

Data Science 210 Final Project

Thomas Kwashnak

Fall 2021

Contents

1	Background	2
1.1	Goals	2
1.2	Mathematic Background	2
1.3	Dataset	3
2	Design	3
2.1	Sigmoid Activation Function	4
2.2	Feed-Forward Algorithm	4
2.3	Loss Function	4
2.4	0th Layer Weight Derivative	5
2.5	1st Layer Weight Derivative	5
2.6	Backwards Propagation	6
2.7	Stochastic Gradient Descent	6
2.8	Training the Network	7
3	Implementation	8
3.1	Data Importing	8
4	Validation	8
5	Reflection	8

1 Background

1.1 Goals

The main goal of this project is to create an algorithm that will predict whether or not a given player is considered winning in any given state of connect-4. To do this, I created a neural network that takes in a connect-4 board state (6 wide x 7 high), and outputs a single number between $[0, 1]$ that represents the predicted confidence that player 1 is winning. To train the network, I will use a stochastic gradient descent algorithm, as well as a back propagation algorithm, in order to update the weights.

1.2 Mathematic Background

The process of modifying the weights is what drives the neural network to learn. Training it is an iterative process; Each iteration, we give the network a random sample of data to train off of. We then take the average "gradient" and modify the weights accordingly.

To describe it best, lets start from the ground up. First, we are given a neural network. A neural network is a "network" that partially mimics a neuron map. The network is divided into layers. For this project, we will use a neural network with 1 hidden layer. If you are confused about how a neural network actually functions, I suggest watching this video (<https://www.youtube.com/watch?v=aircAruvnKk>) by 3blue1brown on Youtube. He does a great job with providing a clean animation to describe what is actually happening.

So we have a neural network, which basically is just a great big equation with a bunch of variables. Those variables, apart from the input variables, are all of our "weights". We store these weights into matrices. Each matrix represents the flow of values from one layer to the next

Now, we need to create a function we want to reduce, or a "loss" function. What this means is that the loss function should be lower if our network results in an answer closer to the expected answer for that input. In short, the higher the loss value, the worse the network. In order to improve the network, we need to modify the weights in such a way that it reduces the loss function.

This is where derivatives come in. Given a partial derivative of a function with respect to a variable, that derivative indicates the direction of greatest ascent in the function. If we inverse it, however, we get the direction of greatest *descent*. Using this, we can calculate the partial derivative of the loss function with respect to each weight, and use that to modify the weights in such a way that improves the accuracy of the network. Each iteration, we reduce the "scale" at which we modify these variables by. Think of this as playing golf on a huge space. You want to get "within a ballpark" on your first shot, then fine tune further shots until you make it in the hole.

Now, while modifying the values each iteration from a single input/output would work, it makes the network favor earlier passes much more than later passes. To eliviate this, we keep the same network for a set of batch. We calculate the weight gradients for each input/output in the batch, and instead of applying each one to the network, we apply the average of the gradients to the network.

To summarize, we first iterate a number of times. Each iteration, we create a batch of randomly selected inputs and expected outputs. We feed those in, and use the loss function to find the partial derivatives for each of the weights. Next, we average out those derivatives for each weight, multiply it by some step factor that decreases over time, and apply those changes to the network. By repeating this, we start to get a network that can potentially properly evaluate our connect-4 state.

1.3 Dataset

In order to train the network to evaluate a state of a connect-4 game, I used a dataset of connect-4 games, which provided a list of connect-4 game states, and the resulting value of the game.

John Tromp. (1995). Connect-4 <https://archive.ics.uci.edu/ml/datasets/Connect-4>.
Irvine, CA: University of California, School of Information and Computer Science.

The dataset contains 67,557 entries, each one representing a state in a game of Connect-4. As decribed on the database page:

This database contains all legal 8-ply positions in the game of connect-4 in which neither player has won yet, and in which the next move is not forced.
x is the first player; o the second.
The outcome class is the game theoretical value for the first player.

The data is stored in a .csv file. Each row of the dataset represents a legal position. Each row is composed of 43 comma separated variables. The first 42 variables represent the board state, and the last one represents the expected value of whether player 1 wins, loses or draws.

For the board position, there are three possible states. 'x' means that the first player has placed a piece in that square. 'o' indicates that the second player has placed their piece in that square. 'b' indicates that the square is vacant.

There are three possible "results" of the board state. "win" means that player 1 is expected to win from this game state. "lose" means that player 2 is expected to win from this game state. "draw" means that neither player is expected to win, or that the board will be completely filled before any player can get 4 in a row.

In using the dataset, we need to convert the letters to numbers. First, for the board position, we will make 'x' = 1, 'o' = -1, and 'b' = 0. Secondly, because we want our network to predict whether or not player 1 will win, we will make 'lose' = 'draw' = 0 and 'win' = 1

2 Design

During this section, a variety of variables and variable syntax will be used. The following is a list that explains each variable. Remember that for a given "layer", the output is considered the 0th layer, the hidden layer is considered the 1st layer, and the input is considered the 2nd layer.

y^n The output vector of the nodes in the n^{th} hidden layer.

$y_{i,0}^n$ The output value of node i in the n^{th} hidden layer.

y By default, if y is left alone, it represents the output value, known as $y_{0,0}^0$

\hat{y} The expected final output as a scalar value

W^n The weight matrix for a given layer n , as it applies to the values Y^{n+1} (the outputs of the previous layer)

$w_{i,j}^n$ The weight value for the value passed from node j of layer $n + 1$ as it is passed to node i of layer n

η The step size of the gradient descent as it nudges the weights.

I've sectioned off the design of the algorithms into the independent algorithms, ordering them to try and explain the progression of deriving the algorithms. This will take some experimentation, but the first idea is to

2.1 Sigmoid Activation Function

In neural networks, an activation function is typically used to condense the output of a node down to a [0-1] scale. In this lab, I used a sigmoid function. Below is the sigmoid function and its derivative.

$$\sigma(n) = \frac{1}{1 + e^{-n}}$$
$$\sigma'(n) = \sigma(n) \cdot (1 - \sigma(n))$$

I used the following algorithm so that if I passed in a matrix, it would just apply the sigmoid function to each value in the matrix.

Algorithm 1: $\sigma(n)$ function for both constants and matrices

Input: A - which can either be a constant, or an array, with the first dimension being length n

Output: $\sigma(A)$ - The result of putting A through the $\sigma()$ function

```
1 if  $A$  is a Matrix (or array) then
2    $B \leftarrow$  new Matrix of 0s in the same shape as  $A$ 
3   for  $i = 0, 1, 2 \dots n$  do
4     // This works recursively until it reaches scalar values
5      $B[i] \leftarrow \sigma(A[i])$ 
6   end
7   return  $B$ 
8 else
9   return  $1/(1 + e^{-A})$ 
10 end
```

2.2 Feed-Forward Algorithm

The Feed-Forward algorithm is the algorithm used to "run" a neural network. This is one of the fundamental steps in how a neural network works. This step is used whenever you want to find the output that a neural network will calculate from a given set of inputs. Since the neural network is split into 3 layers, the transition of the values from one layer to the next can be described as:

$$\vec{y}^n = \sigma(W_n \cdot \vec{y}^{n+1})$$

If we substitute in all of our layers, we can get the following equation as the result of the entire network.

$$y = y_{0,0}^0 = \sigma(W_0 \cdot \sigma(W_1 \cdot \vec{y}^2))_{0,0}$$

2.3 Loss Function

In order to use gradient descent, we will need to create a function that measures the correctness of our network. This is known as the "Loss Function". The premise of the loss function is that the higher it is,

the worse the network is. Therefore, in order to make our network more accurate, we need to minimize this value.

$$L = (\hat{y} - y_{0,0}^0)^2 = (\hat{y} - y)^2$$

2.4 0th Layer Weight Derivative

In back propagation, we need to find the derivative of each of the weights in terms of the loss function. First, we need to find the derivative of the loss function with respect to the weights that connect the hidden layer with the output layer.

Given the following equations:

$$y_{i,0}^1 = \sigma\left(\sum_h w_{i,h}^1 \cdot y_{h,0}^2\right)$$

$$y_{0,0}^0 = \sigma\left(\sum_h w_{0,h}^0 \cdot y_{h,0}^1\right)$$

We can calculate the derivative as follows:

$$\begin{aligned} \frac{\delta L}{\delta w_{0,i}^0} &= \frac{\delta L}{\delta y_{0,0}^0} \cdot \frac{\delta y_{0,0}^0}{\delta \sum_h w_{0,h}^0 \cdot y_{h,0}^1} \cdot \frac{\delta \sum_h w_{0,h}^0 \cdot y_{h,0}^1}{\delta w_{0,i}^0} \\ \frac{\delta L}{\delta w_{0,i}^0} &= 2 \cdot (\hat{y} - y_{0,0}^0) \cdot \sigma'\left(\sum_h w_{0,h}^0 \cdot y_{h,0}^1\right) \cdot y_{i,0}^1 \\ \frac{\delta L}{\delta w_{0,i}^0} &= 2 \cdot (\hat{y} - y_{0,0}^0) \cdot \sigma\left(\sum_h w_{0,h}^0 \cdot y_{h,0}^1\right) \cdot (1 - \sigma\left(\sum_h w_{0,h}^0 \cdot y_{h,0}^1\right)) \cdot y_{i,0}^1 \\ \frac{\delta L}{\delta w_{0,i}^0} &= 2 \cdot (\hat{y} - y_{0,0}^0) \cdot y_{0,0}^0 \cdot (1 - y_{0,0}^0) \cdot y_{i,0}^1 \end{aligned}$$

2.5 1st Layer Weight Derivative

Next, we need to find the derivative of each of the weights that connect the input layer to the hidden layer. Given the following equations:

$$y_{i,0}^1 = \sigma\left(\sum_h w_{i,h}^1 \cdot y_{h,0}^2\right)$$

$$y_{0,0}^0 = \sigma\left(\sum_h w_{0,h}^0 \cdot y_{h,0}^1\right)$$

We can calculate the derivative as the following:

$$\begin{aligned} \frac{\delta L}{\delta w_{i,j}^1} &= \frac{\delta L}{\delta y_{0,0}^0} \cdot \frac{\delta y_{0,0}^0}{\delta \sum_h w_{0,h}^0 \cdot y_{h,0}^1} \cdot \frac{\delta \sum_h w_{0,h}^0 \cdot y_{h,0}^1}{\delta y_{i,0}^1} \cdot \frac{\delta y_{i,0}^1}{\delta \sum_h w_{i,h}^1 \cdot y_{h,0}^2} \cdot \frac{\delta \sum_h w_{i,h}^1 \cdot y_{h,0}^2}{\delta w_{i,j}^1} \\ \frac{\delta L}{\delta w_{i,j}^1} &= 2 \cdot (\hat{y} - y_{0,0}^0) \cdot \sigma'\left(\sum_h w_{0,h}^0 \cdot y_{h,0}^1\right) \cdot w_{0,i}^0 \cdot \sigma'\left(\sum_h w_{i,h}^1 \cdot y_{h,0}^2\right) \cdot y_{j,0}^2 \\ \frac{\delta L}{\delta w_{i,j}^1} &= 2 \cdot (\hat{y} - y_{0,0}^0) \cdot \sigma\left(\sum_h w_{0,h}^0 \cdot y_{h,0}^1\right) \cdot (1 - \sigma\left(\sum_h w_{0,h}^0 \cdot y_{h,0}^1\right)) \cdot w_{0,i}^0 \cdot \sigma\left(\sum_h w_{i,h}^1 \cdot y_{h,0}^2\right) \cdot (1 - \sigma\left(\sum_h w_{i,h}^1 \cdot y_{h,0}^2\right)) \cdot y_{j,0}^2 \\ \frac{\delta L}{\delta w_{i,j}^1} &= 2 \cdot (\hat{y} - y_{0,0}^0) \cdot y_{0,0}^0 \cdot (1 - y_{0,0}^0) \cdot w_{0,i}^0 \cdot y_{i,0}^1 \cdot (1 - y_{i,0}^1) \cdot y_{j,0}^2 \end{aligned}$$

We can then reorder it as..

$$\frac{\delta L}{\delta w_{i,j}^1} = 2 \cdot (\hat{y} - y_{0,0}^0) \cdot y_{0,0}^0 \cdot (1 - y_{0,0}^0) \cdot y_{i,0}^1 \cdot w_{0,i}^0 \cdot (1 - y_{i,0}^1) \cdot y_{j,0}^2$$

And substitute in the 0th layer derivative!

$$\frac{\delta L}{\delta w_{i,j}^1} = \frac{\delta L}{\delta w_{0,i}^0} \cdot w_{0,i}^0 \cdot (1 - y_{i,0}^1) \cdot y_{j,0}^2$$

2.6 Backwards Propagation

Backwards propagation is calculating the "gradient" of weights for a given input. I say "gradient" because in reality, the outputs of this algorithm will be two matrices, which each contain the derivative for each variable. We will be able to use these nudges in order to reduce the overall loss function. I also included an error calculations, such that we can also keep track of how accurate the network is at each batch of inputs.

Algorithm 2: Back Propagation for a 2-layer neural network

Input: Y^2 - The input vector (matrix), dimensions $m \times 1$.
 W^1 - The weight matrix, dimensions $h \times m$, that represents the weights used as values go from $Y^2 \rightarrow Y^1$
 W^0 - The weight matrix, dimensions $1 \times h$, that represents the weights used as values go from $Y^1 \rightarrow Y^0$
 \hat{y} - The expected result from Forward-Feeding

Output: δW^1 - A matrix of the derivatives of the weights in W^1 , as a $h \times m$ matrix.
 δW^0 - A matrix of the derivatives of the weights in W^0 , as a $1 \times h$ matrix.
Err - The absolute error of the network for this given input

```

1  $\delta W^1 \leftarrow$  new matrix of dimensions  $n \times m$ 
2  $\delta W^0 \leftarrow$  new matrix of dimensions  $1 \times h$ 
3  $Y^1 \leftarrow \sigma(W^1 \cdot Y^2)$ 
4  $Y^0 \leftarrow \sigma(W^0 \cdot Y^1)$ 
5 Err  $\leftarrow |\hat{y} - y_{0,0}^0|$ 
6 for  $i \leftarrow (0, 1, \dots, h - 1)$  do
7    $\delta W_{0,i}^0 \leftarrow 2 \cdot (\hat{y} - y_{0,0}^0) \cdot Y_{0,0}^0 \cdot (1 - Y_{0,0}^0) \cdot Y_{i,0}^1$ 
8   for  $j \leftarrow (0, 1, \dots, m - 1)$  do
9      $\delta W_{i,j}^1 \leftarrow \delta W_{0,i}^0 \cdot W_{0,i}^0 \cdot (1 - y_{i,0}^1) \cdot y_{j,0}^2$ 
10  end
11 end
12 return  $\delta W^1, \delta W^0, Err$ 

```

2.7 Stochastic Gradient Descent

For a good portion while working on this project, this was the part I was the most confused on. Stochastic steepest descent, at least my interpretation of it, is taking the accumulation of all of the weight changes taken from the back propagation algorithm, and averaging them before applying them to the network. The general formula is described as follows:

$$W^i = W^i - \frac{\eta}{n} \cdot \sum_i^n \text{Back-Propagation}(\dots)_{W^i}$$

To update each of the weights, we take the average changes in the back propagation, and multiply it by a scalar coefficient η , which represents the step size of the learning. The smaller η is, the less the network

will learn. For a single cycle in gradient descent, we will use the following algorithm to update the weights accordingly:

Algorithm 3: Stochastic Gradient Descent Iteration

Input: input - Length- n list of input vectors, as matrices. Each input vector is of size $m \times 1$
output - Length- n list of all expected values, as scalars. Each output corresponds to the input at that same index
 W^1 - The weight matrix, dimensions $h \times m$, that represents the weights used as values go from $Y^2 \rightarrow Y^1$
 W^0 - The weight matrix, dimensions $1 \times h$, that represents the weights used as values go from $Y^1 \rightarrow Y^0$
 η - The step coefficient to nudge the weights by.

Output: W^1 - The updated weight matrix, dimensions $h \times m$, where values have been modified by their average derivative across all inputs/outputs, scaled by η
 W^0 - The updated weight matrix, dimensions $1 \times h$, where values have been modified by their average derivative across all inputs/outputs, scaled by η
Err - The average absolute error of the network on each input.

```

1 Err Total  $\leftarrow$  0
2  $\Delta W^1 \leftarrow$  matrix of zeros with dimensions  $h \times m$ 
3  $\Delta W^0 \leftarrow$  matrix of zeros with dimensions  $1 \times h$ 
4 for  $i \leftarrow 0, 1, \dots, n - 1$  do
5    $\delta W^1, \delta W^0, \delta \text{Err} \leftarrow$  Back Propagation on input[ $i$ ] and output[ $i$ ]
6    $\Delta W^1 \leftarrow \Delta W^1 + \delta W^1$ 
7    $\Delta W^0 \leftarrow \Delta W^0 + \delta W^0$ 
8   Err Total  $\leftarrow$  Err Total + Err
9 end
10 Err Total  $\leftarrow$  Err Total/ $n$ 
11  $W^1 \leftarrow W^1 - \eta/n \cdot \Delta W^1$ 
12  $W^0 \leftarrow W^0 - \eta/n \cdot \Delta W^0$ 
13 return  $W^1, W^0, \text{Err Total}$ 

```

2.8 Training the Network

Now, in order to train the network, I need to send in batches of data for it to train on, and each batch of data I need to reduce the step function so the network can fine-tune its training. I plan to do this by using a random sample approach. I'll have a set of modifyable parameters, namely the batch size and the iteration count. These will allow me to change how many data points are put in to the network for each iteration, and how many iterations I will have. For step size, I will pass in a function, such that I can experiment with different step size functions. During implementation, I may also allow the ability to allow for keyboard interruption instead, so I can leave it running for however long I want, and come back to a trained network. The algorithm to perform this training is described as follows:

Algorithm 4: Neural Network Training Algorithm

Input: input - Length- n list of input vectors, as matrices. Each input vector is of size $m \times 1$
output - Length- n list of all expected values, as scalars. Each output corresponds to the input at that same index
 W^1 - The weight matrix, dimensions $h \times m$, that represents the weights used as values go from $Y^2 \rightarrow Y^0$
 W^0 - The weight matrix, dimensions $h \times m$, that represents the weights used as values go from $Y^2 \rightarrow Y^0$
 $f(i, I)$ - Step size function, which when passed in a current iteration i , and total iteration count I , returns the step size for that iteration
 k - The total number of iterations to run
 c - The total number of inputs to feed into each iteration
 $seed$ - the random seed to use in the random selection of inputs

Output: W^1 - The updated weight matrix, dimensions $h \times m$, trained on the given inputs and outputs
 W^0 - The updated weight matrix, dimensions $1 \times h$, trained on the given inputs and outputs

```
// Set the random seed to seed here
1 for  $i \leftarrow 0, 1, \dots, k - 1$  do
2    $\eta \leftarrow f(i, I)$ 
3   indexes  $\leftarrow$  list of uniform-random indexes within the range  $[0, n - 1]$  using the seed  $seed$ 
4   input-iteration  $\leftarrow$  list of inputs at each corresponding index from indexes
5   output-iteration  $\leftarrow$  list of outputs at each corresponding index from indexes
6    $W^1, W^0, Err \leftarrow$  Stochastic Gradient Descent of weights and step size for input-iteration and
   output-iteration
7   Log the  $Err$  of the training
8 end
9 return  $W^1, W^0$ 
```

3 Implementation

Implementing the algorithms was mostly straight-forward. The only caveat came with data importing, which is described below.

3.1 Data Importing

Since the data was given in a .csv file type, with strings instead of numbers, another python script was created in order to convert the data into something the network could use. The file `data_import.py` contains the code responsible for importing the data for the data and converting the string values into numbers for the network to perform.

4 Validation

5 Reflection