

# PART I

## Neural Networks

**DATASETS USED:** Iris Flower Dataset, MNIST Dataset

### 1 Basics of Neural Networks

An Artificial Neural Network is a machine learning model loosely based on the structure of interconnected neurons in the human brain. Since the biological brain is well-adapted to handling complex tasks, the inspiration for Artificial Neural Networks is that a model that mimics the brain's structure and basic functioning would prove most capable at solving difficult machine learning problems.

Structurally, a Neural Network is a collection of interconnected units called Perceptrons. The perceptrons are arranged in different layers and the data flows from one layer to the next. This flow of data is controlled by weights that connect adjacent layers, as shown in Figure 1.

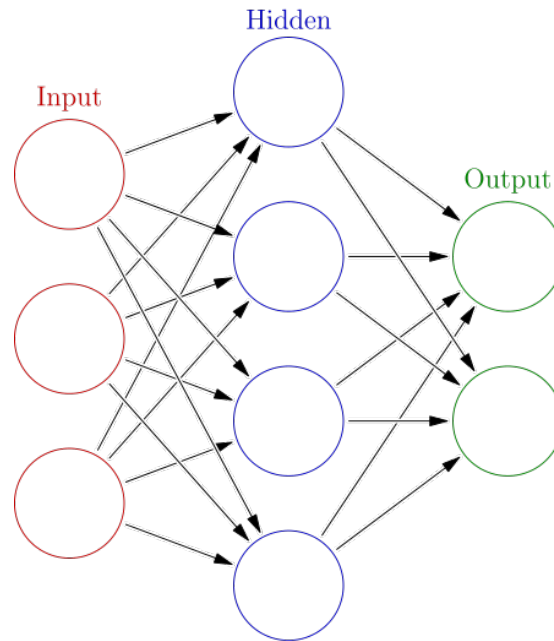


Figure 1: Basic Structure of ANN

In the most common type of neural network, all of the perceptrons in a

given layer receive inputs from all the perceptrons in the previous layer. After receiving these inputs, the perceptron computes its value as follows:

If the inputs to the network are  $X$  and the weights of the perceptron are  $W$  then:

$$X = [x_0, x_1, x_2, \dots] \quad (1)$$

$$W = [w_0, w_1, w_2, \dots] \quad (2)$$

The weighted sum of inputs is calculated and a bias is added:

$$\begin{aligned} Z &= X \cdot W + b \\ \implies Z &= b + x_0 w_0 + x_1 w_1 + \dots \end{aligned} \quad (3)$$

An activation function is used to introduce non-linearity:

$$H = \text{sigmoid}(Z) = \frac{1}{1 + e^{-Z}} \quad (4)$$

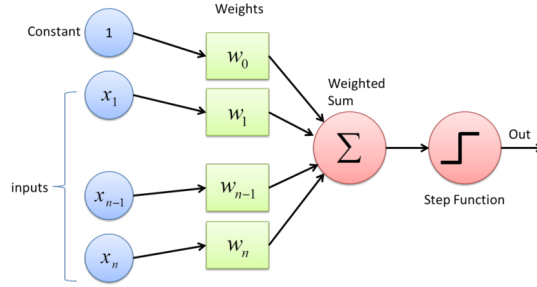


Figure 2: Calculations in a Perceptron

## 2 Functioning of Neural Networks

### 2.1 Vectorization

- In a neural network, the collection of weights of every unit in a layer are vectorized into a matrix  $\theta$ , in which each row holds all the weights of a single perceptron.
- The values of all the units in a given layer are handled together as a vector
- The biases in a given layer are also handled together as a vector.
- The inputs and outputs are given and received in the form of a vector

## 2.2 Notations Used

$L :$	<i>The Number of Layers. (<math>l = 1, 2, \dots, L</math>)</i>
$z_j^i :$	<i>The weighted sum of <math>j^{\text{th}}</math> unit in <math>i^{\text{th}}</math> layer.</i>
$a_j^{(i)} :$	<i>the activated value of <math>j^{\text{th}}</math> unit in <math>i^{\text{th}}</math> layer</i>
$(x^{(m)}, y^{(m)}) :$	<i>input – output pair for the <math>m^{\text{th}}</math> example.</i>
$\theta^{(l)} :$	<i>the weights connecting layer <math>l</math> to the next.</i>

## 2.3 Forward Propagation

Forward propagation uses the parameters of the network to calculate an output vector from an input vector, by stepping through the computations layer by layer as follows:

For a given input vector  $x$ , to calculate the prediction vector  $\hat{y}$ :

Start with:

$$a^{(1)} = x$$

Repeat for  $l = 2$  to  $L$ :

$$a^{(l-1)} = [1; a^{(l-1)}] \quad \text{prepend 1 for bias.}$$

$$z^{(l)} = \theta^{(l-1)} \cdot a^{(l-1)} + b \quad \text{compute the weighted sum } z.$$

$$a^{(l)} = \text{sigmoid}(z^{(l)}) \quad \text{apply the activation function.}$$

End with:

$$\hat{y} = a^{(L)}$$

## 2.4 Cost Function and Optimization

Although many cost functions are possible, the squared error cost function is most common.

$$J = \sum_{\forall j} (\hat{y} - y)^2 \quad (5)$$

For a given dataset with  $m$  examples, the cost is averaged over the dataset.

Training the neural network is the process of reducing this cost function by adjusting the parameters  $\theta$ . We use gradient descent for this purpose.

Therefore, the gradient of the cost function is computed w.r.t every parameter  $\theta_{ij}^{(l)}$  and we use the update step:

$$\theta_{ij}^{(l)} = \theta_{ij}^{(l)} - \alpha \frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta) \quad (6)$$

Gradient Descent is applied on a neural network as per the backpropagation algorithm.

## 2.5 The Backpropagation Algorithm

Training Set  $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$

Set  $\Delta_{ij}^{(l)} = 0$  (for all  $l, i, j$ ).

For  $i = 1$  to  $m$ :

Set  $a^{(1)} = x^{(i)}$

Perform forward propagation to compute  $a^{(l)}$  for  $l = 2, 3, \dots, L$

Using  $y^{(i)}$ , compute  $\delta^{(L)} = a^{(L)} - y^{(i)}$

Backpropagate  $\delta s$ : Compute  $\delta^{(l)} = (\theta^{(l)})^T \delta^{(l+1)}$  for  $l=L-1, \dots, 2$

$\Delta_{ij}^{(l)} := \Delta_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)}$

$$\frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta) := \frac{1}{m} \Delta_{ij}^{(l)}$$

## 2.6 One-Hot Coding

To build a Neural Network for multi-class classification, the labels of the data needs to be encoded into a format that the network can output. Most commonly, the data are one-hot coded.

This means that, for an  $n$ -class problem, each label is a  $n$ -dimensional vector, with a 1 in the place of the correct output class and 0's everywhere else.

For the Iris Data Set, this coding looks like:

<i>Iris setosa</i>	[1 0 0]
<i>Iris virginica</i>	[0 1 0]
<i>Iris versicolor</i>	[0 0 1]

## 2.7 Stochastic Gradient Descent

Stochastic Gradient Descent is a variation of Gradient Descent where each update step is performed with a single training example. An example is randomly chosen from the dataset, the gradient of the cost function w.r.t this example is computed and the update is also done.

The purpose of S.G.D was to speed up the computation for large datasets, where computing the gradient with all examples before each update step can be very time expensive.

A drawback with this approach is that, since the algorithm sees different data for each update, the optimization never appears to converge to the minimum. Instead, it stays in random motion around an area close to the minimum.

Another variation of Gradient Descent is called mini-batch gradient descent, where a random mini-batch is chosen for computing the gradient. It can be thought of as a middle-ground between batch gradient descent and S.G.D.

## 3 Experiment 1

### 3.1 Objective

To build a Neural Network that recognizes handwritten digits by training on the MNIST database using Mini-Batch Gradient Descent.

Use multiple runs successively and tune the hyperparameters manually based on the results of the previous run.

### 3.2 Implementation

- \* A class *ANN* was created that handles Neural Networks of any shape. It wrapped up the functions of parameter initialization, Forward Propagation and Backpropagation.
- \* Created a function *Cost()* that returns the Cost of the network on a given dataset.
- \* Created a function *fit()* that performs mini-batch gradient descent on a given network. Allows the dataset to be split into *pieces* number of pieces (determines size of a batch) and also sets the given learning rate and no. of epochs.
- \* Created functions *acc()* and *ConfMat()* that return the accuracy value and the full confusion matrix based on test data.

### 3.3 Procedure

- \* Initialized a Neural Network *m* (using *ANN* class) with 2 hidden layers of sized *h1*=183 and *h2*=42. The sizes were determined by inserting 2 geometric means b/w input size(784) and output size(10)
- \* Each example in the MNIST dataset was flattened into a 1-dimensional vector and passed into the network.
- \* **Run1:** Using *fit()*, the network was trained on the dataset for 10 epochs, with batch size chosen as 1000 and learning rate as 0.003
- \* The graph of Cost for Run1 suggested that the learning rate was too high and/or the batch size was too low.
- \* **Run2:** Using *fit()*, the network was trained on the dataset for 10 epochs, with batch size chosen as 2000 and learning rate as 0.001
- \* For Run2, once again, the graph of Cost vs Iterations suggested that the learning rate was too high and/or the batch size was too low.
- \* **Run3:** Using *fit()*, the network was trained on the dataset for 10 epochs, with batch size chosen as 4000 and learning rate as 0.0003
- \* After these 3 Runs, the trained network was tested on the 10,000 test examples and a Confusion Matrix was tabulated.

### 3.4 Results and Analysis

#### 3.4.1 Run 1

batch\_size = 1000

$\alpha = 0.003$

epochs = 10

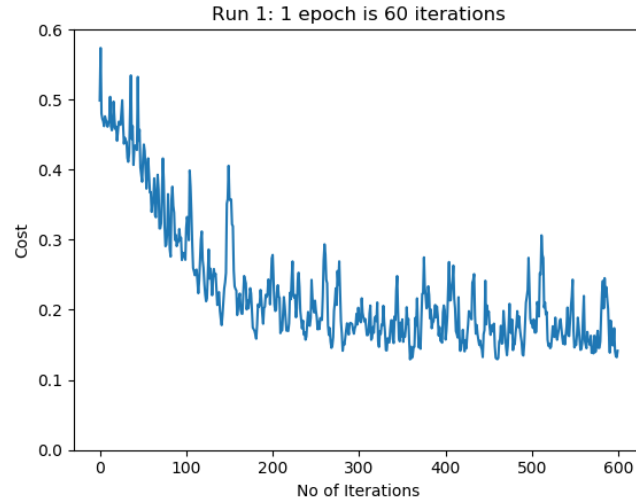


Figure 3: Fall of Cost for Run 1

- \* The fall of cost plateaued out in the first half of the training itself. There was no meaningful improvement for nearly the latter half of the iterations.
- \* The curve was very rough and the cost increased and decreased by a significant amount intermittently.
- \* This trend of roughness in the curve could be due to a step size too large. Hence,  $\alpha$  was reduced in the next run.
- \* This trend might also have arisen because the mini-batch was too small. Hence, batch\_size was increased in the next run.

Initial Cost = 0.498

Final Cost = 0.141

#### 3.4.2 Run 2

batch\_size = 2000

$\alpha = 0.001$

epochs = 10

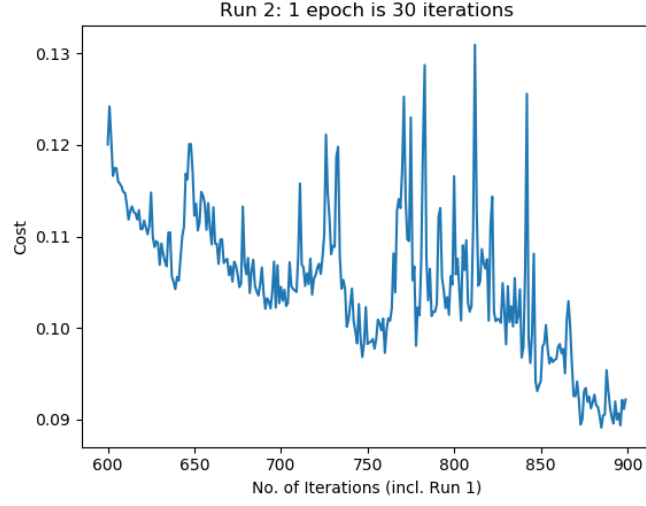


Figure 4: Fall of Cost for Run 2

- \* The overall fall of cost was smaller compared to Run 1.
- \* The trend in the graph was for the cost to increase and decrease drastically in intermittent steps. The cost even exceeded the starting cost for a few iterations.
- \* The reduction of  $\alpha$  and increase of batch\_size from Run 1 proved effective in the beginning. But by the end of Run 2, the cost had stopped decreasing to the desired degree.
- \* Hence,  $\alpha$  was further decreased and batch\_size was further increased for the next run.

Initial Cost = 0.120

Final Cost = 0.0921

### 3.4.3 Run 3

batch\_size = 4000

$\alpha = 0.0003$

epochs = 10

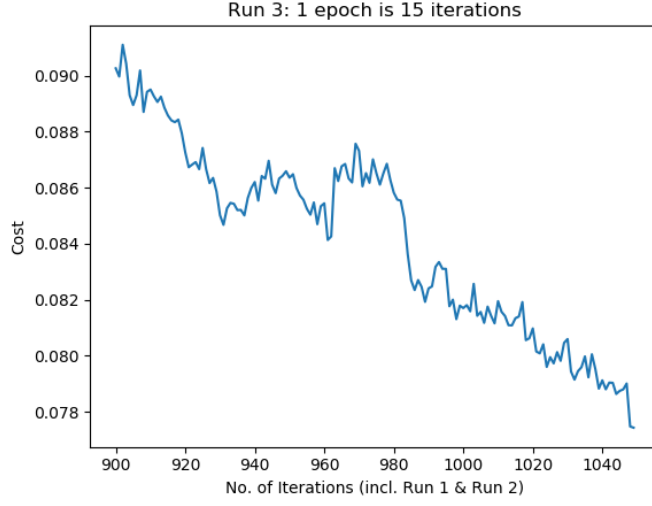


Figure 5: Fall of Cost for Run 3

Initial Cost = 0.0902

Final Cost = 0.0774

A confusion matrix was plotted with the model trained after 3 Runs:

<i>Predicted Class</i> → <i>Actual Class</i>	'0'	'1'	'2'	'3'	'4'	'5'	'6'	'7'	'8'	'9'
'0'	949	0	2	3	1	9	8	3	4	1
'1'	0	1107	3	2	1	1	4	3	14	0
'2'	12	5	903	34	17	2	14	19	21	5
'3'	5	2	30	879	1	24	2	15	45	7
'4'	0	0	3	1	856	4	19	1	9	89
'5'	23	2	3	54	13	721	13	8	45	10
'6'	22	3	7	1	12	14	889	1	6	3
'7'	6	20	22	4	9	0	1	936	4	26
'8'	8	2	9	17	11	38	9	6	849	25
'9'	6	5	5	9	36	9	5	15	11	908

Table 1: Summary of the model's classifications on test data

### 3.5 Conclusions

- \* If  $\alpha$  is too high, the model may not converge close to the minimum. As such,  $\alpha$  might need to be lowered as training progresses.
- \* The higher the batch size, the closer the model can converge to the minimum. As training progresses, the batch size can be increased to enable a convergence closer to the optimum than possible with lower batch sizes.



## 4 Experiment 2

## 5 Objective

To build a Neural Network that classifies the Iris Data Set - the shape of the Neural Network has to be selected. Find the optimal size of the hidden layer by varying the its size - choose the one with highest classification accuracy.

## 6 Procedure

- \* An ANN class was created that handles Neural Networks of any shape. It wrapped up the functions of parameter initialization, Forward Propagation and Backpropagation.
- \* Created a function *basic\_fit()* that uses ANN class to train on the Iris Data Set, with 1000 iterations of backpropagation.
- \* Created a function *rept()* that, with a fixed shape of the hidden layer, trains 100 networks on the data set(using *basic\_fit()*) and returns the average accuracy for the 100 trained networks. This step ensured that the results for a given shape are consistent across many runs.
- \* A function *linsearch()* was created which varies the size of the hidden layer from 1 to 20 and gathers the metrics for each shape using *rept()*.
- \* Comparing the values of accuracy (on test data), the best Neural Network shape was chosen.

## 7 Results and Analysis

For a 3-layered Neural Network, the following graph shows the variation of accuracy (on test data and training data) vs number of units in the hidden layer:

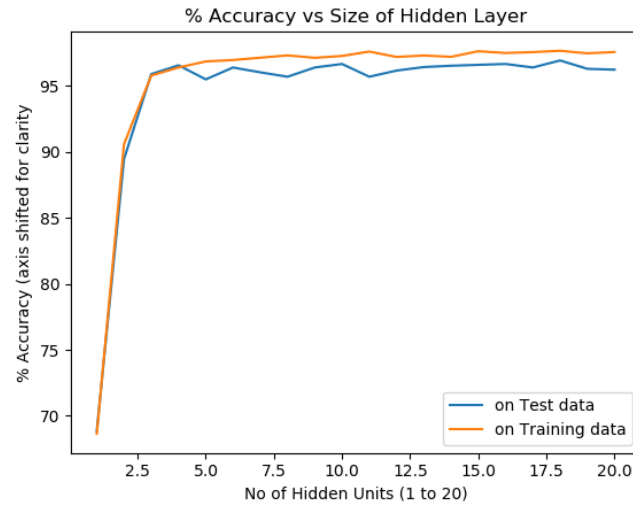


Figure 6: Results of training with various hidden layer sizes

Highest Test Accuracy: 96.93%

Hidden Layer Size: 18

- \* Thus, a Hidden Layer with 18 units was found to be optimal.

## 8 Observations Drawn

- \* If we had chosen the H.L. size based on training accuracy, we would still get 18 units.
  - This indicates that a H.L. size of 18 is ideally suited to the problem since it outperforms other shapes in fitting to the data (evidenced by Training accuracy) and generalizing to test data.
- \* For the smallest hidden layers, the accuracy is relatively very less.
  - This indicates that these networks are underfitting the data.
- \* For higher values of H.L. size, we begin to see testing accuracy fall below training accuracy. (although slightly).
  - This indicates that these Networks are overfitting the data .

## 9 Concluions

- \* For a given problem, there exists some shape of Neural Network that performs the best. It is worthwhile to search and find the right shape.
- \* If the Network is too complex for the problem, it might overfit the data to some degree.
- \* If the network is not complex enough to learn the patterns in the data, accuracy will be very low and the network underfits.