

HLCV Assignment 2

2576831 Maria Sargsyan,

2571618 Hui-Syuan Yeh

2576656 Shivam Sharma

Course-Semester

May 10, 2019

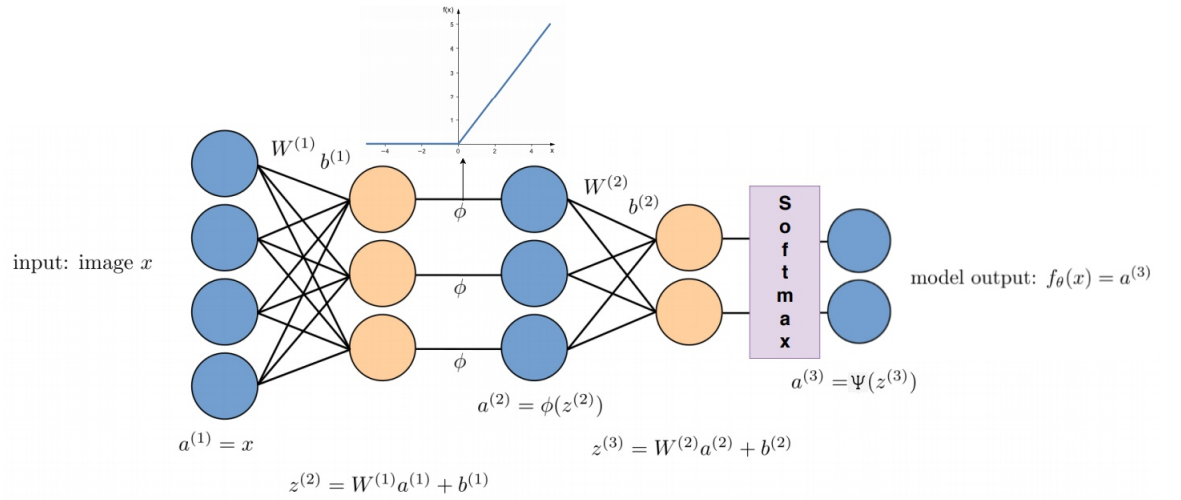


Figure 1: Visualisation of the two layer fully connected network, used in Q1-Q3

The equations from the ex2

$$\varphi(x) = \max(0, x) \quad (1)$$

$$\psi(u)_i = \frac{\exp u_i}{\sum_{j=1}^n \exp u_j} \quad (2)$$

$$a^{(1)} = x \quad (3)$$

$$z^{(2)} = W^{(1)}a^{(1)} + b^{(1)} \quad (4)$$

$$a^{(2)} = \varphi(z^{(2)}) \quad (5)$$

$$z^{(3)} = W^{(2)}a^{(2)} + b^{(2)} \quad (6)$$

$$f_\theta(x) := a^{(3)} = \psi(z^{(3)}) \quad (7)$$

$$J(\theta, x_i, y_i) = -\log P(Y = y_i, X = x_i) \quad (8)$$

$$= -\log f_\theta(x_i)_{y_i} \quad (9)$$

$$= -\log \psi(z^{(3)})_{y_i} \quad (10)$$

$$J(\theta, x_i, y_i) = -\log \frac{\exp z_{y_i}^{(3)}}{\sum_{j=1}^K \exp z_{y_j}^{(3)}} \quad (11)$$

$$J(\theta, (x_i, y_i)_{i=1}^N) = -\frac{1}{N} \sum_{i=1}^N \log \frac{\exp z_{y_i}^{(3)}}{\sum_{j=1}^K \exp z_{y_j}^{(3)}} \quad (12)$$

$$\tilde{J}(\theta) = J(\theta, (x_i, y_i)_{i=1}^N) + \lambda(\|W^{(1)}\|_2^2 + \|W^{(2)}\|_2^2) \quad (13)$$

Problem 2a. Verify that the loss function defined in Eq.12 has the gradient w.r.t $z^{(3)}$ as below.

$$\frac{\partial J}{\partial z^{(3)}}((x_i, y_i)_{i=1}^N) = \frac{1}{N}(\varphi(z^{(3)}) - \Delta) \quad (14)$$

$$\Delta_{ij} = \begin{cases} 1 & \text{if } y_i = j \\ 0 & \text{else} \end{cases}$$

Solution Let's first make the calculations for the case when $y_i \neq j$

$$\begin{aligned} \frac{\partial J}{\partial z_j^{(3)}} &= -\frac{1}{N} \frac{1}{a_{y_i}^{(3)}} \cdot -\frac{\exp z_{y_i}^{(3)} \cdot \exp z_j^{(3)}}{(\sum_{k=1}^K \exp z_{y_k}^{(3)})^2} \\ &= \frac{1}{N} \frac{1}{a_{y_i}^{(3)}} \cdot \frac{\exp z_{y_i}^{(3)}}{\sum_{k=1}^K \exp z_{y_k}^{(3)}} \cdot \frac{\exp z_j^{(3)}}{\sum_{k=1}^K \exp z_{y_k}^{(3)}} \\ &= \frac{1}{N} \frac{1}{a_{y_i}^{(3)}} \cdot a_{y_i}^{(3)} \cdot a_j^{(3)} = \frac{1}{N} a_j^{(3)} = \frac{1}{N} \psi_j^{(3)} \end{aligned}$$

If $y_i = j$

$$\begin{aligned}\frac{\partial J}{\partial z_j^{(3)}} &= -\frac{1}{N} \frac{1}{a_{y_i}^{(3)}} \cdot \left(a_{y_i}^{(3)} - \left(\frac{\exp z_{y_i}^{(3)}}{\sum_{k=1}^K \exp z_{y_k}^{(3)}} \right)^2 \right) \\ &= -\frac{1}{N} \frac{1}{a_{y_i}^{(3)}} \cdot (a_{y_i}^{(3)} - (a_{y_i}^{(3)})^2) = \frac{1}{N} \cdot (a_{y_i}^{(3)} - 1)\end{aligned}$$

Problem 2b. To compute the effect of the weight matrix W_2 on the loss in Eq.12 incurred by the network, we compute the partial derivatives of the loss function with respect to W_2 . This is done by applying the chain rule. Verify that the partial derivative of the loss w.r.t W_2 is :

$$\frac{\partial J}{\partial W^{(2)}}((x_i, y_i)_{i=1}^N) = \frac{\partial J}{\partial z^{(3)}} \cdot \frac{\partial z^{(3)}}{\partial W^{(2)}} \quad (15)$$

$$= \frac{1}{N} (\psi(z^{(3)}) - \Delta) a^{(2)'} \quad (16)$$

Similarly verify that:

$$\frac{\partial \tilde{J}}{\partial W^{(2)}} = \frac{1}{N} (\psi(z^{(3)}) - \Delta) a^{(2)'} + 2\lambda W^{(2)} \quad (17)$$

Solution Let us calculate the derivatives for a single entry in the matrix W_2 , namely w_{ij} . Also, note that $\frac{\partial J}{\partial z^{(3)}}$ and $\frac{\partial z^{(3)}}{\partial W_{ij}^{(2)}}$ have the same dimensions. Now

$$\frac{\partial z_{ek}^{(3)}}{\partial W_{ij}^{(2)}} = \begin{cases} a_{ej}^{(2)} & \text{if } k = j \\ 0 & \text{else} \end{cases}$$

Similarly,

$$\frac{\partial \lambda \|W^{(2)}\|_2^2}{\partial W_{ij}^{(2)}} = 2\lambda W_{ij}^{(2)} \quad (18)$$

Now,

$$\begin{aligned}\frac{\partial J}{\partial W_{ij}^{(2)}}((x_k, y_k)_{k=1}^N) &= \frac{\partial J}{\partial z^{(3)}} \cdot \frac{\partial z^{(3)}}{\partial W_{ij}^{(2)}} \\ &= \sum_{e=1}^N \sum_{c=1}^{H_2} \frac{\partial J}{\partial z_{ec}^{(3)}} \cdot \frac{\partial z_{ec}^{(3)}}{\partial W_{ij}^{(2)}} = \sum_{e=1}^N \frac{\partial J}{\partial z_{ei}^{(3)}} \cdot \frac{\partial z_{ei}^{(3)}}{\partial W_{ij}^{(2)}} = \sum_{e=1}^N \frac{\partial J}{\partial z_{ei}^{(3)}} \cdot a_{ej}^{(2)}\end{aligned}$$

Thus,

$$\frac{\partial J}{\partial W_{ij}^{(2)}}((x_k, y_k)_{k=1}^N) = a^{(2)\top} \cdot \frac{\partial J}{\partial z^{(3)}} \quad (19)$$

Finally,

$$\frac{\partial \tilde{J}}{\partial W^{(2)}} = a^{(2)\top} \cdot \frac{1}{N} (\psi(z^{(3)}) - \Delta) + 2\lambda W^{(2)} \quad (20)$$

Problem 2c. Derive the expressions for the derivatives of the regularized loss in w.r.t other variables

Solution Let's calculate $\frac{\partial \tilde{J}}{\partial b^{(2)}}$ which is the same as $\frac{\partial J}{\partial b^{(2)}}$

$$\frac{\partial z_{ec}^{(3)}}{\partial b_j^{(2)}} = \begin{cases} 1 & \text{if } c = j \\ 0 & \text{else} \end{cases}$$

$$\begin{aligned} \frac{\partial J}{\partial b_j^{(2)}}((x_k, y_k)_{k=1}^N) &= \frac{\partial J}{\partial z^{(3)}} \cdot \frac{\partial z^{(3)}}{\partial b_j^{(2)}} \\ &= \sum_{e=1}^N \sum_{c=1}^{H_2} \frac{\partial J}{\partial z_{ec}^{(3)}} \cdot \frac{\partial z_{ec}^{(3)}}{\partial b_j^{(2)}} = \sum_{e=1}^N \frac{\partial J}{\partial z_{ej}^{(3)}} \cdot \frac{\partial z_{ej}^{(3)}}{\partial b_j^{(2)}} = \sum_{e=1}^N \frac{\partial J}{\partial z_{ej}^{(3)}} \end{aligned}$$

Thus,

$$\frac{\partial J}{\partial b^{(2)}}((x_k, y_k)_{k=1}^N) = \mathbf{1}^\top \frac{\partial J}{\partial z^{(3)}}, \quad (21)$$

where $\mathbf{1} \in \mathbb{R}^N$ is a vector full of ones Now, to calculate the rest of the gradients we need to calculate first $\frac{\partial J}{\partial z^{(2)}}$ since

$$\frac{\partial J}{\partial b^{(1)}} = \frac{\partial J}{\partial z^{(2)}} \cdot \frac{\partial z^{(2)}}{\partial b^{(1)}}$$

$$\frac{\partial J}{\partial W^{(1)}} = \frac{\partial J}{\partial z^{(2)}} \cdot \frac{\partial z^{(2)}}{\partial W^{(1)}} \quad (22)$$

So,

$$\frac{\partial J}{\partial z^{(2)}} = \frac{\partial J}{\partial z^{(3)}} \cdot \frac{\partial z^{(3)}}{\partial a^{(2)}} \cdot \frac{\partial a^{(2)}}{\partial z^{(2)}} \quad (23)$$

$$\begin{aligned} \frac{\partial J}{\partial a_{ij}^{(2)}} &= \frac{\partial J}{\partial z^{(3)}} \cdot \frac{\partial z^{(3)}}{\partial a_{ij}^{(2)}} = \sum_{e=1}^N \sum_{c=1}^{H_2} \frac{\partial J}{\partial z_{ec}^{(3)}} \cdot \frac{\partial z_{ec}^{(3)}}{\partial a_{ij}^{(2)}} \\ &= \sum_{c=1}^{H_2} \frac{\partial J}{\partial z_{ic}^{(3)}} \cdot \frac{\partial z_{ic}^{(3)}}{\partial a_{ij}^{(2)}} = \sum_{c=1}^{H_2} \frac{\partial J}{\partial z_{ic}^{(3)}} \cdot W_{jc}^{(2)} \end{aligned}$$

Thus

$$\frac{\partial J}{\partial a^{(2)}} = \frac{\partial J}{\partial z^{(3)}} \cdot W^{(2)\top} \quad (24)$$

$$\frac{\partial a_{ec}^{(2)}}{\partial z_{ij}^{(2)}} = \begin{cases} \phi'(z_{ij}^{(2)}) & \text{if } e = i \text{ and } c = j \\ 0 & \text{else} \end{cases}$$

$$\phi'(x) = \begin{cases} 1 & \text{if } x \geq 0 \\ 0 & \text{else} \end{cases}$$

Now,

$$\frac{\partial J}{\partial z^{(2)}_{ec}} = \frac{\partial J}{\partial a^{(2)}} \cdot \frac{\partial a^{(2)}}{\partial z^{(2)}_{ec}} = \frac{\partial J}{\partial a^{(2)}_{ec}} \cdot \frac{\partial a^{(2)}_{ec}}{\partial z^{(2)}_{ec}} \quad (25)$$

Which is just element-wise multiplication

Following similar arguments to 21

$$\frac{\partial J}{\partial b^{(1)}}((x_k, y_k)_{k=1}^N) = \mathbf{1}^\top \frac{\partial J}{\partial z^{(3)}} \cdot W^{(2)} \cdot \phi'(z^{(2)}) \quad (26)$$

$$\frac{\partial J}{\partial W^{(1)}}((x_k, y_k)_{k=1}^N) = a^{(1)\top} \left(\frac{\partial J}{\partial z^{(3)}} \cdot W^{(2)} \cdot \phi'(z^{(2)}) \right) \quad (27)$$

Finally

$$\frac{\partial \tilde{J}}{\partial W^{(2)}} = \frac{\partial J}{\partial W^{(2)}} + 2\lambda W^{(2)} = a^{(2)\top} \cdot \frac{1}{N} (\psi(z^{(3)}) - \Delta) + 2\lambda W^{(2)} \quad (28)$$

$$\frac{\partial \tilde{J}}{\partial W^{(1)}} = \frac{\partial J}{\partial W^{(1)}} + 2\lambda W^{(1)} = a^{(1)\top} \cdot \left(\frac{\partial J}{\partial z^{(3)}} \cdot W^{(2)} * \phi'(z^{(2)}) \right) + 2\lambda W^{(1)} \quad (29)$$

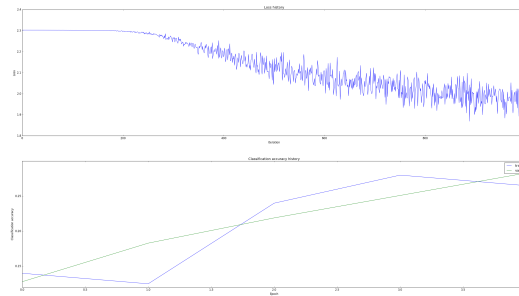
$$\frac{\partial \tilde{J}}{\partial b^{(2)}} = \frac{\partial J}{\partial b^{(2)}} = \mathbf{1}^\top \frac{\partial J}{\partial z^{(3)}} \quad (30)$$

$$\frac{\partial \tilde{J}}{\partial b^{(1)}} = \frac{\partial J}{\partial b^{(1)}} = \mathbf{1}^\top \frac{\partial J}{\partial z^{(3)}} \cdot W^{(2)} * \phi'(z^{(2)}) \quad (31)$$

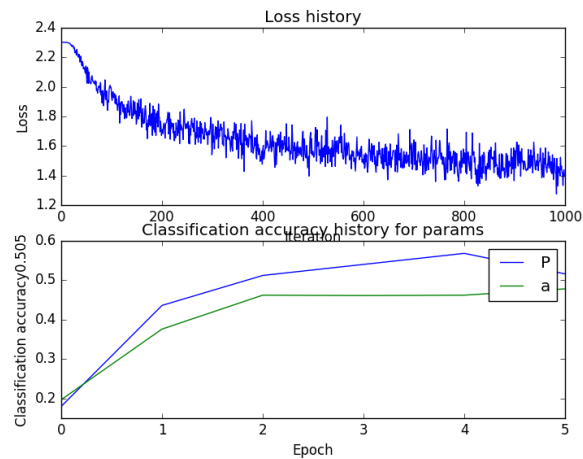
$$\frac{\partial J}{\partial z^{(3)}}((x_i, y_i)_{i=1}^N) = \frac{1}{N} (\phi(z^{(3)}) - \Delta) \quad (32)$$

Problem 3b. There are several pointers provided in the comments in the ex2 FCnet.py to help you understand why the network might be underperforming (Line 224-250). Once you have tuned your hyper parameters, and get validation accuracy greater than 48 report the performance

As, it was suggested in the code. From the graph above, it is clearly seen the loss during the first iterations is decreasing very slowly. Actually, only by changing the learning rate from 0.0001 to 0.001, the accuracy increased up to 47 percents.



We could increase the accuracies up to 0.505 percents by increasing the number of hidden layers to 80 and setting batch size to 250, learning decay to 0.99 and regularization to 0.12 and test accuracy is 0.483



```
(Validation data shape: ', (1000, 3072))
(Validation labels shape: ', (1000,))
(Test data shape: ', (1000, 3072))
(Test labels shape: ', (1000,))
iteration 0 / 1000: loss 2.302849
iteration 100 / 1000: loss 1.926791
iteration 200 / 1000: loss 1.696625
iteration 300 / 1000: loss 1.576330
iteration 400 / 1000: loss 1.576058
iteration 500 / 1000: loss 1.562522
iteration 600 / 1000: loss 1.482224
iteration 700 / 1000: loss 1.642342
iteration 800 / 1000: loss 1.482232
iteration 900 / 1000: loss 1.558295
(Validation accuracy: ', 0.505)
(Test accuracy: ', 0.483)
ex02/ex2/python-code-assignment master X
```

Problem N. Now that you can train the two layer network to achieve reasonable performance, try increasing the network depth to see if you can improve the performance. Experiment with networks of at least 2, 3, 4, and 5 layers, of your chosen configuration. Report the training and validation accuracies for these models and discuss your observations. Run the evaluation on the test set with your best model and report the test accuracy. (report, 4 points)

Solution As far as we have tried, models with several fully connected layers were not learning. For example, given the dimensions of hidden layers [80,50,50,50] validation accuracy was increasing very slowly. Presumably, appending more fully connected layers may decrease the spatial information contained in the image which is better captured in the first layer.

```

Validataion accuracy is: 52 %
Epoch [26/30], Step [100/245], Loss: 1.0116
Epoch [26/30], Step [200/245], Loss: 0.9814
Validataion accuracy is: 52 %
Epoch [27/30], Step [100/245], Loss: 0.9067
Epoch [27/30], Step [200/245], Loss: 0.9393
Validataion accuracy is: 52 %
Epoch [28/30], Step [100/245], Loss: 0.9395
Epoch [28/30], Step [200/245], Loss: 0.9157
Validataion accuracy is: 52 %
Epoch [29/30], Step [100/245], Loss: 0.8019
Epoch [29/30], Step [200/245], Loss: 0.9372
Validataion accuracy is: 53 %
Epoch [30/30], Step [100/245], Loss: 1.0146
Epoch [30/30], Step [200/245], Loss: 1.0075
Validataion accuracy is: 52 %

ex02/ex2/python-code-assignment master X
> python ex2_pytorch.py
Using device: cpu
Accuracy of the network on the 1000 test images: 54 %
Epoch [4/30], Step [100/245], Loss: 2.3030
Epoch [4/30], Step [200/245], Loss: 2.3027
Validataion accuracy is: 7 %
Epoch [5/30], Step [100/245], Loss: 2.3026
Epoch [5/30], Step [200/245], Loss: 2.3029
Validataion accuracy is: 7 %
Epoch [6/30], Step [100/245], Loss: 2.3020
Epoch [6/30], Step [200/245], Loss: 2.3032
Validataion accuracy is: 7 %
Epoch [7/30], Step [100/245], Loss: 2.3030
Epoch [7/30], Step [200/245], Loss: 2.3028
Validataion accuracy is: 7 %
Epoch [8/30], Step [100/245], Loss: 2.3023
Epoch [8/30], Step [200/245], Loss: 2.3031
Validataion accuracy is: 7 %
Epoch [9/30], Step [100/245], Loss: 2.3028
Epoch [9/30], Step [200/245], Loss: 2.3028
Validataion accuracy is: 7 %
Epoch [10/30], Step [100/245], Loss: 2.3026
Epoch [10/30], Step [200/245], Loss: 2.3026
Validataion accuracy is: 7 %
Epoch [11/30], Step [100/245], Loss: 2.3026
Epoch [11/30], Step [200/245], Loss: 2.3028

```

The first image refers to the final results. The second refers to the training the net with more than 2 layers