# Brief Summary

In this paper, I implemented a feed forward neural network using only numpy. Different activation functions have their own classes, where the forward and backward method implements the feedforward and backpropagation step. Similarly as an output of the final layer, we have separate class for Softmax operation. Another class is designed for maintaining the linear hidden layers.

In the final model class, we can choose the type and structure of the hidden layers, activation functions between the hidden layers and the final output for the final layers. In addition we also need to provide the cost function which will be used to generate loss for backpropagation.

I used MNIST dataset to classify the digit image. This is a classification problem, so we used softmax function at the final layer which will produce a probability for each of the output layer neuron. For primary investigation, we used a network of 2 hidden layers with 128 and 64 neurons. Just by changing the layers attribute of the model class, we can easily modify our network structures. We used Leaky ReLU as activation function as they do not suffer from the vanishing gradient (like sigmoid) or dead neuron (like ReLU) problem. Similarly to avoid vanishing gradient problem, we used He weight initialization. Multi-class cross entropy is used as the loss function. Our model is capable to use the L2 regularization and dropout. Similarly it is designed to incorporate SGD, Nesterov Momentum, RMSProp and Adam optimizer.

Backpropagation starts from the cost function, then it follows backward to the Softmax, hidden layers and activation function. So other than the cost function, every backward method takes the input of the gradient wrt the next layer and outputs the gradient wrt the previous layer.

In primary investigation with only SGD and a learning rate of 0.5, the model completed 10 epoch in 5 min with a validation loss of 0.13.

# Export Data

```
1 from sklearn.datasets import fetch_openml
2 import numpy as np
3 from sklearn.model_selection import train_test_split
4 import time
5 import matplotlib.pyplot as plt
6
```

```
1 x_orig, y_orig = fetch_openml('mnist_784', version=1, return_X_y=True)
2
```

```
1 inputs = np.array((x_orig/255).astype('float32'))
```

```
2 y = np.array(y_orig, dtype = 'int')
3 targets = np.eye(10)[y]
```

## Activation Functions

```python
1 class Sigmoid():
2     def forward(self, x):
3         self.sigmoid = 1. / (1 + np.exp(-x))
4         return self.sigmoid
5
6     def backward(self, gradient):
7         return self.sigmoid * (1 - self.sigmoid) * gradient
8
9 class ReLU():
10     def forward(self, x):
11         self.x = x
12         return (x>0)*1
13
14     def backward(self, gradient):
15         return (self.x > 0) * gradient
16
17 class LeakyReLU():
18     def __init__(self, alpha = 0.01):
19         self.alpha = alpha
20
21     def forward(self, x):
22         self.x = x
23         return (x>0)*x + (x<0)*self.alpha*x
24
25     def backward(self, gradient):
26         return ( (self.x>0) + (self.x<0)*self.alpha ) * gradient
```

Double-click (or enter) to edit

## Output Layer

```python
1 class Softmax():
```

```
2    def forward(self, x):
3        exps = np.exp(x - x.max())
4        self.output = exps / np.sum(exps, axis = 1)[:,None]
5        return self.output
6
7    def backward(self, gradient):
8        return self.output * (gradient -
9                            (gradient * self.output).sum(axis=1)[:,None] )
10
```

### Loss Function

```
1 class CrossEntropy():
2    def forward(self, y_pred, y_true):
3        self.y_pred = y_pred
4        self.y_true = y_true
5        return np.sum(-1. * y_true * np.log(y_pred)) / y_pred.shape[0]
6
7    def backward(self):
8        return -1. * self.y_true / self.y_pred
```

### Linear Layers

```
1 (np.random.rand(5) < 0.5) * np.array([1,2,3,4,5])

   array([1, 2, 0, 4, 0])
```

```
1 class Linear():
2    def __init__(self, n_input, n_output, dropout_p = 0):
3        # He weight initialization as default activation is ReLU
4        self.Ws = np.random.randn(n_input, n_output) * np.sqrt(2/n_input)
5        self.bs = np.zeros(n_output)
6        # parameters for momentum optimizer
7        self.vw = np.zeros((n_input, n_output))

8        self.vb = np.zeros(n_output)
9        # parameters for RMSprop optimizer
10       self.grad_sq_w = np.zeros((n_input, n_output))
11       self.grad_sq_b = np.zeros(n_output)
```

```
11          self.grad_sq_b = np.zeros(n_output)
12          # parameters for adam optimizer
13          self.moment_1_w = np.zeros((n_input, n_output))
14          self.moment_2_w = np.zeros((n_input, n_output))
15          self.moment_1_b = np.zeros(n_output)
16          self.moment_2_b = np.zeros(n_output)
17          self.dropout_p = dropout_p # dropout probability
18
19      def forward(self, x):
20          self.x = x
21          # dropout layer: 0 means it will be dropped away
22          dropout_layer = np.random.rand(len(self.bs)) > self.dropout_p
23          return (np.dot(x, self.Ws) + self.bs) * dropout_layer
24
25      def backward(self, gradient):
26          self.grad_b = gradient.mean(axis=0)
27          self.grad_W = (self.x[:,:,None] @ gradient[:,None,:]).mean(axis=0)
28          return np.dot(gradient, self.Ws.transpose())
```

## Full Model

```
1 class Model():
2     def __init__(self, layers, cost_func, n_epoch=10, lr=0.01,
3                  reg_lambda=0, beta = 0, decay_rate = 0,
4                  beta1 = 0, beta2 = 0, dropout_p = 0):
5         self.layers = layers
6         self.cost = cost_func
7         self.n_epoch = n_epoch
8         self.lr = lr # learning rate
9         self.beta = beta # momentum parameter
10        self.reg_lambda = reg_lambda # regularization parameter
11        self.decay_rate = decay_rate # rmsprop decay parameter
12        # Adam moment parameter
13        self.beta1 = beta1
14        self.beta2 = beta2
15        self.dropout_p = dropout_p # dropout probability
16
17     def forward(self, x):
18         for layer in self.layers:
19             if type(layer) == Linear:
20                 layer.dropout_p = self.dropout_p
21                 x = layer.forward(x)
```

```
21              x = layer.forward(x)
22          return x
23
24      def loss(self, input, y_true):
25          y_pred = self.forward(input)
26          return self.cost.forward(y_pred, y_true)
27
28      def backward(self):
29          gradient = self.cost.backward()
30          n_layer = len(self.layers)
31          for i in range(n_layer-1, -1, -1):
32              gradient = self.layers[i].backward(gradient)
33
34      def make_minibatch(self, x_data, y_data, mb_size):
35          minibatch_data = [(x_data[i:i+mb_size], y_data[i:i+mb_size]) \
36                  for i in range(0, len(x_data), mb_size)]
37          return minibatch_data
38
39      def train(self, features, targets, mb_size, test_size = 0.20,
40              optimizer = 'nesterov'):
41          x_train, x_val, y_train, y_val = train_test_split(features, targets, test_size=test_size,
42                                                  random_state=42)
43          minibatch_data = self.make_minibatch(x_train, y_train, mb_size)
44          self.train_loss = []
45          self.test_loss = []
46          for epoch in range(self.n_epoch):
47              current_loss = 0
48              n_batch = len(minibatch_data)
49              for minibatch in minibatch_data:
50                  inputs, labels = minibatch
51                  current_loss += self.loss(inputs, labels)
52                  self.backward()
53                  # paramter updates
54                  for layer in self.layers:
55                      if type(layer) == Linear:
56                          if optimizer == 'nesterov':
57                              self.nesterov_momentum(layer)
58                          elif optimizer == 'rmsprop':
59                              self.rmsprop(layer)
60                          elif optimizer == 'adam':
61                              self.adam(layer, epoch)
62              self.train_loss.append(current_loss/n_batch)
63              self.test_loss.append(self.loss(x_val, y_val))
64              if (epoch+1) % 5 == 0:
```

```
64              IT (epoch+I) % 5 == 0:
65                  print(f'Epoch {epoch+1}/{self.n_epoch}: loss = {self.train_loss[epoch]}')
66
67      def loss_curve(self):
68          plt.plot(range(1, self.n_epoch+1), self.train_loss, label = 'train')
69          plt.plot(range(1, self.n_epoch+1), self.test_loss, label = 'test')
70          plt.legend()
71          plt.xlabel('Epoch')
72          plt.ylabel('Cross Entropy Loss')
73          plt.title('Loss Curve')
74
75      def nesterov_momentum(self, layer):
76          '''
77          for beta = 0, it becomes a normal SGD
78          for reg_alpha = 0, there will be no regularization
79          '''
80          vw_old = layer.vw
81          # velocity update with regularization (weight decay)
82          layer.vw = self.beta * vw_old - self.lr * layer.grad_W \
83                      - self.lr * self.reg_lambda * layer.Ws
84          # position Update
85          layer.Ws += (1+self.beta) * layer.vw - self.beta * vw_old
86
87          vb_old = layer.vb
88          # velocity update
89          layer.vb = self.beta * vb_old - self.lr * layer.grad_b
90          # position Update
91          layer.bs += (1+self.beta) * layer.vb - self.beta * vb_old
92
93      def rmsprop(self, layer):
94          layer.grad_sq_w = self.decay_rate * layer.grad_sq_w \
95                  + (1 - self.decay_rate) * layer.grad_W**2
96          layer.Ws += - self.lr * layer.grad_W / (np.sqrt(layer.grad_sq_w) + 1e-7)
97
98          layer.grad_sq_b = self.decay_rate * layer.grad_sq_b \
99                  + (1 - self.decay_rate) * layer.grad_b**2
100         layer.bs -= self.lr * layer.grad_b / (np.sqrt(layer.grad_sq_b) + 1e-7)
101
102     def adam(self, layer, epoch):
103         # moment update
104         layer.moment_1_w = self.beta1 * layer.moment_1_w \
105                 + (1 - self.beta1) * layer.grad_W
106         layer.moment_2_w = self.beta2 * layer.moment_2_w \
107                 + (1 - self.beta2) * layer.grad_W**2
```

```
107                    + (1 - self.beta2) * layer.grad_w**2
108            # bias correction
109            unbias_moment_1_w = layer.moment_1_w / (1 - self.beta1 ** (epoch+1))
110            unbias_moment_2_w = layer.moment_2_w / (1 - self.beta2 ** (epoch+1))
111            # update parameter
112            layer.Ws -= self.lr * unbias_moment_1_w / (np.sqrt(unbias_moment_2_w) + 1e-7)
113
114            # moment update
115            layer.moment_1_b = self.beta1 * layer.moment_1_b \
116                    + (1 - self.beta1) * layer.grad_b
117            layer.moment_2_b = self.beta2 * layer.moment_2_b \
118                    + (1 - self.beta2) * layer.grad_b**2
119            # bias correction
120            unbias_moment_1_b = layer.moment_1_b / (1 - self.beta1 ** (epoch+1))
121            unbias_moment_2_b = layer.moment_2_b / (1 - self.beta2 ** (epoch+1))
122            # update parameter
123            layer.bs -= self.lr * unbias_moment_1_b / (np.sqrt(unbias_moment_2_b) + 1e-7)
124
```

## Testing Working Model

## SGD

```
1 model = Model([Linear(784,128), LeakyReLU(), Linear(128,64), LeakyReLU(), Linear(64,10), Softmax()], CrossEntropy())
2
3 model.n_epoch = 1
4 model.lr = 0.1
5 model.reg_lambda = 0
6
7 mb_size = 50
8 model.train(inputs, targets, mb_size,
9             optimizer='nesterov')
```

```
1 model.n_epoch = 1
2 model.test_loss[-1], model.train_loss[-1]
```

```
    (0.18076048213141493, 0.3187320765505537)
```

```
1 # we will use test dataset for the unbiased evaluation of the final model
```

```
1 # we will use test dataset for the unbiased evaluation of the final model
2 x_train, x_test, y_train, y_test = train_test_split(inputs, targets, test_size
3                                                             random_state=42)
```
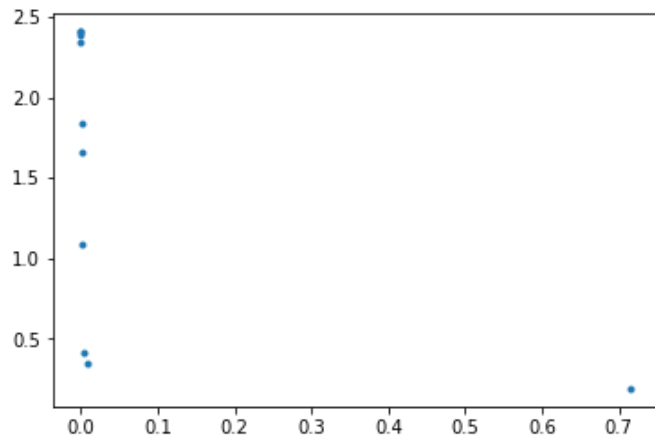
```
1 # random search 1-fold cross_validation
2 n_iterations = 10
3 validation_loss  = []
4 lr_s = []
5 for _ in range(n_iterations):
6     model = Model([Linear(784,128), LeakyReLU(), Linear(128,64), LeakyReLU(),
7                     Linear(64,10), Softmax()], CrossEntropy())
8     model.lr = 10**np.random.uniform(-6,1)
9     model.n_epoch = 2
10    mb_size = 50
11    model.train(x_train, y_train, mb_size,
12                optimizer='nesterov')
13    validation_loss.append(model.test_loss[-1])
14    lr_s.append(model.lr)
15
```

```
1 plt.plot(lr_s, validation_loss,'.')
```

```
[<matplotlib.lines.Line2D at 0x7f464def8510>]
```



```
1 model = Model([Linear(784,128), LeakyReLU(), Linear(128,64), LeakyReLU(),
2                 Linear(64,10), Softmax()], CrossEntropy())
3 model.lr = 0.5
4 model.n_epoch = 10
5 mb_size = 50
```
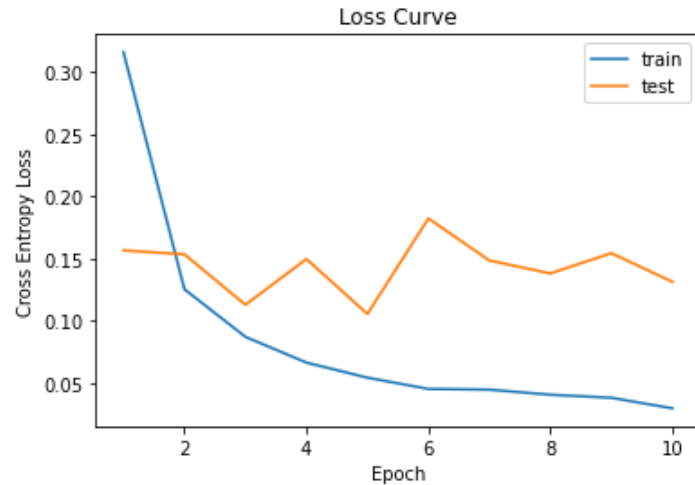
```
5 mb_size = 50
6 model.train(x_train, y_train, mb_size,
7              optimizer='nesterov')
8 model.loss_curve()
```

```
Epoch 5/10: loss = 0.05425715980876315
Epoch 10/10: loss = 0.02961947304835689
```


Loss Curve

```
1 model.train_loss[-1], model.test_loss[-1]
```

```
(0.02961947304835689, 0.13130076481722563)
```

## ▾ SGD + Momentum

```
1 model = Model([Linear(784,128), ReLU(), Linear(128,64), ReLU(), Linear(64,10), Softmax()], CrossEntropy())
2
3 model.n_epoch = 50
4 model.lr = 0.2
5 model.reg_lambda = 0
6 model.beta = 0.9
7
8 mb_size = 50
9 model.train(inputs[:1000], targets[:1000], mb_size,
10             optimizer='nesterov')
```
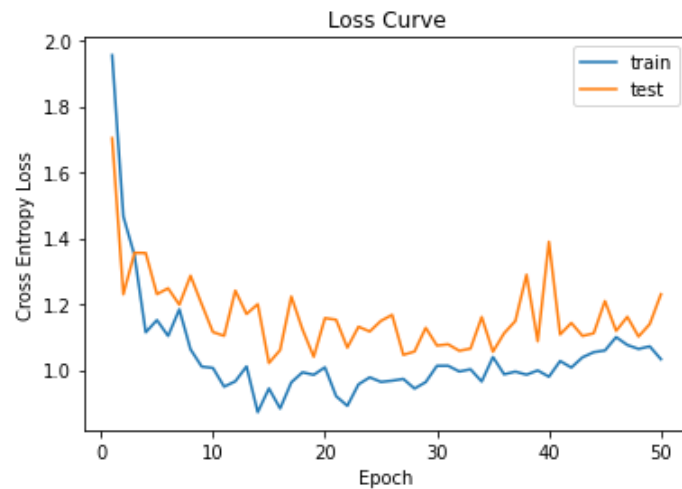
```
Epoch 5/50: loss = 1.1517014898087208
```

```
Epoch 5/50: loss = 1.1317014898687268
Epoch 10/50: loss = 1.00662743040652
Epoch 15/50: loss = 0.9443916198321197
Epoch 20/50: loss = 1.0084874563668824
Epoch 25/50: loss = 0.9643326139300877
Epoch 30/50: loss = 1.0130861220831628
Epoch 35/50: loss = 1.0398071627464067
Epoch 40/50: loss = 0.980003936129946
Epoch 45/50: loss = 1.0596051365987094
Epoch 50/50: loss = 1.0330920507013976
```

```
1 model.loss_curve()
```


Loss Curve

## RMSprop

```
1 model = Model([Linear(784,128), ReLU(), Linear(128,64), ReLU(), Linear(64,10), Softmax()], CrossEntropy())
2
3 model.n_epoch = 50
4 model.lr = 0.1
5
6 model.decay_rate = 0.99
7 mb_size = 50

8 model.train(inputs[:1000], targets[:1000], mb_size,
9             optimizer='rmsprop')
```
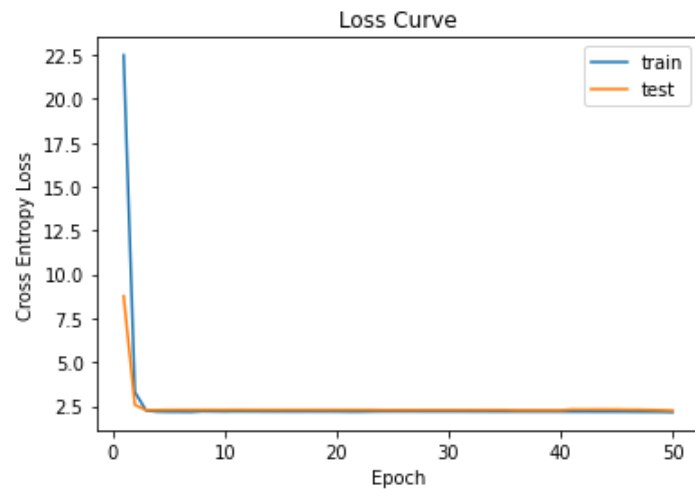
```
Epoch 5/50: loss = 2.182079555828885
```

```
Epoch 5/50:  loss = 2.182079555828885
Epoch 10/50: loss = 2.1988231122133968
Epoch 15/50: loss = 2.1985008803422024
Epoch 20/50: loss = 2.194965930453911
Epoch 25/50: loss = 2.2040840284984777
Epoch 30/50: loss = 2.201275018452942
Epoch 35/50: loss = 2.1953194477842866
Epoch 40/50: loss = 2.1929040052329642
Epoch 45/50: loss = 2.186529289103248
Epoch 50/50: loss = 2.161922242154096
```
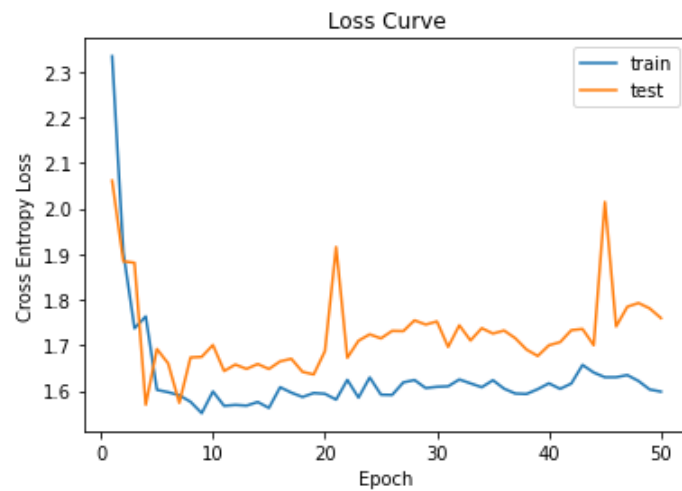
```
1 model.loss_curve()
```



## ▾ Adam

```
 1 model = Model([Linear(784,128), ReLU(), Linear(128,64), ReLU(), Linear(64,10), Softmax()], CrossEntropy())
 2
 3 model.n_epoch = 50
 4 model.lr = 0.1
 5
 6 model.beta1 = 0.9
 7 model.beta2 = 0.999
 8
 9 mb_size = 50
10 model.train(inputs[:1000], targets[:1000], mb_size,
11              optimizer='adam')
```
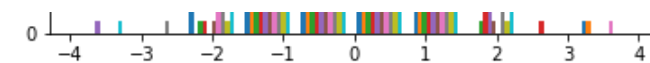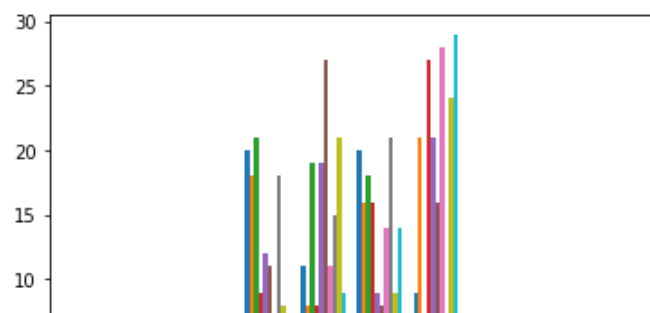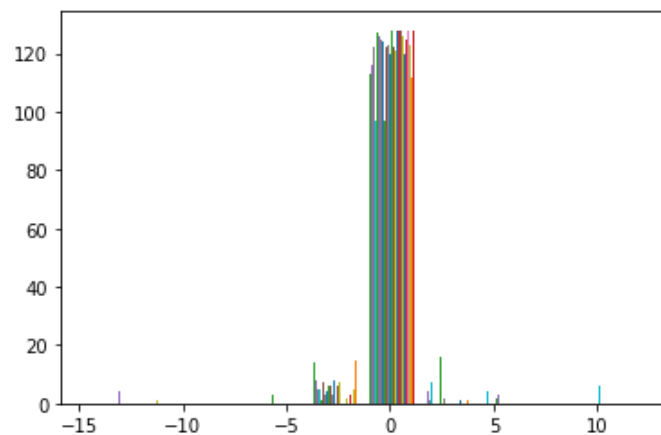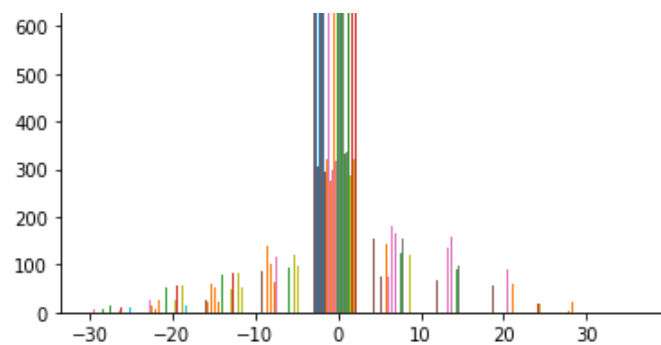
```
Epoch 5/50: loss = 1.6021393644126982
Epoch 10/50: loss = 1.5985715699298186
Epoch 15/50: loss = 1.5624256906820824
Epoch 20/50: loss = 1.593509593272557
Epoch 25/50: loss = 1.5915735457146003
Epoch 30/50: loss = 1.608829267389039
Epoch 35/50: loss = 1.62321705817252
Epoch 40/50: loss = 1.6164474601976628
Epoch 45/50: loss = 1.6298084917409386
Epoch 50/50: loss = 1.5982098766245372
```

```
1 model.loss_curve()
```



```
1 for layer in model.layers:
2     if type(layer) == Linear:
3         plt.hist(layer.Ws)
4         plt.show()
```

1

0s    completed at 3:01 PM