

Name: Nattapat Juthaprachakul, Student ID: 301350117

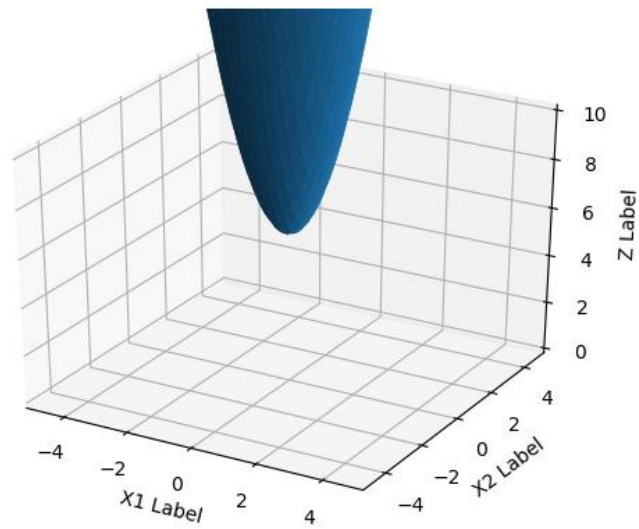
CMPT-726 Machine Learning: Assignment 2

1. Softmax for multi-class classification

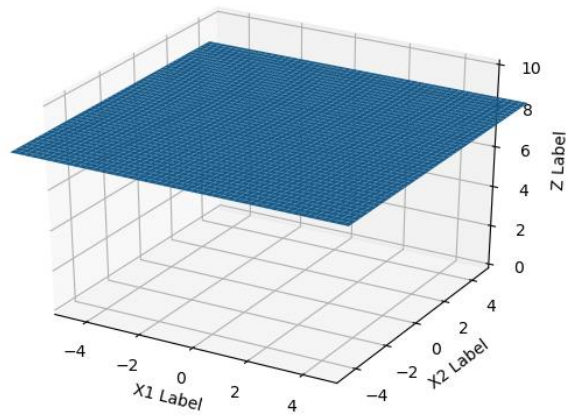
- 1.1 ANS: The probability at the green point for each class is equally likely to be any of these 3 classes (0.33 percent for each class).
- 1.2 ANS: The probability of input (green point) depends on which direction the green dot heads to (moving along the line in this case). For example, if the green dot moves downward along the red line, the probability of input x being class 1 and 3 (region 1 and 3) is more likely than class 2 (region 2). Also, the more the green dot moves downward, the more unlikely the input is classified as class 2 (region 2). In sum, when green dot moves along red line downward, the probability of input classified to be either class 1 or 3 is equally likely but unlikely to be class 2. This logic applies to both moving leftward and rightward as well. (Moving leftward probability of input being classified to be either class 2 and 3 is equally likely but unlikely to be class 1 and moving rightward probability of input being classified to be either class 1 and 2 is equally likely but unlikely to be class 3.)
- 1.3 ANS: The probability of inputs being classified as the following class depends on the region that the green dot is in. For example, if the green dot is in region 1, the input is likely to be classified as class 1 than classified as class 2 and 3.

2. Generalized Linear model for classification

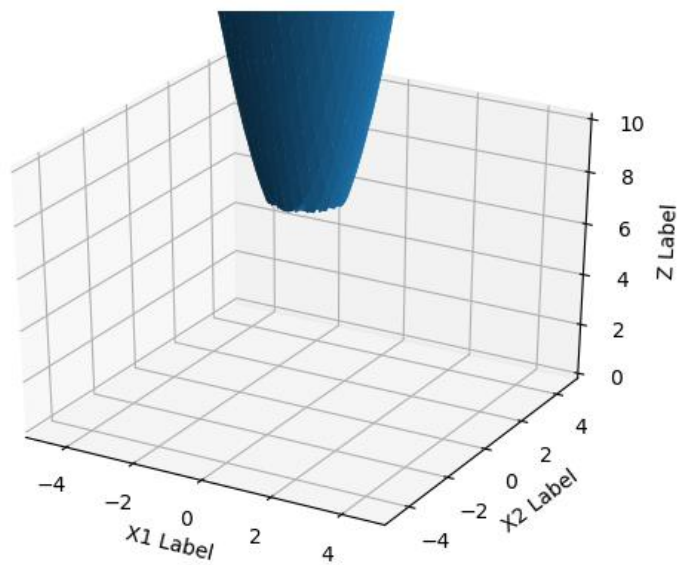
2.1



2.2



2.3



3. 3.1)ANS:

$$a_2^{(2)} = w_{21}^{(1)} x_1 + w_{22}^{(1)} x_2 + w_{23}^{(1)} x_3, \quad z_2^{(2)} = h(a_2^{(2)})$$

$$E_n(w) = \frac{1}{2} (y(x_n, w) - t_n)^2 = \frac{1}{2} (z^{(4)} - t_n)^2$$

$$\rightarrow \frac{\partial E_n(w)}{\partial a_1^{(4)}} = \frac{\partial E_n(w)}{\partial z^{(4)}} \frac{\partial z^{(4)}}{\partial a_1^{(4)}} = (z^{(4)} - t_n) [h'(a_1^{(4)})]$$

$$\rightarrow \frac{\partial E_n(w)}{\partial w_{12}^{(3)}} = \frac{\partial E_n(w)}{\partial z^{(4)}} \frac{\partial z^{(4)}}{\partial a_1^{(4)}} \frac{\partial a_1^{(4)}}{\partial w_{12}^{(3)}} = \square \times \frac{\partial a_1^{(4)}}{\partial w_{12}^{(3)}}$$

since $\frac{\partial a_1^{(4)}}{\partial w_{12}^{(3)}} = w_{11}^{(3)} z_{11}^{(3)} + w_{12}^{(3)} z_{12}^{(3)} + w_{13}^{(3)} z_{13}^{(3)}$

$$\frac{\partial a_1^{(4)}}{\partial w_{12}^{(3)}} = 0 + z_{12}^{(3)} + 0 = z_{12}^{(3)}$$

$$\frac{\partial E_n(w)}{\partial w_{12}^{(3)}} = (z^{(4)} - t_n) (h'(a_1^{(4)})) (z_{12}^{(3)})$$

3.2)ANS:

$$\rightarrow \frac{\partial E_n(w)}{\partial a_1^{(3)}} = \frac{\partial E_n(w)}{\partial z^{(4)}} \frac{\partial z^{(4)}}{\partial a_1^{(4)}} \frac{\partial a_1^{(4)}}{\partial z^{(3)}} \frac{\partial z^{(3)}}{\partial a_1^{(3)}}$$

$$= \delta_1^{(4)} h'(a_1^{(4)}) (w_{11}^{(3)}) h'(a_1^{(3)})$$

$$\rightarrow \frac{\partial E_n(w)}{\partial w_{11}^{(2)}} = \delta_1^{(4)} h'(a_1^{(4)}) w_{11}^{(3)} h'(a_1^{(3)}) \frac{\partial a_1^{(3)}}{\partial w_{11}^{(2)}}$$

$$= \delta_1^{(4)} h'(a_1^{(4)}) w_{11}^{(3)} h'(a_1^{(3)}) z_{11}^{(2)}$$

3.3)ANS:

$$\rightarrow \frac{\partial E_n(w)}{\partial a_3^{(2)}} = \sum_k \frac{\partial E_n}{\partial a_k} \frac{\partial a_k}{\partial a_3^{(2)}} \quad \text{recall } \delta_j = \sum_k \frac{\partial E_n}{\partial a_k} \frac{\partial a_k}{\partial a_j}$$

$$= \delta_1^{(3)} \frac{\partial a_1^{(3)}}{\partial a_3^{(2)}} + \delta_2^{(3)} \frac{\partial a_2^{(3)}}{\partial a_3^{(2)}} + \delta_3^{(3)} \frac{\partial a_3^{(3)}}{\partial a_3^{(2)}}$$

$$= \delta_1^{(3)} \frac{\partial a_1^{(3)}}{\partial z_3^{(2)}} \frac{\partial z_3^{(2)}}{\partial a_3^{(2)}} + \delta_2^{(3)} \frac{\partial a_2^{(3)}}{\partial z_3^{(2)}} \frac{\partial z_3^{(2)}}{\partial a_3^{(2)}} + \delta_3^{(3)} \frac{\partial a_3^{(3)}}{\partial z_3^{(2)}} \frac{\partial z_3^{(2)}}{\partial a_3^{(2)}}$$

$$= \delta_1^{(3)} w_{13}^{(2)} h'(a_3^{(2)}) + \delta_2^{(3)} w_{23}^{(2)} h'(a_3^{(2)}) + \delta_3^{(3)} w_{33}^{(2)} h'(a_3^{(2)})$$

$$\rightarrow \frac{\partial E_n(w)}{\partial w_{11}^{(1)}} = \frac{\partial E_n(w)}{\partial a_1^{(2)}} \frac{\partial a_1^{(2)}}{\partial w_{11}^{(1)}}$$

$$= \left[\left(\delta_1^{(3)} w_{11}^{(2)} + \delta_2^{(3)} w_{21}^{(2)} + \delta_3^{(3)} w_{31}^{(2)} \right) h'(a_1^{(2)}) \right] x_1$$

4. 4.1)ANS

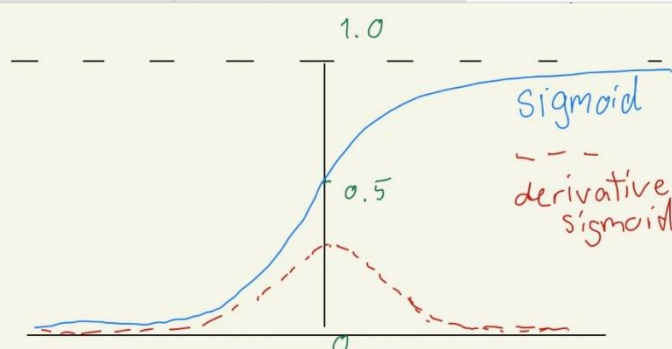
Q4)

$$\begin{aligned}
 4.1) \frac{\partial E_n(w)}{\partial w_{11}^{(1)}} &= \frac{\partial E_n(w)}{\partial z_1^{(153)}} \frac{\partial z_1^{(153)}}{\partial a_1^{(153)}} \frac{\partial a_1^{(153)}}{\partial w_{11}^{(152)}} \frac{\partial w_{11}^{(152)}}{\partial z_1^{(154)}} \\
 &\quad \frac{\partial z_1^{(152)}}{\partial a_1^{(152)}} \frac{\partial a_1^{(152)}}{\partial w_{11}^{(151)}} \dots \frac{\partial a_1^{(2)}}{\partial w_{11}^{(1)}} \\
 &= \frac{\partial E(w)}{\partial z_1^{(153)}} \left[\prod_{i=2}^{153} \frac{\partial z_1^{(i)}}{\partial a_1^{(i)}} \frac{\partial a_1^{(i)}}{\partial w_{11}^{(i-1)}} \right] \quad \text{---}
 \end{aligned}$$

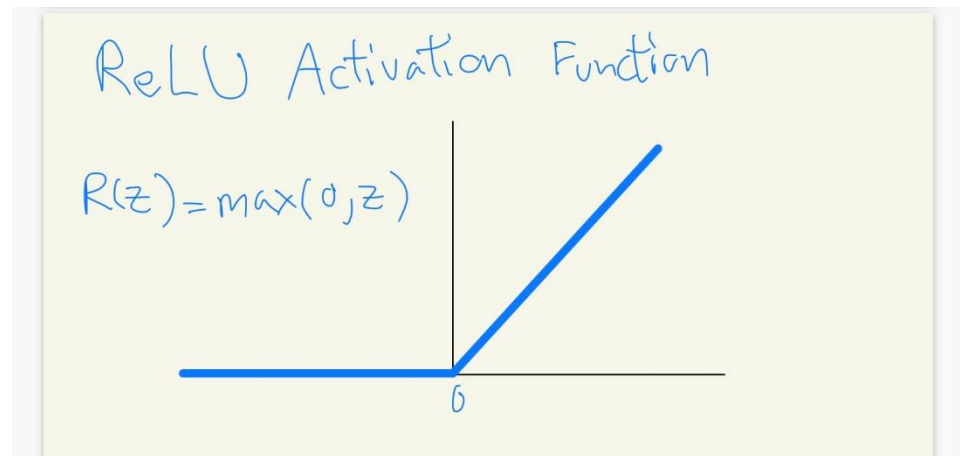
4.2)ANS: The learning will be very slow if the learning rate is very small and the area of the update is in the top and bottom curve (very flat area). The update is slow because the derivative is very small (small rate of change/update). For the softmax function, the gradient of weight could be small or zero when we do backpropagate to modify the weight to minimize the cost function through many layers and many connected nodes as the majority of derivative value of sigmoid lies between 0 and 0.25 (dotted red curve in the graph below), the multiplication of number between 0 and 1 many times will make the values smaller over time and become zero in some connection.

$$\begin{aligned}
 \text{Sol}^n \frac{\partial G(a)}{\partial a} &= \frac{\partial}{\partial a} \left[\frac{1}{1+e^{-a}} \right] = \frac{\partial}{\partial a} \left[(1+e^{-a})^{-1} \right] \\
 &= -(1+e^{-a})^{-2} (-e^{-a})
 \end{aligned}$$

$$\begin{aligned}
 &\downarrow \\
 &= \frac{e^{-a}}{(1+e^{-a})^2} = \frac{1}{1+e^{-a}} \cdot \frac{e^{-a}}{1+e^{-a}} \\
 &= \frac{1}{1+e^{-a}} \times \frac{(1+e^{-a})-1}{1+e^{-a}} \\
 &= \frac{1}{1+e^{-a}} \left[\frac{1+e^{-a}}{1+e^{-a}} - \frac{1}{1+e^{-a}} \right] \\
 &= \left(\frac{1}{1+e^{-a}} \right) \left(1 - \frac{1}{1+e^{-a}} \right) \\
 &= G(a)(1-G(a)) \quad \text{---}
 \end{aligned}$$



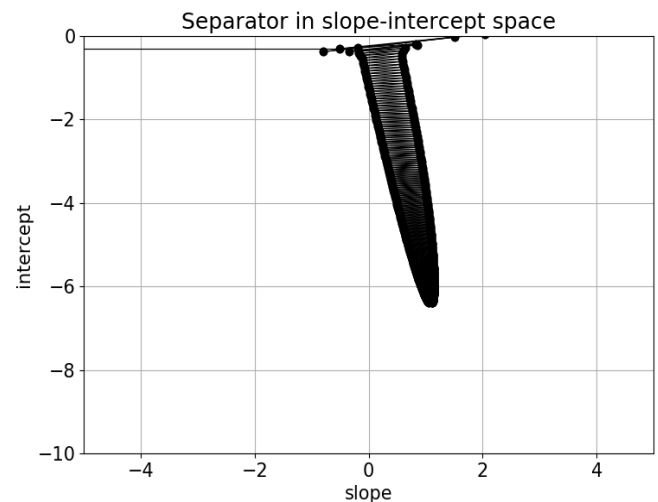
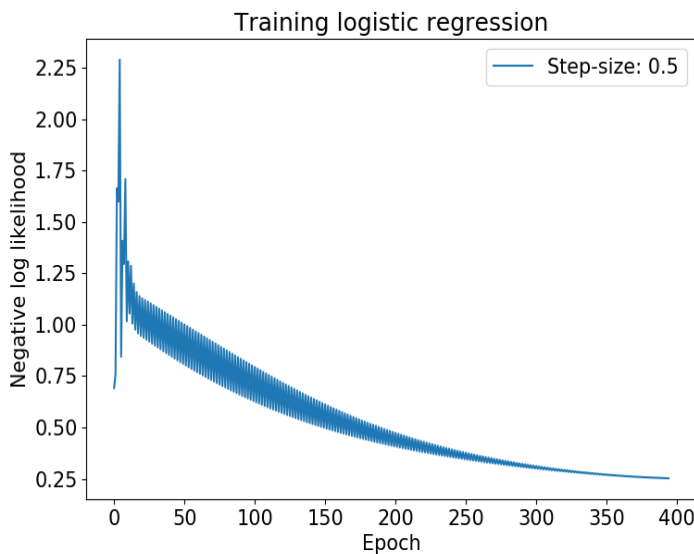
4.3)ANS: When the inputs to ReLU is equal or smaller than zero (the large negative bias term is learned which makes weighted sum of inputs becomes zero or negative), the output of ReLU will be zero, making the derivative become zero as well.



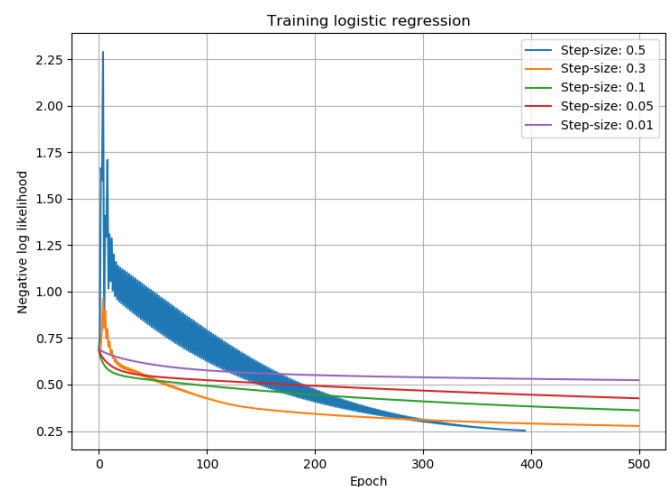
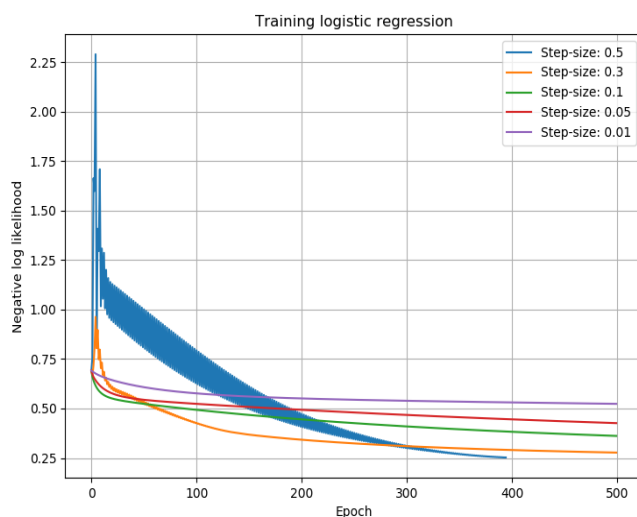
4.4)ANS: The gradient of weight could be zero when we do update/backpropagate to modify the weight to minimize the cost function through many layers and many connected nodes (bipartite). This might be the result when the majority of output from ReLU could become zero as the summation of weight term and bias is zero or negative. The function gradient at zero becomes zero as well. (the gradient descent learning will not alter the weights) This situation is called 'dead' ReLU.

5. 5.1)ANS: When learning rate is too large, it causes drastic weight updates which lead to divergent behaviors (the oscillating curves) observed in the left and right pictures.

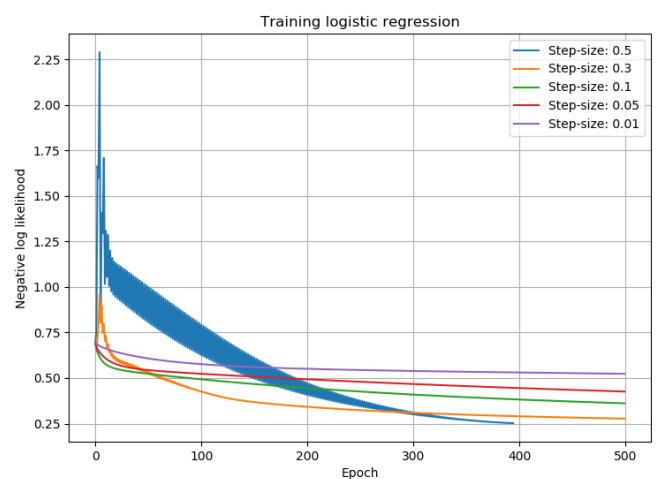
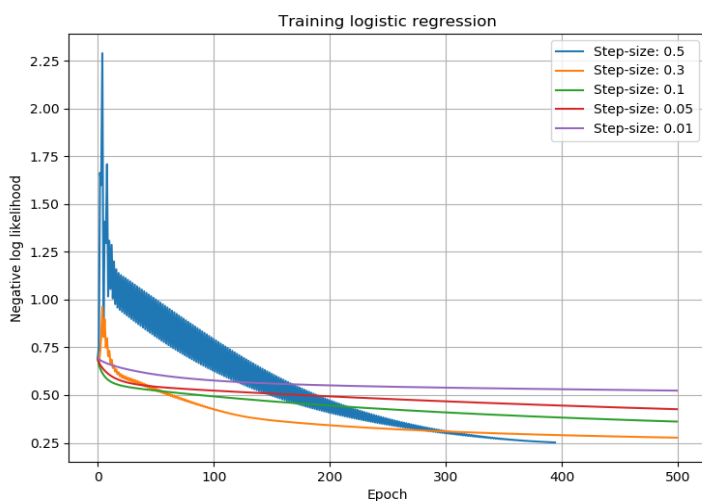
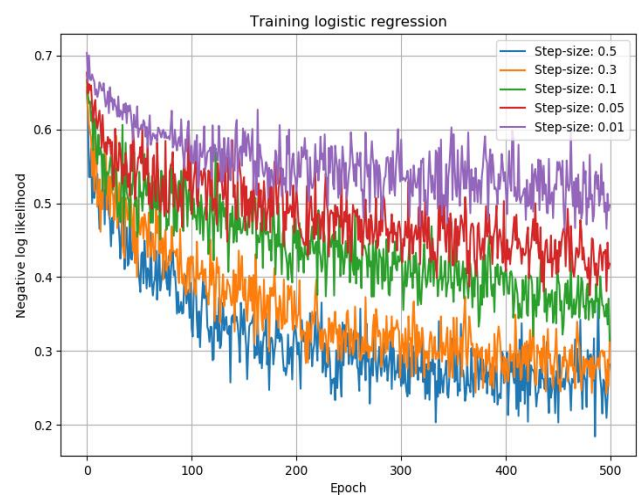
The performance of the model (such as its loss on the training dataset) will oscillate over training epochs. Oscillating performance is caused by weights that diverge.



5.2)ANS: in our case, the 0.5 learning rate finds the lowest training error at the fewest number of epoch. As we can see the larger the learning rate is, the faster the training becomes (given the same training error).



5.3)ANS: in our example, it is not so obvious that which update technique is better than another in term of speed(number of epoch) and performance(error rate). We can see that for the SGD graph is very oscillating since it calculates one data point then makes an update. However, in the setting with big data sets, Gradient Descent takes time to calculate cost or gradient as it needs to sum over all data points. Nonetheless, we do not need to have exact gradient to minimize the cost function in a given iteration. The approximation of gradient is enough; therefore, the Stochastic Gradient descent approximates the gradient using just only one data point at a time which in turn saves lots of time compared to summing over all data.



6. Fine-Tuning a Pre-trained Network

Settings: gaming laptop with a single Nvidia GTX 1050Ti

6.1) Main task that I do:

-Write a Python function to be used at the end of training that generates HTML output showing each test image and its classification scores. You could produce an HTML table output for example. (You can convert the HTML output to PDF or use screen shots.)

6.2) Other tasks:

-Try applying L2 regularization to the coefficients in the small networks we added.

-Try modifying the structure of the new layers that were added on top of ResNet20.

The main code is to evaluate the models and to save loss and classification score (max value from softmax output) for each test image. The next code generates each test image with its classification score and saves them all in our specified folder. The final code is to generate HTML table and convert it into PDF. I attached the filename: image_score_html.html and test_image.pdf.

I also created new additional layers (modifying the original assignment code of just one layer) on top of ResNet20 as follow: [feed forward layer -> batch normalization -> feed forward layer -> softmax output] with dropout technique, L2 regularization(weight decay) and specifiable number of hidden nodes.

Results: different epochs on training set of 50k inputs and 1 epoch on test set of 10k inputs.

Note: Original feed-forward layer(64 nodes) denotes as fc(64)

Additional feed-forward layer(256 nodes) denotes as fc(256)

1. 1 epoch training, ResNet20 + fc(64)
Running optimization on: fc(64)
Accuracy: 65.10%
2. 10 epoch training, ResNet20 + fc(64)
Running optimization on: fc(64)
Accuracy: 68.55%
3. 10 epoch training, ResNet20 + fc(64), L2 regularization(0.001)
Running optimization on: fc(64)
Accuracy: 68.65%
4. 10 epoch training, ResNet20 + fc(64) + fc(256), L2 regularization(0.001)
Running optimization on: fc(256)
Accuracy: 68.52%
5. 10 epoch training, ResNet20 + fc(64)+ fc(256)
Running optimization on: ResNet20 + fc(256)
Accuracy: 84.46%
6. 10 epoch training, ResNet20 + fc(64)
Running optimization on: ResNet20 + fc(64)

Accuracy: 84.44%

7. 10 epoch training, ResNet20 + fc(64)

Running optimization on: ResNet20

Accuracy: 80.33%

8. 10 epoch training, ResNet20 + fc(64), dropout at 0.2

Running optimization on: ResNet20 + fc(64)

Accuracy: 84.55%

