

# DD2424 Deep Learning in Data Science

## Assignment 2

Rui Shi  
srui@kth.se

April 25, 2022

## 1 Introduction

In this assignment, all the functions and codes are implemented on Python. In this section, I trained and tested a two layer network with multiple outputs to classify images from the CIFAR-10 datasets.

## 2 Check gradients

In this section, I compared the difference between the analytical gradient and numerical gradient. Here are part of the results. From the results, we can see that the difference is tiny.

```
Step3:check gradients

xmini=xtr[1:20,1:10]
Ymini=Ytr[1:20,1:10]
ymini=ytr[1:20,1:10]

W1,W2, b1,b2= ini_parameters(xmini, Ymini)

gw1,gw2, gwb1,gwb2=ComputeGradients(x=xmini, Y= Ymini,y=ymini, W1=W1, W2=W2, b1=b1,b2=b2, h=5,lambda=0)
gw1,gw2, gwb1,gwb2=compute_gradient(x=xmini, Y= Ymini, W1=W1,W2=W2,b1=b1,b2=b2,lambda=0,sym=False)
print(gwb1-gwb2)
print(gwb2-gwb1)
print(gw1-gw2)
print(gw2-gw1)

[0] ✓ 0.1s

Output exceeds the size limit. Open the full output data in a text editor
[[[-1.61027788e-11]
 [ 3.68125253e-11]
 [-1.19584898e-11]
 [-1.63271341e-11]
 [-3.49435064e-12]
 [-1.78358439e-11]
 [ 1.65516986e-11]
 [ 1.6553537e-11]
 [ 3.13862099e-11]
 [-2.93533253e-12]]
 [[-1.61027788e-11]
 [ 3.68125253e-11]
 [-1.19584898e-11]
 [-1.63271341e-11]
 [-3.49435064e-12]
 [-1.78358439e-11]
 [ 1.65516986e-11]
 [ 1.6553537e-11]
 [ 3.13862099e-11]
 [-2.93533253e-12]]
 [[4.54702506e-08 3.65514955e-08 5.79807417e-08 7.48985690e-08
 6.89843662e-08 6.89961162e-08 6.81278095e-08 6.88846869e-08
 6.55761434e-08 6.65172836e-08 6.58879982e-08 6.67281631e-08
 6.84739266e-08 6.63658246e-08 6.42886711e-08 5.17356119e-08
 4.93928332e-08 5.22488186e-08 4.43581961e-08]
 ...
 [2.2389994e-07 5.72628818e-07 4.27336836e-07 2.18899279e-07
 3.44825762e-08 8.82799449e-07 9.16236819e-08 1.37844138e-07
 3.82368920e-07 3.23278423e-08 8.33228309e-08 3.71444815e-07
 6.48983177e-07 1.16864279e-06]]
```

Figure 1: Difference between analytical gradient and numerical gradient

After that, I move to sanity check. I tried and trained my network on 100 examples with regularization turned off. From the results, I can overfit to the training data and get a very low loss on the training data after training for a sufficient number of epochs. Here are the graphs representing the training and validation loss. For the training data, the loss could be less than 0.003, and the accuracy could be around 0.99.

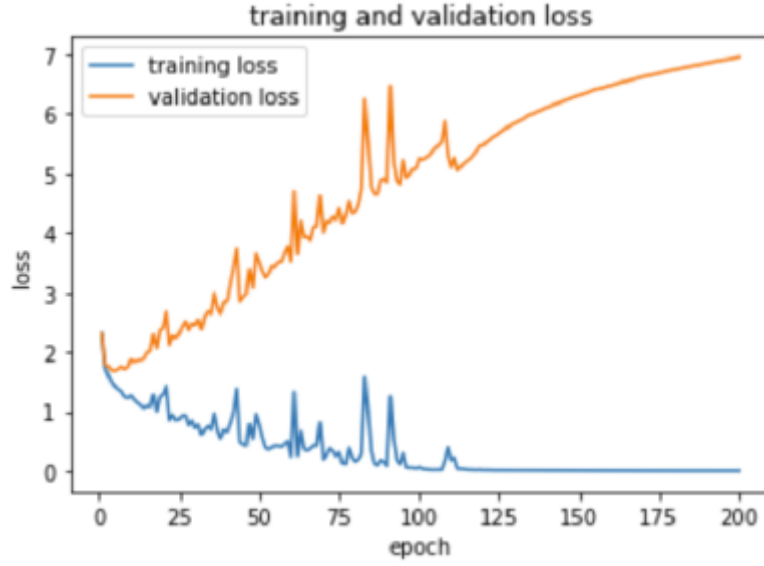


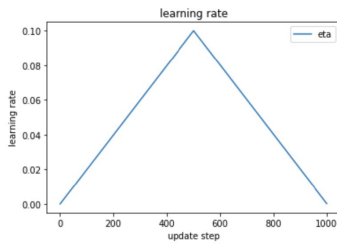
Figure 2: Training and validation loss

### 3 Cycling Learning Rate

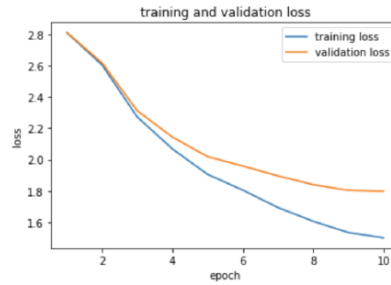
#### 3.1 One cycle training

The parameters are defined as below:

$\eta_{\min} = 1e-5$ ,  $\eta_{\max} = 1e-1$ ,  $\lambda = 0.01$ ,  $ns = 500$ . Here are the graphs representing learning rate, training and validation loss, training and validation accuracy.



(a) Learning rate



(b) Training and validation loss



(c) Training and validation accuracy

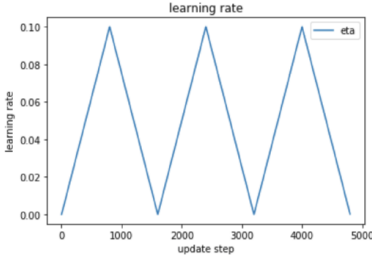
Figure 3: One cycle training

Finally, the accuracy could be up to 44.97%.

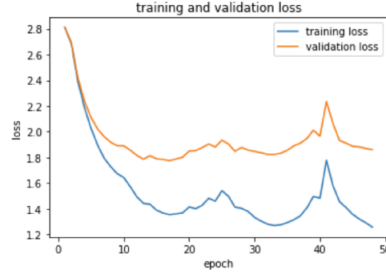
#### 3.2 Three cycle training

The parameters are defined as below:

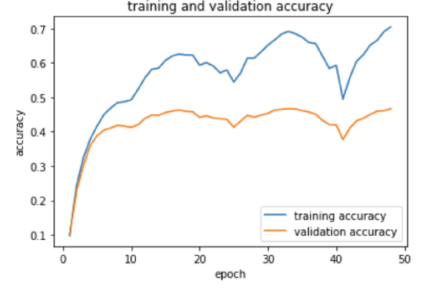
$\eta_{\min} = 1e-5$ ,  $\eta_{\max} = 1e-1$ ,  $\lambda = 0.01$ ,  $ns = 800$ . Here are the graphs representing learning rate, training and validation loss, training and validation accuracy.



(a) Learning rate



(b) Training and validation loss



(c) Training and validation accuracy

Figure 4: One cycle training

Finally, the accuracy could be up to 46.75%.

## 4 Lambda searching

### 4.1 Coarse search

The parameters are defined as below:

$l_{\min} = -5$ ,  $l_{\max} = -1$ ,  $ns = 200$ ,  $\eta_{\min} = 1e-5$ ,  $\eta_{\max} = 1e-1$ ,  $n_{\text{cycle}} = 2$ . After coarse searching for lambda, the parameters of the three best networks are:

lambda	Test accuracy
$10^{(-3.196)}$	51.42%
$10^{(-2.325)}$	51.31%
$10^{(-2.826)}$	51.48%

### 4.2 Fine search

The parameters are defined as below:

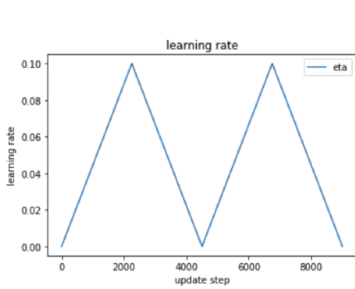
$l_{\min} = -4$ ,  $l_{\max} = -2$ ,  $ns = 200$ ,  $\eta_{\min} = 1e-5$ ,  $\eta_{\max} = 1e-1$ ,  $n_{\text{cycle}} = 2$ . After searching for lambda, the parameters of the best networks are:

$\lambda = 10^{(-2.826380)}$

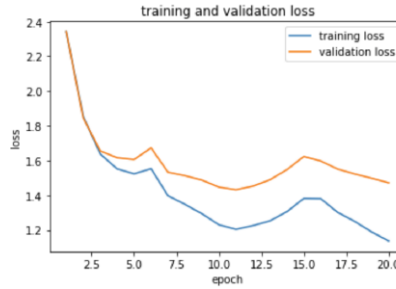
## 5 Best setting training

The parameters are defined as below:

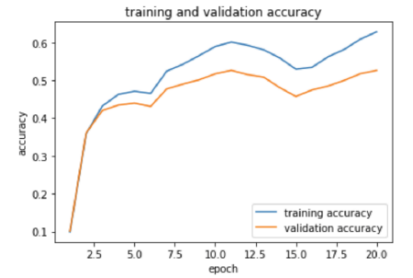
$\lambda = 10^{(-2.826380)}$ ,  $ns = 500$ ,  $n_{\text{epochs}} = 20$ ,  $n_{\text{batch}} = 100$ ,  $\eta_{\min} = 1e-5$ ,  $\eta_{\max} = 1e-1$ .



(a) Learning rate



(b) Training and validation loss



(c) Training and validation accuracy

Figure 5: One cycle training

The final accuracy on test data could be 51.43%.