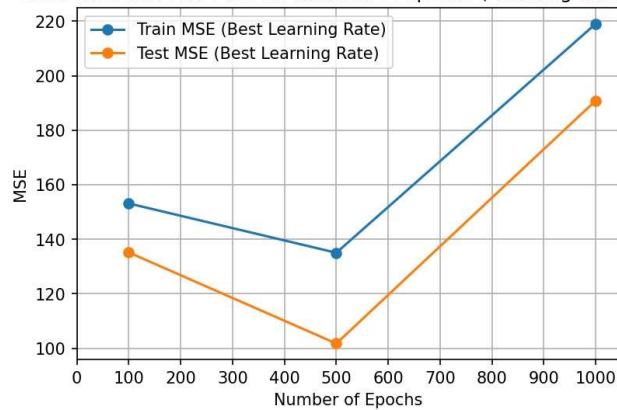



```

tanh Learning Rate: 0.001, Num Epochs: 100, Train Loss: 231.46, Train Accuracy: 12.65, Test Loss: 209.87, Test Accuracy: 12.30
-----
tanh Learning Rate: 0.001, Num Epochs: 500, Train Loss: 287.76, Train Accuracy: 15.06, Test Loss: 260.48, Test Accuracy: 14.69
-----
tanh Learning Rate: 0.001, Num Epochs: 1000, Train Loss: 269.16, Train Accuracy: 14.09, Test Loss: 243.44, Test Accuracy: 13.61
-----
tanh Learning Rate: 0.01, Num Epochs: 100, Train Loss: 253.81, Train Accuracy: 13.12, Test Loss: 215.33, Test Accuracy: 12.39
-----
tanh Learning Rate: 0.01, Num Epochs: 500, Train Loss: 289.12, Train Accuracy: 13.80, Test Loss: 257.49, Test Accuracy: 13.25
-----
tanh Learning Rate: 0.01, Num Epochs: 1000, Train Loss: 266.76, Train Accuracy: 12.69, Test Loss: 214.48, Test Accuracy: 11.67
-----
tanh Learning Rate: 0.1, Num Epochs: 100, Train Loss: 153.14, Train Accuracy: 9.77, Test Loss: 135.16, Test Accuracy: 9.70
-----
tanh Learning Rate: 0.1, Num Epochs: 500, Train Loss: 135.05, Train Accuracy: 9.00, Test Loss: 101.75, Test Accuracy: 7.83
-----
tanh Learning Rate: 0.1, Num Epochs: 1000, Train Loss: 219.07, Train Accuracy: 11.91, Test Loss: 190.89, Test Accuracy: 11.52
-----
tanh Learning Rate: 0.5, Num Epochs: 100, Train Loss: 872.79, Train Accuracy: 27.97, Test Loss: 869.75, Test Accuracy: 27.63
-----
tanh Learning Rate: 0.5, Num Epochs: 500, Train Loss: 754.34, Train Accuracy: 26.19, Test Loss: 782.79, Test Accuracy: 26.32
-----
tanh Learning Rate: 0.5, Num Epochs: 1000, Train Loss: 286.96, Train Accuracy: 14.88, Test Loss: 258.81, Test Accuracy: 14.09
-----
Best Mean squared error (MSE):101.75
Best parameters: {'learning_rate': 0.1, 'num_epochs': 500}
-----

```

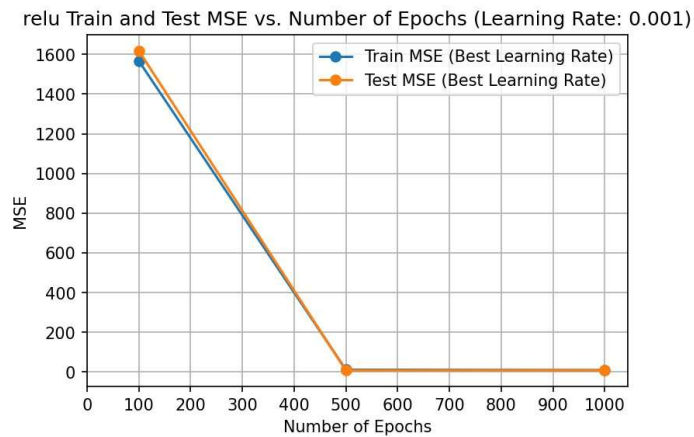
tanh Train and Test MSE vs. Number of Epochs (Learning Rate: 0.1)



```

relu Learning Rate: 0.001, Num Epochs: 100, Train Loss: 1566.87, Train Accuracy: 34.75, Test Loss: 1618.12, Test Accuracy: 35.82
-----
relu Learning Rate: 0.001, Num Epochs: 500, Train Loss: 11.60, Train Accuracy: 2.53, Test Loss: 7.62, Test Accuracy: 2.09
-----
relu Learning Rate: 0.001, Num Epochs: 1000, Train Loss: 9.60, Train Accuracy: 2.35, Test Loss: 8.87, Test Accuracy: 2.32
-----
relu Learning Rate: 0.01, Num Epochs: 100, Train Loss: 10.77, Train Accuracy: 2.41, Test Loss: 10.99, Test Accuracy: 2.32
-----
relu Learning Rate: 0.01, Num Epochs: 500, Train Loss: 28.62, Train Accuracy: 4.58, Test Loss: 30.10, Test Accuracy: 4.71
-----
relu Learning Rate: 0.01, Num Epochs: 1000, Train Loss: 20.25, Train Accuracy: 3.55, Test Loss: 23.41, Test Accuracy: 3.83
-----
relu Learning Rate: 0.1, Num Epochs: 100, Train Loss: 92955287.86, Train Accuracy: 8018.12, Test Loss: 96327574.86, Test Accuracy: 8404.64
-----
relu Learning Rate: 0.1, Num Epochs: 500, Train Loss: 2294.41, Train Accuracy: 46.92, Test Loss: 2144.32, Test Accuracy: 45.35
-----
relu Learning Rate: 0.1, Num Epochs: 1000, Train Loss: 2823269717.64, Train Accuracy: 30560.29, Test Loss: 3348023156.30, Test Accuracy: 31785.55
-----
relu Learning Rate: 0.5, Num Epochs: 100, Train Loss: 366332254753.51, Train Accuracy: 388048.93, Test Loss: 322842765523.19, Test Accuracy: 392196.52
-----
relu Learning Rate: 0.5, Num Epochs: 500, Train Loss: 5781450939.20, Train Accuracy: 59280.89, Test Loss: 5450243051.49, Test Accuracy: 59599.63
-----
relu Learning Rate: 0.5, Num Epochs: 1000, Train Loss: 67321445363041.20, Train Accuracy: 5931569.51, Test Loss: 71302520400159.77, Test Accuracy: 6156249.34
-----
Best Mean squared error (MSE):7.62
Best parameters: {'learning_rate': 0.001, 'num_epochs': 500}
-----
Execution time: 24.97 seconds
Current date and time: 2023-07-03 17:07:40.041273

```

Test MSE of MPG (miles per gallon) comparison											
Activations	sigmoid			tanh			relu			LR by gradient descent	LR by linear regression library 10.71
(R) epochs (D) learning rate	100	500	1000	100	500	1000	100	500	1000	1000	
0.001	413.88	299.86	231.58	209.87	260.48	243.44	1618.12	7.62	8.87		
0.01	291.38	227.94	223.47	215.33	257.49	214.48	10.99	30.10	23.41		
0.1	229.88	58.62	58.62	135.16	101.75	190.89	>10 ⁷	2144.32	>10 ⁷	10.73	
0.5	2336.21	452.05	226.90	869.75	782.79	258.81	>10 ⁷	>10 ⁷	>10 ⁷		

We aim to compare the results with those obtained in Assignment 1. For the prediction of MPG (miles per gallon), the test mean squared error (MSE) achieved using the linear regression library is 10.71. When applying the gradient descent algorithm with 1000 epochs and a learning rate of 0.1, the test MSE is slightly higher at 10.73.

Surprisingly, the ReLU activation function yields a test MSE of 7.62, which not only matches our Assignment 1 result but is even lower. This suggests that the ReLU activation function is effective for this particular task.

However, the test MSE results for MPG using the Sigmoid and Tanh activation functions are unstable. The inconsistency in the test results suggests that these activation functions may not be as suitable for this specific prediction task.

By comparing and analyzing these outcomes, we can gain valuable insights into the performance of different methods and activation functions for the MPG prediction problem.

We have identified a few potential reasons for the observed behavior:

- a. The dataset used for training is relatively small, containing only 392 instances. Limited data can sometimes lead to suboptimal model performance due to insufficient patterns for the algorithm to learn from.
- b. The learning rate used for both the Sigmoid and ReLU activation functions is set to 0.1. This relatively high learning rate may cause the optimization process to converge to local minima instead of the global minimum, impacting the model's overall performance.
- c. It is possible that some features in the dataset may not be highly informative or relevant for the target variable. Utilizing domain knowledge and performing feature selection techniques can potentially improve the model's performance by focusing on the most important features.
- d. The algorithm implemented from scratch is a simple model. Alternatively, incorporating a penalty term, such as regularization techniques like L1 or L2 regularization, could potentially help in reducing the mean squared error (MSE) further and improving the model's generalization capability.

Consider exploring these suggestions to address the observed issues and enhance the performance of the model.