

---

# Assignment 2: Neural Networks

---

By Ahoo Saeifar

ECE421-Winter 2021

March 1, 2021

# 1. Neural Networks using Numpy

## 1.1 Helper Functions

1. ReLu:

```
def relu(x):  
    return np.maximum(0.0,x)
```

2. softmax:

```
def softmax(x):  
    maxVal=np.amax(x)  
    x=x-maxVal  
    sigma=np.exp(x)/np.sum(np.exp(x),axis=1,keepdims=True)  
    return sigma
```

3. computeLayer:

```
def computeLayer(X, W, b):  
    return np.matmul(X,W)+b
```

4. CE:

```
def CE(target, prediction,dataShape):  
    averageCE=-np.sum(target * np.log(prediction))/dataShape  
    return averageCE
```

5. gradCE:

a. Code Snippet

```
def gradCE(target, prediction):  
    return (prediction - target)/N
```

## b. Derivation of Gradient of Cross Entropy Loss

Gradient of softmax function with respect to its input  $z$ :

$$\begin{aligned} \text{if } i = j : \frac{\partial p_i}{\partial z_i} &= \frac{\partial \frac{e^{z_i}}{\sum e^{z_i}}}{\partial z_i} = \frac{e^{z_i} \sum e^{z_i} - e^{z_i} e^{z_i}}{\sum (e^{z_i})^2} = \frac{e^{z_i}}{\sum e^{z_i}} \frac{\sum e^{z_i} - e^{z_i}}{\sum e^{z_i}} = \frac{e^{z_i}}{\sum e^{z_i}} \left(1 - \frac{e^{z_i}}{\sum e^{z_i}}\right) = p_i(1 - p_i) \\ \text{if } i \neq j : \frac{\partial p_i}{\partial z_j} &= \frac{\partial \frac{e^{z_i}}{\sum e^{z_j}}}{\partial z_j} = \frac{0 - e^{z_i} e^{z_j}}{\sum (e^{z_j})^2} = -\frac{e^{z_i}}{\sum e^{z_i}} \frac{e^{z_j}}{\sum e^{z_j}} = -p_i p_j \end{aligned}$$

Using the results from above we can find the gradient of cross entropy loss:

(sensitivity of  $L$  to layer  $l$  input:  $\delta_i^l$ )

$$\begin{aligned} \frac{\partial L}{\partial z_i} &= - \sum_{k=1}^K \frac{\partial y_k \log(p_k)}{\partial z_i} = - \sum_{k=1}^K y_k \frac{\partial \log(p_k)}{\partial z_i} = - \sum_{k=1}^K y_k \frac{1}{p_k} \frac{\partial p_k}{\partial z_i} \\ &= - \frac{y_i}{p_i} \frac{\partial p_i}{\partial z_i} - \sum_{k \neq i}^K \frac{y_k}{p_k} \frac{\partial p_k}{\partial z_i} = - \frac{y_i}{p_i} p_i(1 - p_i) - \sum_{k \neq i}^K \frac{y_k}{p_k} (-p_k p_i) \\ &= -y_i + y_i p_i + \sum_{k \neq i}^K y_k p_i = -y_i + \sum_{k=1}^K y_k p_i = -y_i + p_i \sum_{k=1}^K y_k \\ &= p_i - y_i \end{aligned}$$

## 1.2 Backpropagation Derivation

We can use the partial derivatives below with chain rule to find the required derivatives:

$$\begin{aligned} \frac{\partial L}{\partial z_i} &= p_i - y_i \text{ (from 1.1.5)} & \frac{\partial z_i}{\partial W_o} &= \frac{\partial W_o h + b_o}{\partial W_o} = h \\ \frac{\partial z}{\partial b_o} &= \frac{\partial W_o x + b_o}{\partial b_o} = 1 & \frac{\partial z}{\partial h} &= W_o & \frac{\partial h}{\partial r} &= \begin{cases} 1 & \text{if } r > 0 \\ 0 & \text{else} \end{cases} \\ \frac{\partial r}{\partial W_h} &= x & \frac{\partial r}{\partial b_h} &= 1 \end{aligned}$$

\*Note:  $h = \text{ReLU}(r)$  and  $p = \text{softmax}(z)$

$$\begin{aligned} 1. \quad \frac{\partial L}{\partial W_o} &= \frac{\partial L}{\partial z_i} \frac{\partial z_i}{\partial W_o} \\ \frac{\partial L}{\partial W_o} &= (p_i - y_i)h \end{aligned}$$

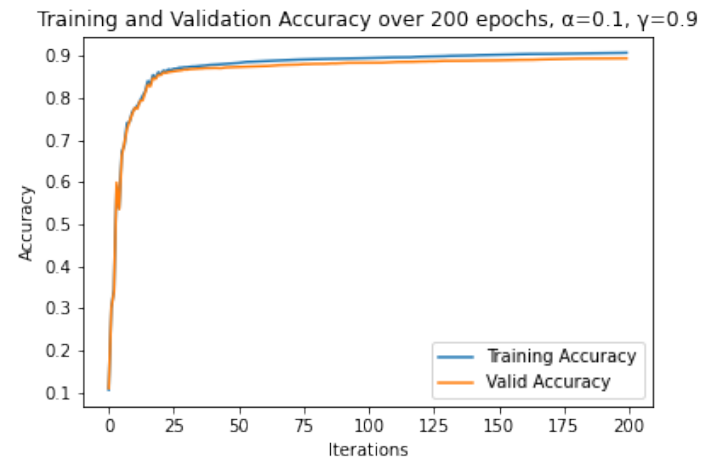
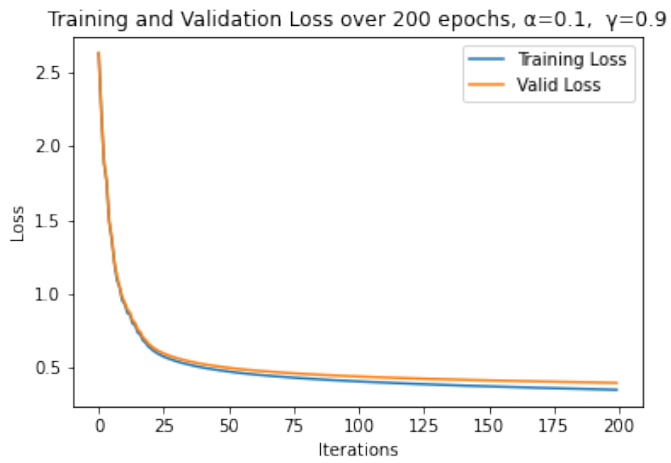
$$\begin{aligned} 2. \quad \frac{\partial L}{\partial b_o} &= \frac{\partial L}{\partial p} \frac{\partial p}{\partial z} \frac{\partial z}{\partial b_o} \\ \frac{\partial L}{\partial b_o} &= (p_i - y_i) \end{aligned}$$

$$\begin{aligned} 3. \quad \frac{\partial L}{\partial W_h} &= \frac{\partial L}{\partial p} \frac{\partial p}{\partial z} \frac{\partial z}{\partial h} \frac{\partial h}{\partial r} \frac{\partial r}{\partial W_h} \\ \frac{\partial L}{\partial W_h} &= (p_i - y_i)(W_o)(h > 0)(x) \end{aligned}$$

$$\begin{aligned} 4. \quad \frac{\partial L}{\partial b_h} &= \frac{\partial L}{\partial p} \frac{\partial p}{\partial z} \frac{\partial z}{\partial h} \frac{\partial h}{\partial r} \frac{\partial r}{\partial b_h} \\ \frac{\partial L}{\partial b_h} &= (p_i - y_i)(W_o)(h > 0) \end{aligned}$$

\*Code snippets are in the Appendix

## 1.3 Learning



# Appendix

## 1.1 Code Snippet for Gradients required in Backpropagation

```
# the gradient of the loss with respect to the output layer weights
def gradLossW_o(target, prediction, hiddenLayer):
    CE=gradCE(target, prediction)
    return np.matmul(np.transpose(hiddenLayer), CE)

# the gradient of the loss with respect to the output layer bias
def gradLossb_o(target, prediction, hiddenLayer):
    ones=np.ones((1,N))
    CE=np.matmul(ones, gradCE(target, prediction))
    return CE

# the gradient of the loss with respect to the hidden layer weights
def gradLossW_h(target, prediction, hiddenLayer, x_in, W_out): # x_in is train_data
    CE=gradCE(target, prediction)
    loss_W_h=np.matmul(np.transpose(x_in), der_relu(hiddenLayer)*np.matmul(CE, np.transpose(W_out)))
    return loss_W_h

# the gradient of the loss with respect to the hidden layer biases
def gradLossb_h(target, prediction, hiddenLayer, x_in, W_out):
    CE=gradCE(target, prediction)
    loss_W_h=np.matmul(CE, np.transpose(W_out))*der_relu(hiddenLayer)
    ones=np.ones((1,N))
    return np.matmul(ones, loss_W_h)
```