

# ECE 285 – Assignment #2

## Backpropagation

*Written by Sneha Gupta, Shobhit Trehan and Charles Deledalle.*

This assignment focuses on Multiclass Classification on the MNIST Dataset. The MNIST Dataset consists of labeled images of handwritten digits from 0 to 9. Each image is 28 by 28 pixels (784 pixels). Read about the dataset here <http://yann.lecun.com/exdb/mnist/>

We will use a shallow artificial neural network to classify the digits. The network will have an input layer, a **relu** hidden layer and a **softmax** output layer. We will use logistic regression with cross-entropy loss as the objective function in order to optimize the weights and biases. The optimization will be performed by gradient descent with backprop as seen in class. We will rely on a more general matrix form for backprop given by

$$\begin{aligned} \mathbf{W}_k^{t+1} &= \mathbf{W}_k^t - \gamma \nabla_{\mathbf{W}_k} E^t \\ \mathbf{b}_k^{t+1} &= \mathbf{b}_k^t - \gamma \nabla_{\mathbf{b}_k} E^t \\ \text{with } \nabla_{\mathbf{W}_k} E &= \boldsymbol{\delta}_k \mathbf{h}_{k-1}^T \quad \text{where } \mathbf{h}_0 = \mathbf{x} \\ \text{and } \nabla_{\mathbf{b}_k} E &= \boldsymbol{\delta}_k \mathbf{1}_N \\ \text{and } \boldsymbol{\delta}_k &= \left[ \frac{\partial g_k(\mathbf{a}_k)}{\partial \mathbf{a}_k} \right]^T \times \mathbf{e}_k \quad \text{where } \mathbf{e}_k = \begin{cases} \nabla_{\mathbf{y}} E & \text{if } k \text{ is an output layer} \\ \mathbf{W}_{k+1}^T \boldsymbol{\delta}_{k+1} & \text{otherwise} \end{cases} \end{aligned}$$

where the notations are:

- $t$ : iteration index of gradient descent,
- $\gamma$ : learning rate of gradient descent,
- $\mathbf{W}_k$ : matrix of weights at layer  $k$ ,
- $\mathbf{b}_k$ : vector of biases at layer  $k$ ,
- $\mathbf{x}$ : matrix with the  $N$  training input vectors in column,
- $\mathbf{h}_k = g_k(\mathbf{a}_k)$ : matrix with all hidden outputs at layer  $k$  in column,
- $\mathbf{a}_k = \mathbf{W}_k \mathbf{h}_{k-1} + \mathbf{b}_k$ : matrix with all weighted sums at layer  $k$  in column,
- $g_k$ : activation function at layer  $k$ ,
- $\mathbf{1}_N$ : column vector of size  $N$  containing only ones.

In the following, for simplicity, we will drop the indices  $k$  and  $t$ .

## 1 Getting started

First of all, connect to DSMLP (**ieng6** server), start a *pod* and connect to your Jupyter Notebook from your web browser (go back to Lab #1 for more details). Create a new notebook **Lab2.ipynb** and import

```
import numpy as np
from matplotlib import pyplot
```

For the following questions, please write your code and answers directly in your notebook. Organize your notebook with headings, markdown and code cells (following the numbering of the questions).

## 2 Read MNIST Data

The MNIST dataset is already installed on DSMLP in the directory `/datasets/MNIST/`. The data can be loaded using a small dedicated package `/datasets/ee285f-public/MNISTtools.py`. From your terminal, create a symbolic link on this package into your working/current directory

```
$ ln -s /datasets/ee285f-public/MNISTtools.py .
```

Back to your Jupyter Notebook, you can now import its functions as:

```
import MNISTtools
```

This package contains two functions: `load` and `show`. You can display their description using the Python command `help`:

```
help(MNISTtools.load)
help(MNISTtools.show)
```

You can also type `MNISTtools.load` and press `Shift + Tab`.

1. Using `MNISTtools.load`, store the images and labels from the training datasets into two variables, respectively, `xtrain` and `ltrain`. What are the shapes of both variables? What is the size of the training dataset? What is the feature dimension?
2. Display the image of index 42 and check that its content corresponds to its label.
3. What is the range of `xtrain` (minimum and maximum values)? What is the type of `xtrain`?
4. Create a function

```
def normalize_MNIST_images(x):
```

that takes a collection of images (such as `xtrain`) and return a modified version in the range  $[-1, 1]$  of type `float64`. Update `xtrain` accordingly.

*Hint: convert a Numpy array `x` from `int8` to `float64` using `x.astype(np.float64)`.*

5. Using integer array indexing, complete the following function

```
def label2onehot(lbl):
    d = np.zeros((lbl.max() + 1, lbl.size))
    d[COMPLETE, np.arange(0, lbl.size)] = 1
    return d
```

such that `dtrain = label2onehot(ltrain)` will create a Numpy array `dtrain` of one-hot codes stacked in columns (with shape `(10, 60000)`). Make sure that the one-hot code `dtrain[:,42]` corresponds to `ltrain[42]`.

6. Complete the following function

```
def onehot2label(d):
    lbl = d.argmax(axis=COMPLETE)
    return lbl
```

such that `ltrain == onehot2label(dtrain)`.

### 3 Activation functions

For our digit multiclass classification, we will use the **softmax** activation function for the output layer (with 10 units). Given a vector  $\mathbf{a} \in \mathbb{R}^{10}$  (obtained by forward propagation), **softmax** will output a vector  $\mathbf{y} \in \mathbb{R}^{10}$ , where each element  $y_i$  represents the probability that  $\mathbf{x} \in \mathbb{R}^{784}$  is in class  $i$ . The relation between  $\mathbf{a}$  and  $\mathbf{y}$  is given for all  $i \in [1, 10]$  by

$$y_i = g(\mathbf{a})_i = \frac{\exp(a_i)}{\sum_{j=1}^{10} \exp(a_j)}$$

7. When using exponential functions, one should make sure that the input won't be too large or you may observe numerical issues (apparition of **Inf** and subsequently of **NaN**). A simple trick to get rid of this problem is based on the following observation

$$y_i = g(\mathbf{a})_i = \frac{\exp(a_i - M)}{\sum_{j=1}^{10} \exp(a_j - M)} \quad \text{where} \quad M = \max_{j=1 \dots 10} a_j$$

As all inputs of the exponentials will be smaller than 0 and one of them will be exactly 0, you won't encounter numerical issues. Based on this trick, create a function

```
def softmax(a):
```

that returns an array whose columns are the 60,000 predictions  $\mathbf{y}$  from an array whose columns are the 60,000 vectors  $\mathbf{a}$ . *Hint: use Broadcasting and the methods `.max(axis=0)` and `.sum(axis=0)`.*

8. Show that  $\frac{\partial g(\mathbf{a})_i}{\partial a_i} = g(\mathbf{a})_i(1 - g(\mathbf{a})_i)$ .
9. Show that  $\frac{\partial g(\mathbf{a})_i}{\partial a_j} = -g(\mathbf{a})_i g(\mathbf{a})_j$  for  $j \neq i$ .
10. Given a vector  $\mathbf{e} \in \mathbb{R}^{10}$  (obtained during backward propagation), backprop algorithm have to compute

$$\boldsymbol{\delta} = \left[ \frac{\partial g(\mathbf{a})}{\partial \mathbf{a}} \right]^T \times \mathbf{e}$$

where  $\times$  denotes here the matrix vector product. The matrix  $\frac{\partial g(\mathbf{a})}{\partial \mathbf{a}}$  is the Jacobian matrix defined as

$$\frac{\partial g(\mathbf{a})}{\partial \mathbf{a}} = \begin{pmatrix} \frac{\partial g(\mathbf{a})_1}{\partial a_1} & \frac{\partial g(\mathbf{a})_1}{\partial a_2} & \cdots & \frac{\partial g(\mathbf{a})_1}{\partial a_{10}} \\ \frac{\partial g(\mathbf{a})_2}{\partial a_1} & \frac{\partial g(\mathbf{a})_2}{\partial a_2} & \cdots & \frac{\partial g(\mathbf{a})_2}{\partial a_{10}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial g(\mathbf{a})_{10}}{\partial a_1} & \frac{\partial g(\mathbf{a})_{10}}{\partial a_2} & \cdots & \frac{\partial g(\mathbf{a})_{10}}{\partial a_{10}} \end{pmatrix}$$

When the activation is element-wise ( $g(\mathbf{a})_i = g(a_i)$ ), the Jacobian is diagonal and the update of  $\boldsymbol{\delta}$  is simply an element wise product between  $g'(a_i)$  and  $e_i$  (as seen in class). But **softmax** is not an element-wise function and its Jacobian is not diagonal. Nevertheless, it enjoys some interesting properties. From the previous question, deduce that the Jacobian of **softmax** is symmetric and that

$$\boldsymbol{\delta} = g(\mathbf{a}) \otimes \mathbf{e} - \langle g(\mathbf{a}), \mathbf{e} \rangle g(\mathbf{a})$$

where  $\otimes$  is the element-wise product. (Remark that the expression is not that much different from the derivative of the sigmoid logistic function  $g'(a) = g(a)(1 - g(a))$ . This is because softmax is a multivariate generalization of the sigmoid logistic function.)

Based on this formula, write a function

```
def softmax(a, e):
```

that, given an array whose columns are the vectors  $\mathbf{a}$  and an array whose columns are vectors  $\mathbf{e}$ , returns an array whose columns are the vectors  $\delta$ .

*Hint: call `softmax(a)`.*

- Since the Jacobian is symmetric, the update of  $\delta$  corresponds to the directional derivative of  $g$  at point  $\mathbf{a}$  in the direction  $\mathbf{e}$ :

$$\delta = \frac{\partial g(\mathbf{a})}{\partial \mathbf{a}} \times \mathbf{e} = \lim_{\epsilon \rightarrow 0} \frac{g(\mathbf{a} + \epsilon \mathbf{e}) - g(\mathbf{a})}{\epsilon}$$

Based on this formula, you can check your implementation by numerical approximations (finite difference). Complete the following script to check your function `softmax` as follows

```
eps          = 1e-6                # finite difference step
a            = np.random.randn(10, 200) # random inputs
e            = np.random.randn(10, 200) # random directions
diff         = softmax(a, e)
diff_approx  = COMPLETE
rel_error    = np.abs(diff - diff_approx).mean() / np.abs(diff_approx).mean()
print(rel_error, 'should be smaller than 1e-6')
```

- For the hidden layers, we will be using  $\text{ReLU}(\mathbf{a})_i = \max(a_i, 0)$ . Write two functions:

```
def relu(a):
    COMPLETE

def relup(a, e):
    COMPLETE
```

implementing ReLU and its directional derivative. Check your implementation based on numerical approximations.

## 4 Backpropagation

We are now ready to implement our shallow network. The (single) hidden layer between the input and output will consist of  $N_h = 64$  units with the `relu` activation function. The output layer will consist of  $N_o = 10$  units with the `softmax` activation function. The input layer is of dimension  $N_i = 784$ .

- Use the following function to create/initialize your shallow network as follows

```
def init_shallow(Ni, Nh, No):
    b1 = np.random.randn(Nh, 1) / np.sqrt((Ni+1.)/2.)
    W1 = np.random.randn(Nh, Ni) / np.sqrt((Ni+1.)/2.)
    b2 = np.random.randn(No, 1) / np.sqrt((Nh+1.))
    W2 = np.random.randn(No, Nh) / np.sqrt((Nh+1.))
    return W1, b1, W2, b2
```

```

Ni = xtrain.shape[0]
Nh = 64
No = dtrain.shape[0]
netinit = init_shallow(Ni, Nh, No)

```

These types of initializations are called He and Xavier initializations, respectively, and will be explained later during the class.

14. Complete the function `forwardprop_shallow` to evaluate the prediction of our initial network:

```

def forwardprop_shallow(x, net):
    W1 = net[0]
    b1 = net[1]
    W2 = net[2]
    b2 = net[3]

    a1 = W1.dot(x) + b1
    COMPLETE

    return y

yinit = forwardprop_shallow(xtrain, netinit)

```

that produces an array whose columns are vectors  $\mathbf{y}$  corresponding to the predictions obtained for each column  $\mathbf{x}$  given in arguments.

15. Complete the function `eval_loss`:

```

def eval_loss(y, d):
    COMPLETE

print(eval_loss(yinit, dtrain), 'should be around .26')

```

that given your predictions  $\mathbf{y}$  and the desired one-hot codes  $\mathbf{d}$ , computes the average cross-entropy loss (averaged over both the training samples and the vector dimension). Recall that the cross-entropy for  $K$  classes of a vector  $\mathbf{y}$  against a vector  $\mathbf{d}$  is

$$E = - \sum_{i=1}^{10} d_i \log y_i$$

16. Complete the function `eval_perfs`:

```

def eval_perfs(y, lbl):
    COMPLETE

print(eval_perfs(yinit, ltrain))

```

that given your predictions  $\mathbf{y}$  and the desired labels  $\mathbf{lbl}$ , computes the percentage of misclassified samples. Interpret the result.

*Hint: use the function `onehot2label` and do not use loops.*

17. Complete the following function `update_shallow`

```
def update_shallow(x, d, net, gamma=.05):
    W1 = net[0]
    b1 = net[1]
    W2 = net[2]
    b2 = net[3]
    Ni = W1.shape[1]
    Nh = W1.shape[0]
    No = W2.shape[0]

    gamma = gamma / x.shape[1] # normalized by the training dataset size

    COMPLETE

    return W1, b1, W2, b2
```

such that it performs one backpropagation update for your shallow neural network. Recall that the inputs  $\mathbf{x}$  and  $\mathbf{d}$  are going to be arrays with 60,000 columns corresponding to the vectors of images and the vectors of one-hot codes respectively. Show that

$$(\nabla_{\mathbf{y}} E)_i = -\frac{d_i}{y_i}$$

*Hint: Takes inspiration of the code described during the class. Use the functions `softmax`, `softmaxp`, `relu` and `relu_p`.*

18. Using `update_shallow`, complete the function `backprop_shallow`

```
def backprop_shallow(x, d, net, T, gamma=.05):
    lbl = onehot2label(d)
    for t in range(0, T):
        COMPLETE TO UPDATE NET
        COMPLETE TO DISPLAY LOSS AND PERFS
    return net
```

that performs  $T$  updates of the network and **print** the loss and the percentage of training errors at each iteration of backprop. Start testing it using  $T = 2$  iterations

```
nettrain = backprop_shallow(xtrain, dtrain, netinit, 2)
```

When your code starts working, increase the number of iterations to  $T = 5$ , and if it is still working try  $T = 20$ . The loss and training errors should decrease (with some fluctuations).

What percentage of training errors do you reach? (Feel free to increase even more the number of iterations: with  $T = 100$  you should reach about 13% of training errors)

19. Load the testing dataset into two variables `xtest` and `ltest`. What is the size of the testing set? Evaluate the performance of your network on the testing dataset.
20. We will now implement a variant of backpropagation based on stochastic/minibatch gradient descent (we will explain later this variant during the class). The algorithm is very similar to backprop but

instead of updating the weights based on the 60,000 data points at once, we will decompose our training sets on blocks of size  $B = 100$  (called minibatches) and update successively the weights and biases based on the error of each minibatch. The number of updates is then  $TN/B$ . An epoch corresponds to a succession of updates where all minibatches have been processed. The parameter  $T$  is then referred to as the number of epochs. Using `update_shallow`, complete and run for 5 epochs the function `backprop_minibatch_shallow`:

```
def backprop_minibatch_shallow(x, d, net, T, B=100, gamma=.05):
    N = x.shape[1]
    lbl = onehot2label(d)
    for t in range(0, T):
        for l in range(0, (N+B-1)/B):
            idx = np.arange(B*l, min(B*(l+1), N))
            COMPLETE TO UPDATE NET
            y = forwardprop_shallow(x, net)
            COMPLETE TO DISPLAY LOSS AND PERFS
    return net

netminibatch = backprop_minibatch_shallow(xtrain, dtrain, netinit, 5, B=100)
```

*Hint: use integer array indexing.*

21. Compare the performance of this new network on the testing dataset.

## 5 Experiment with network topology and learning parameters

22. Try with  $N_h = 16$  and  $N_h = 256$  number of hidden units. Look at the training and testing errors. Interpret the results.
23. Try different step sizes:  $\gamma = .02$  and  $\gamma = .08$ . Look at the training and testing errors. Interpret the results.
24. Try with minibatches of sizes:  $B = 50$  and  $B = 200$ . Look at the training and testing errors. Interpret the results.
25. Try minibatch gradient descent with more epochs. What is the best testing error that you can achieved?

## 6 Bonus

26. Write generic functions that can deal with an arbitrary number  $L$  of hidden layers.
27. Increase the number of hidden layers, e.g., use two hidden layers instead of one. Create a new architecture that uses two hidden layers of equal size and has approximately the same number of parameters as the previous network with one hidden layer. By that, we mean it should have roughly the same total number of weights and biases. Study the loss, training and testing errors vs. the number of epochs. Repeat with four, eight, sixteen, ... layers.