



國立陽明交通大學

NATIONAL YANG MING CHIAO TUNG UNIVERSITY

Deep Learning

Lab1: back-propagation

Name: 許子駿

Student ID: 311551166

Institute of Computer Science and Engineering

2023 年 3 月 18 日

Table of Contents

1	Introduction	2
2	Experiment setup	3
2.1	Initialization method	3
2.2	Activation functions	3
2.3	Optimizers	5
2.4	Loss function	5
2.5	Data structure of linear layer	6
3	Analysis and comments	7
3.1	Different optimizers	7
3.2	Different activation functions	10
3.3	Different learning rates	12
3.4	Different dimension of hidden layers	14
4	Conclusion	16

Chapter 1

Introduction

This task is to build a deep neural network with 2 hidden layers to classify 2 datasets which each has 2 classes (See Figure 1.1).

One dataset can be approximately separated by an $y = x$ linear classifier. The other dataset is like an X, the data points on the $y = x$ line belong to the same class, and data points on the $xy = 1$ line belong to the other class (See Figure 1.2).

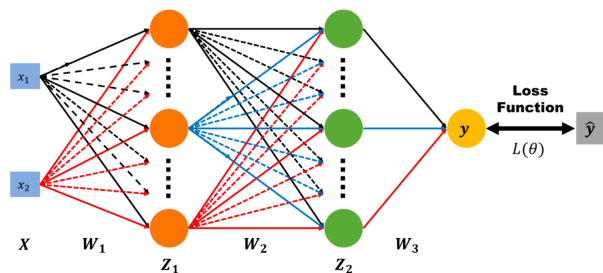


Figure 1.1: Illustration of deep neural network with 2 hidden layers.

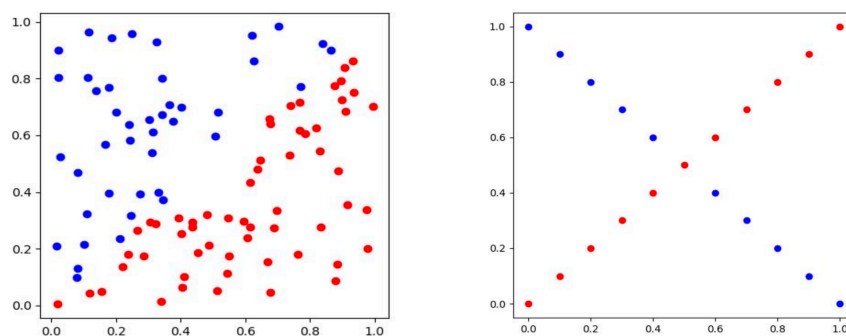


Figure 1.2: Illustration of two datasets.

Chapter 2

Experiment setup

Experiment setup is listed in this chapter including 5 parts, that is, initialization method, activation functions, optimizers, loss function, and the data structure of linear layer.

2.1 Initialization method

He normal initialization is used in this experiment (See Listing 2.1). This initialization method performs well empirically when neural networks include [ReLU](#) activation function by many researchers.

```
1 dim_in = layers[i]
2 dim_out = layers[i + 1]
3 a = np.sqrt(6. / (dim_in + dim_out))
4 std = np.sqrt(2. / ((1 + a ** 2) * dim_in))
5
6 W = np.random.normal(loc=0., scale=std, size=[dim_out, dim_in])
7 b = np.random.normal(loc=0., scale=std, size=[1, dim_out])
```

Listing 2.1: He normal initialization code.

2.2 Activation functions

There are many activation functions, e.g., [LeakyReLU](#), [ReLU](#), [Sigmoid](#), and etc. In this experiment, 3 activation functions, [LeakyReLU](#), [ReLU](#) and [Sigmoid](#) (See Figures 2.1, 2.2, and 2.3), are implemented.

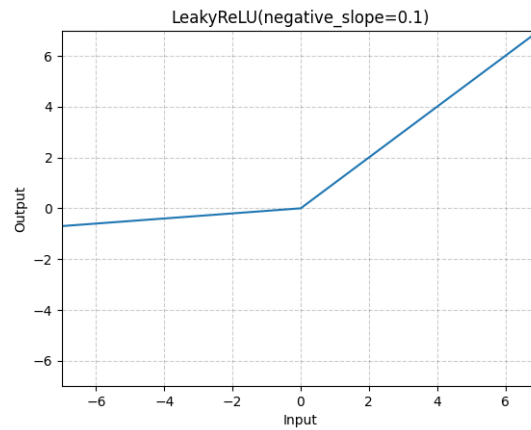


Figure 2.1: Illustration of LeakyReLU activation function.

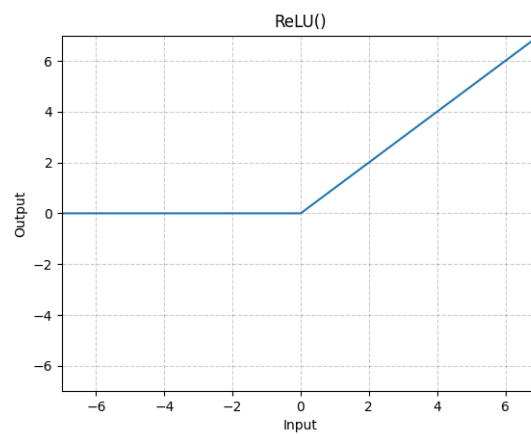


Figure 2.2: Illustration of ReLU activation function.

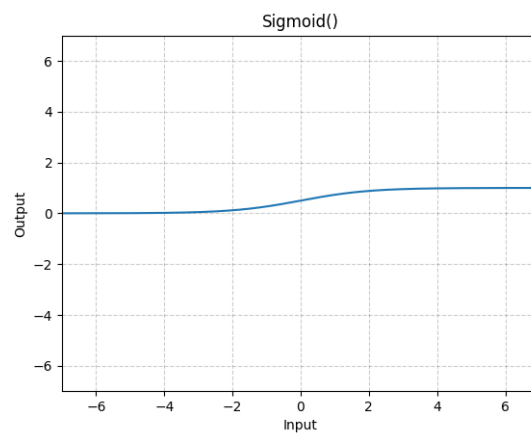


Figure 2.3: Illustration of Sigmoid activation function.

2.3 Optimizers

There are many optimizers, e.g., [SGD](#), [Momentum](#), [AdaGrad](#), [Adam](#), and etc. In this experiment, 2 optimizers, [SGD](#) (vanilla gradient descent in this case) and [Adam](#) (See Figure 2.4) optimizers, are implemented.

[Adam](#) optimizer is mostly used, and it has the great parts of [Momentum](#) and [AdaGrad](#).

```

input :  $\gamma$  (lr),  $\beta_1, \beta_2$  (betas),  $\theta_0$  (params),  $f(\theta)$  (objective)
          $\lambda$  (weight decay), amsgrad, maximize
initialize :  $m_0 \leftarrow 0$  ( first moment),  $v_0 \leftarrow 0$  (second moment),  $\widehat{v}_0^{max} \leftarrow 0$ 

```

```

for  $t = 1$  to ... do
  if maximize :
     $g_t \leftarrow -\nabla_{\theta} f_t(\theta_{t-1})$ 
  else
     $g_t \leftarrow \nabla_{\theta} f_t(\theta_{t-1})$ 
  if  $\lambda \neq 0$ 
     $g_t \leftarrow g_t + \lambda \theta_{t-1}$ 
   $m_t \leftarrow \beta_1 m_{t-1} + (1 - \beta_1) g_t$ 
   $v_t \leftarrow \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$ 
   $\widehat{m}_t \leftarrow m_t / (1 - \beta_1^t)$ 
   $\widehat{v}_t \leftarrow v_t / (1 - \beta_2^t)$ 
  if amsgrad
     $\widehat{v}_t^{max} \leftarrow \max(\widehat{v}_t^{max}, \widehat{v}_t)$ 
     $\theta_t \leftarrow \theta_{t-1} - \gamma \widehat{m}_t / (\sqrt{\widehat{v}_t^{max}} + \epsilon)$ 
  else
     $\theta_t \leftarrow \theta_{t-1} - \gamma \widehat{m}_t / (\sqrt{\widehat{v}_t} + \epsilon)$ 

```

```

return  $\theta_t$ 

```

Figure 2.4: Pseudo code of Adam optimizer.

2.4 Loss function

Cross entropy (See Equation 2.1) is used in this experiment, and it is also mostly used for classification tasks.

$$L(y, \hat{y}) = -[y \log \hat{y} + (1 - y) \log(1 - \hat{y})] \quad (2.1)$$

where y are ground truth labels and \hat{y} are predicted labels.

2.5 Data structure of linear layer

Data structure of linear layer is implemented using a class with only two member variables. It is well-defined, since W (weight matrix) and b (bias vector) are almost used at the same situations.

```
1  class Linear(object):
2      def __init__(self, W: np.ndarray, b: np.ndarray) -> None:
3          self._W = W
4          self._b = b
5
6      @property
7      def W(self) -> np.ndarray:
8          return self._W
9
10     @property
11     def b(self) -> np.ndarray:
12         return self._b
13
14     @W.setter
15     def W(self, W: np.ndarray) -> None:
16         self._W = W
17
18     @b.setter
19     def b(self, b: np.ndarray) -> None:
20         self._b = b
21
22     def __str__(self) -> str:
23         return f'<Linear W: {self._W.shape} b: {self._b.shape}>'
```

Listing 2.2: Data structure code of linear layer.

Chapter 3

Analysis and comments

Different optimizers, activation functions, learning rates, and dimension of hidden layers are tested in this chapter. Experiments of this chapter are done with common setup below:

- Seed: 42.
- Number of hidden layers: 2.

3.1 Different optimizers

Experiments of this section are done with common setup below:

- Activation function: [LeakyReLU](#).
- Learning rate: 10^{-2} .
- Dimension of hidden layers: 4.

From Figures 3.1 and 3.2, we can see that using [Adam](#) optimizer makes the model converge faster in this case.

```
(dl) root@5b48b3907782:/workspace/lab1# python main.py --is_writer_logging True --is_image_saving True --lr 1e-2 --optim adam --activation leaky_relu --dataset xor
0%|          | 0/30000 [00:00<?, 71it/s]
[INFO]: Accuracy reaches 1.0. Early stop at epoch 100. Final prediction:
0%|          | 99/30000 [00:00<02:04, 240.68it/s]
(dl) root@5b48b3907782:/workspace/lab1# python main.py --is_writer_logging True --is_image_saving True --lr 1e-2 --optim adam --activation leaky_relu --dataset linear
0%|          | 144/30000 [00:00<00:20, 1434.93it/s]
[INFO]: Accuracy reaches 1.0. Early stop at epoch 200. Final prediction:
1%|          | 199/30000 [00:00<01:53, 263.64it/s]
```

Figure 3.1: Results with Adam optimizer. In this case, the model converges in less than 100 and 200 epochs, respectively.


```

(dl) root@5b48b3907782:/workspace/lab1# python main.py --is_writer_logging True --is_image_saving True --lr 1e-2 --optim vanilla --activation leaky_relu --dataset xor
2% | 492/30000 [00:00<00:23, 1277.15it/s]
Epoch: 500, loss: 0.688860, acc: 0.523810.
3% |
Epoch: 1000, loss: 0.688079, acc: 0.523810.
5% |
Epoch: 1500, loss: 0.687263, acc: 0.523810.
7% |
Epoch: 2000, loss: 0.686030, acc: 0.523810.
8% |
Epoch: 2500, loss: 0.683972, acc: 0.523810.
10% |
Epoch: 3000, loss: 0.680284, acc: 0.523810.
12% |
Epoch: 3500, loss: 0.674193, acc: 0.523810.
13% |
Epoch: 4000, loss: 0.661351, acc: 0.523810.
14% |
Epoch: 4500, loss: 0.632023, acc: 0.666667.
16% |
Epoch: 5000, loss: 0.569772, acc: 0.809524.
18% |
Epoch: 5500, loss: 0.495491, acc: 0.984762.
19% |
[INFO]: Accuracy reaches 1.0. Early stop at epoch 5800. Final prediction:
19% |
(dl) root@5b48b3907782:/workspace/lab1# python main.py --is_writer_logging True --is_image_saving True --lr 1e-2 --optim vanilla --activation leaky_relu --dataset linears
1% | 346/30000 [00:00<00:21, 1392.44it/s]
Epoch: 500, loss: 0.659871, acc: 0.590000.
3% |
Epoch: 1000, loss: 0.628831, acc: 0.590000.
5% |
Epoch: 1500, loss: 0.517747, acc: 0.780000.
6% |
Epoch: 2000, loss: 0.376159, acc: 0.910000.
8% |
Epoch: 2500, loss: 0.269853, acc: 0.940000.
10% |
Epoch: 3000, loss: 0.207850, acc: 0.950000.
11% |
Epoch: 3500, loss: 0.168794, acc: 0.960000.
13% |
Epoch: 4000, loss: 0.141918, acc: 0.970000.
15% |
Epoch: 4500, loss: 0.122475, acc: 0.970000.
16% |
Epoch: 5000, loss: 0.107567, acc: 0.980000.
18% |
Epoch: 5500, loss: 0.095971, acc: 0.990000.
19% |
[INFO]: Accuracy reaches 1.0. Early stop at epoch 5600. Final prediction:
19% |

```

Figure 3.2: Results with SGD optimizer. In this case, the model converges in less than 5800 and 5600 epochs, respectively.

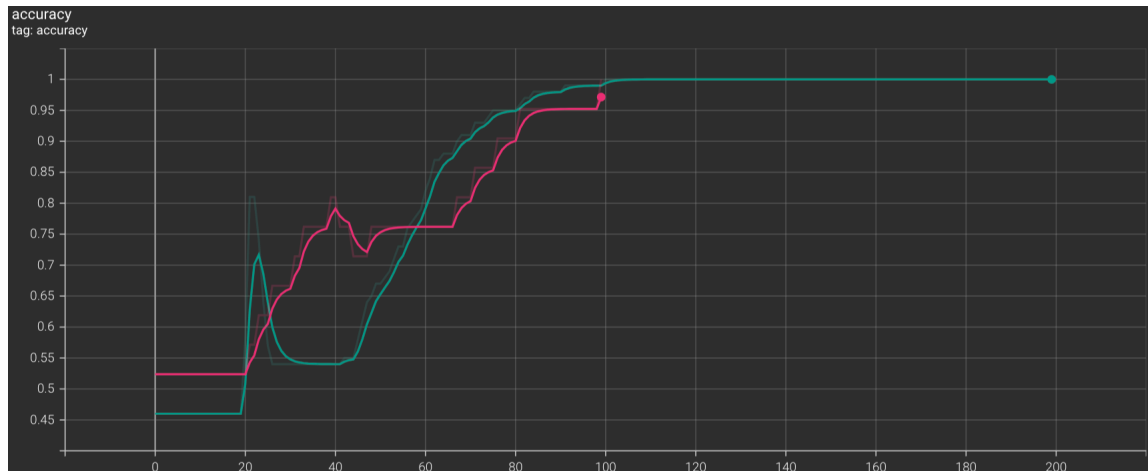


Figure 3.3: Accuracy with Adam optimizer.

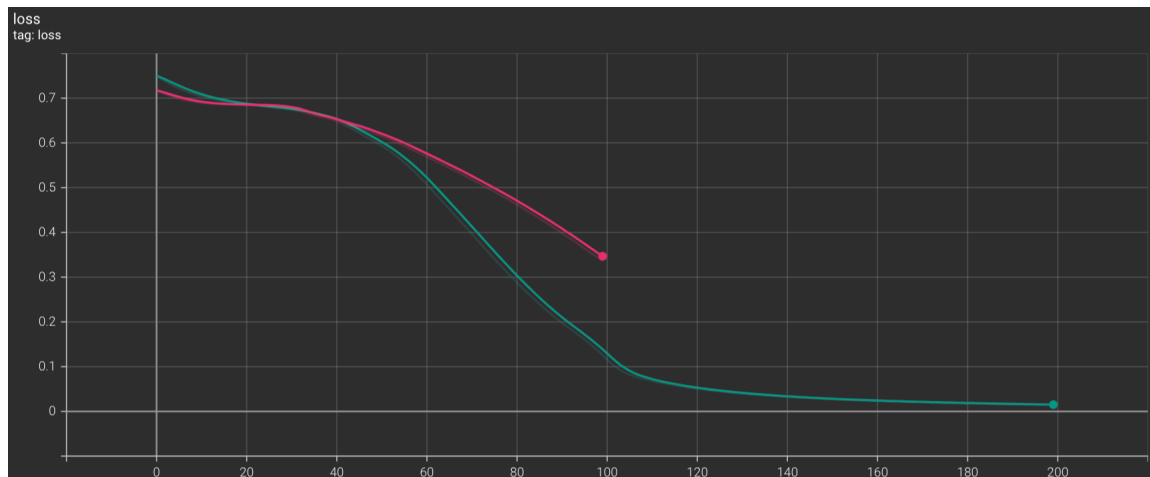


Figure 3.4: Loss with Adam optimizer.

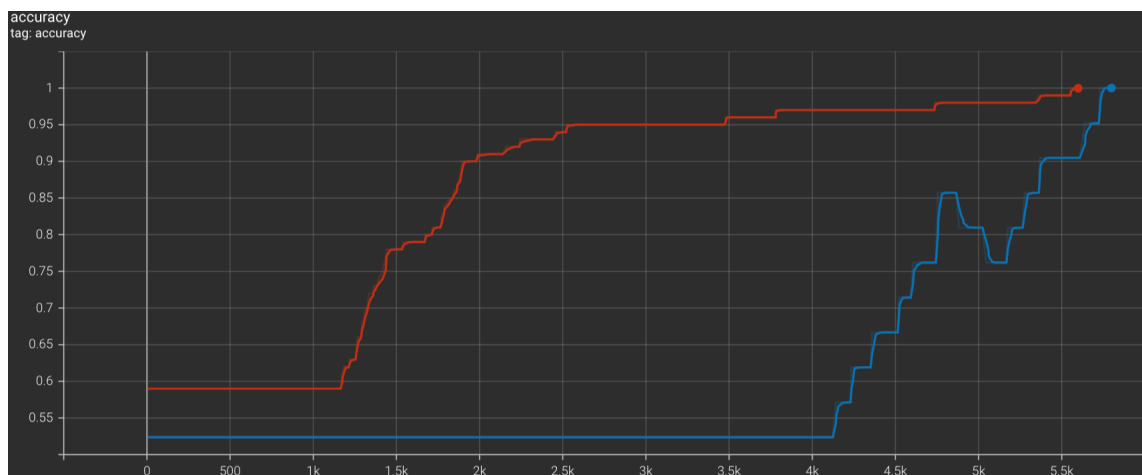


Figure 3.5: Accuracy with SGD optimizer.

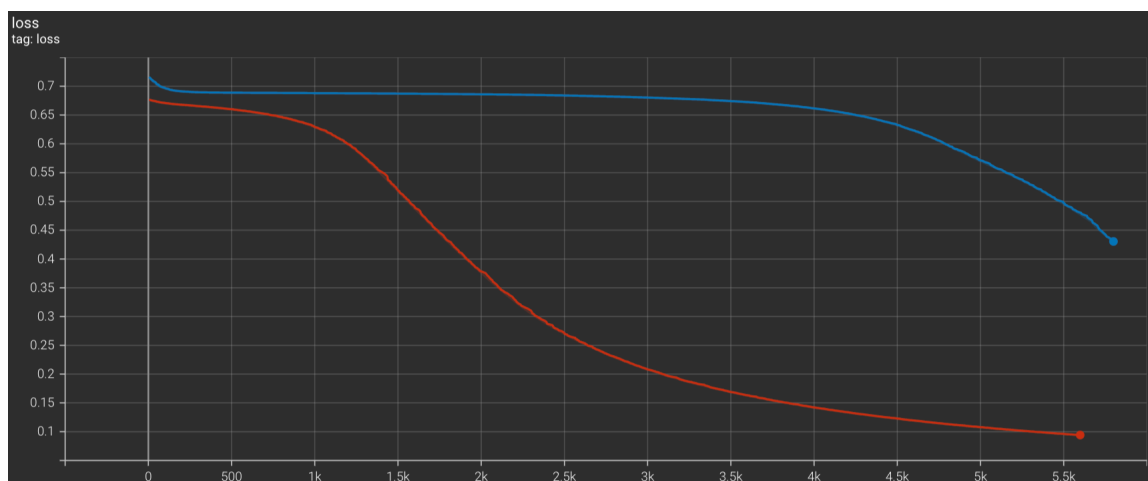


Figure 3.6: Loss with SGD optimizer.

3.2 Different activation functions

Experiments of this section are done with common setup below:

- Optimizer: [Adam](#).
- Learning rate: 10^{-2} .
- Dimension of hidden layers: 4.

From Figures 3.7 and 3.8, we can see that using [LeakyReLU](#) optimizer makes the model converge faster in this case.

```
(d1) root@5b48b3907782:/workspace/lab1# python main.py --is_writer_logging True --is_image_saving True --lr 1e-2 --optim adam --activation leaky_relu --dataset xor
0% | 0/30000 [00:00<7, 71it/s]
[INFO]: Accuracy reaches 1.0. Early stop at epoch 100. Final prediction:
0% | 99/30000 [00:00<02:04, 240.68it/s]
(d1) root@5b48b3907782:/workspace/lab1# python main.py --is_writer_logging True --is_image_saving True --lr 1e-2 --optim adam --activation leaky_relu --dataset linear
0% | 144/30000 [00:00<00:20, 1454.93it/s]
[INFO]: Accuracy reaches 1.0. Early stop at epoch 200. Final prediction:
1% | 199/30000 [00:00<01:53, 263.64it/s]
```

Figure 3.7: Results with LeakyReLU. In this case, the model converges in less than 100 and 200 epochs, respectively.

```
(d1) root@5b48b3907782:/workspace/lab1# python main.py --is_writer_logging True --is_image_saving True --lr 1e-2 --optim adam --activation relu --dataset xor
1% | 175/30000 [00:00<00:17, 1746.43it/s]
[INFO]: Accuracy reaches 1.0. Early stop at epoch 200. Final prediction:
1% | 199/30000 [00:00<01:27, 340.51it/s]
(d1) root@5b48b3907782:/workspace/lab1# python main.py --is_writer_logging True --is_image_saving True --lr 1e-2 --optim adam --activation relu --dataset linear
1% | 340/30000 [00:00<00:18, 1619.33it/s]
[INFO]: Accuracy reaches 1.0. Early stop at epoch 400. Final prediction:
1% | 399/30000 [00:00<00:55, 530.30it/s]
```

Figure 3.8: Results with ReLU. In this case, the model converges in less than 200 and 400 epochs, respectively.

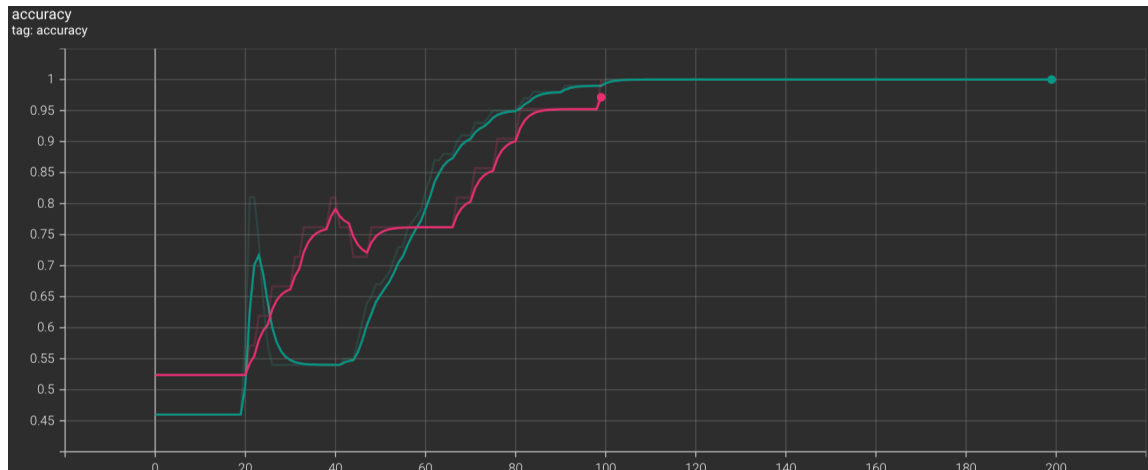


Figure 3.9: Accuracy with LeakyReLU.

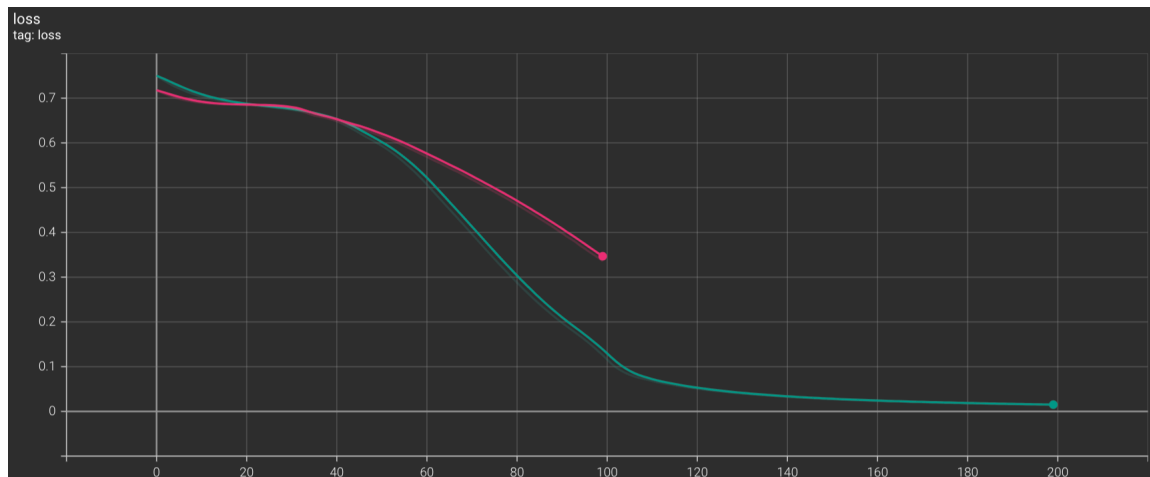


Figure 3.10: Loss with LeakyReLU.

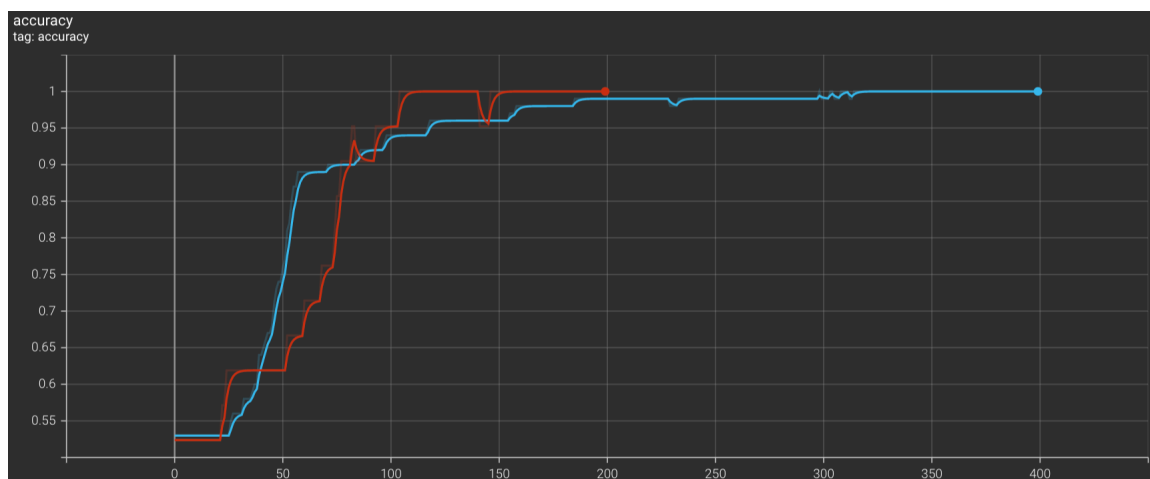


Figure 3.11: Accuracy with ReLU.

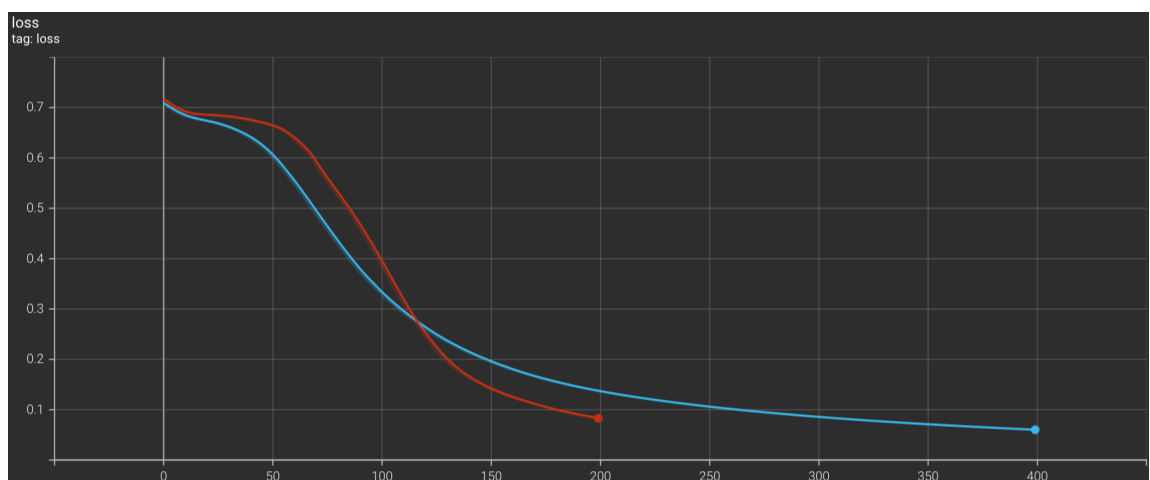


Figure 3.12: Loss with ReLU.

3.3 Different learning rates

Experiments of this section are done with common setup below:

- Optimizer: [Adam](#).
- Activation function: [LeakyReLU](#).
- Dimension of hidden layers: 4.

From Figures 3.13 and 3.14, we can see that using learning rate 10^{-2} makes the model converge faster in this case.

```
(dl) root@5b48b3907782:/workspace/lab1# python main.py --is_writer_logging True --is_image_saving True --lr 1e-2 --optim adam --activation leaky_relu --dataset xor
0% | 0/30000 [00:00<?, ?it/s]
[INFO]: Accuracy reaches 1.0. Early stop at epoch 100. Final prediction:
0% | 99/30000 [00:00<02:04, 240.68it/s]
(dl) root@5b48b3907782:/workspace/lab1# python main.py --is_writer_logging True --is_image_saving True --lr 1e-2 --optim adam --activation leaky_relu --dataset linear
0% | 144/30000 [00:00<00:20, 1434.93it/s]
[INFO]: Accuracy reaches 1.0. Early stop at epoch 200. Final prediction:
1% | 199/30000 [00:00<01:53, 263.64it/s]
```

Figure 3.13: Results with learning rate 10^{-2} . In this case, the model converges in less than 100 and 200 epochs, respectively.

```
(dl) root@5b48b3907782:/workspace/lab1# python main.py --is_writer_logging True --is_image_saving True --lr 1e-4 --optim adam --activation leaky_relu --dataset xor
7% | 1982/30000 [00:01<00:27, 1004.15it/s]
Epoch: 2000, loss: 0.680860, acc: 0.523810.
13% | 3952/30000 [00:03<00:24, 1066.33it/s]
Epoch: 4000, loss: 0.577232, acc: 0.761905.
20% | 5913/30000 [00:05<00:26, 917.96it/s]
Epoch: 6000, loss: 0.484097, acc: 0.952381.
26% | 7918/30000 [00:07<00:18, 1197.48it/s]
Epoch: 8000, loss: 0.365170, acc: 0.952381.
29% | 8609/30000 [00:08<00:21, 1014.14it/s]
[INFO]: Accuracy reaches 1.0. Early stop at epoch 8700. Final prediction:
29% | 8699/30000 [00:08<00:21, 983.20it/s]
(dl) root@5b48b3907782:/workspace/lab1# python main.py --is_writer_logging True --is_image_saving True --lr 1e-4 --optim adam --activation leaky_relu --dataset linear
6% | 1939/30000 [00:01<00:29, 966.26it/s]
Epoch: 2000, loss: 0.657030, acc: 0.550000.
13% | 3933/30000 [00:03<00:24, 1050.34it/s]
Epoch: 4000, loss: 0.443382, acc: 0.950000.
20% | 5966/30000 [00:05<00:19, 1213.01it/s]
Epoch: 6000, loss: 0.260349, acc: 0.970000.
26% | 7933/30000 [00:07<00:24, 901.06it/s]
Epoch: 8000, loss: 0.155702, acc: 0.970000.
33% | 9932/30000 [00:09<00:18, 1057.95it/s]
Epoch: 10000, loss: 0.092129, acc: 0.990000.
34% | 10140/30000 [00:10<00:20, 951.93it/s]
[INFO]: Accuracy reaches 1.0. Early stop at epoch 10200. Final prediction:
34% | 10199/30000 [00:10<00:21, 938.73it/s]
```

Figure 3.14: Results with learning rate 10^{-4} . In this case, the model converges in less than 8700 and 10200 epochs, respectively.

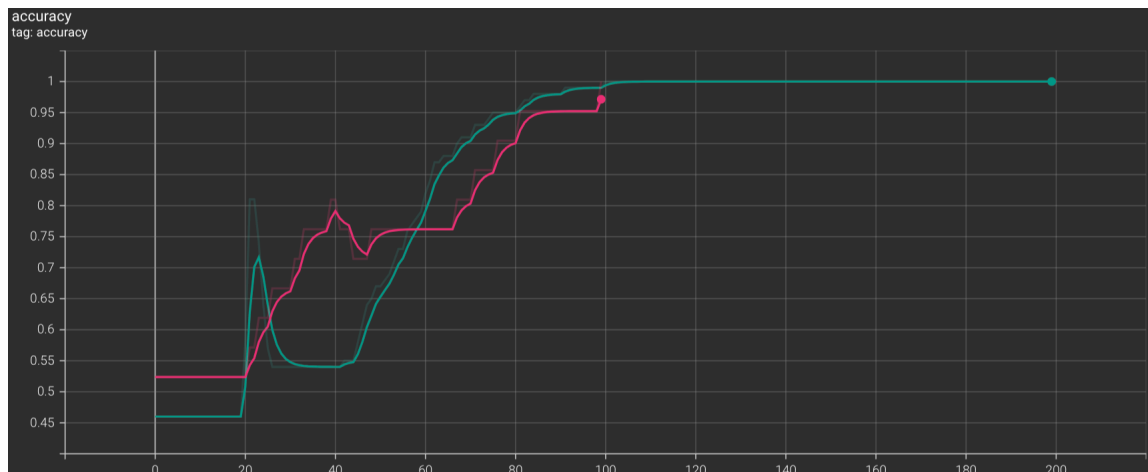
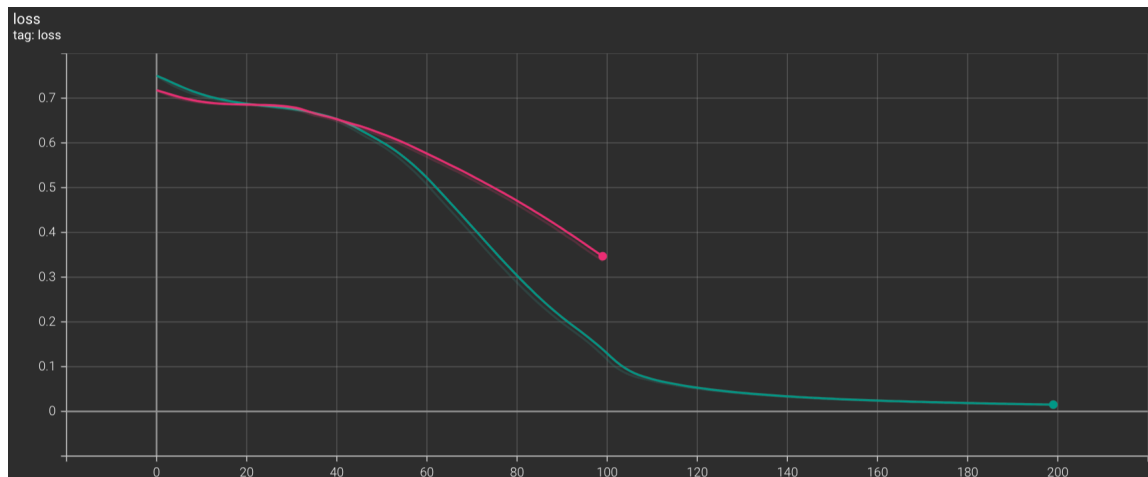
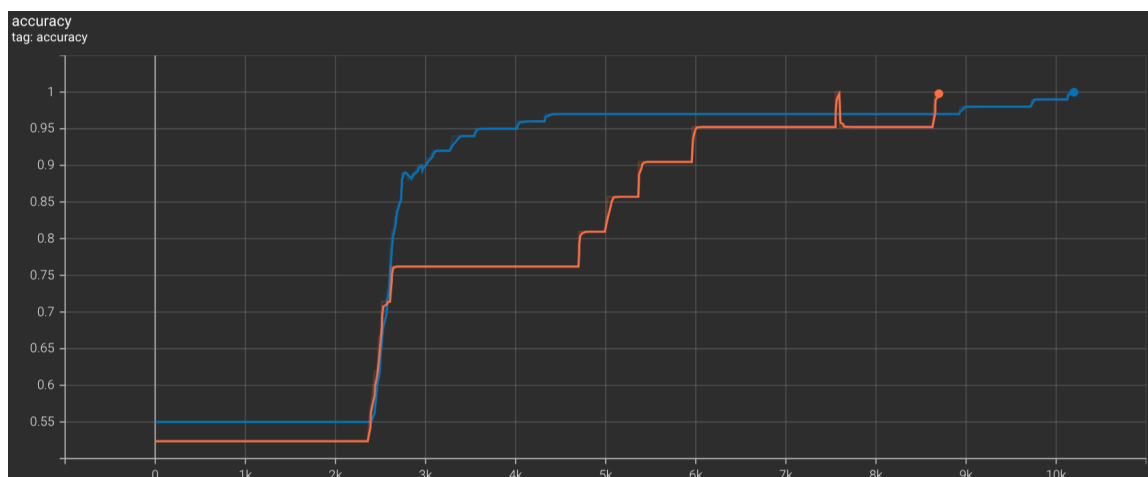
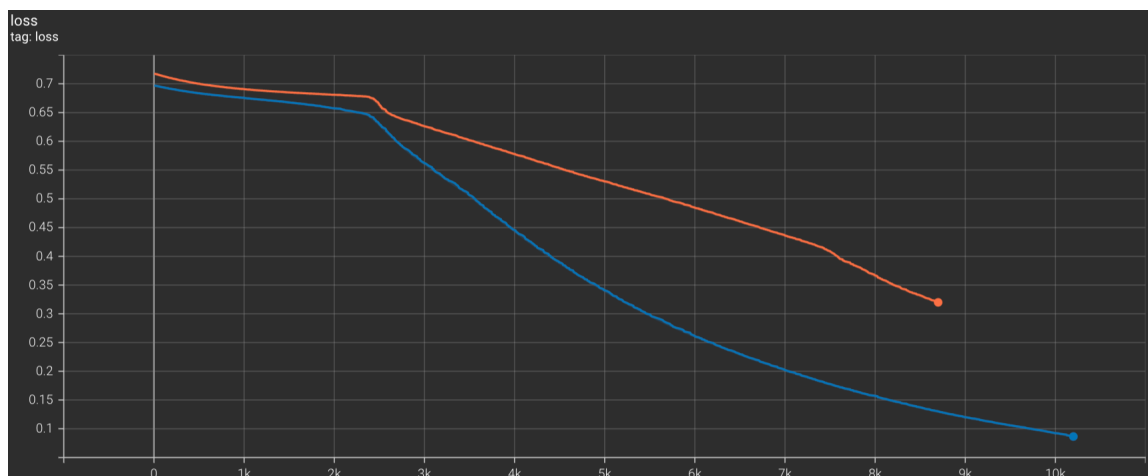


Figure 3.15: Accuracy with learning rate 10^{-2} .

Figure 3.16: Loss with learning rate 10^{-2} .Figure 3.17: Accuracy with learning rate 10^{-4} .Figure 3.18: Loss with learning rate 10^{-4} .

3.4 Different dimension of hidden layers

Experiments of this section are done with common setup below:

- Optimizer: [Adam](#).
- Activation function: [LeakyReLU](#).
- Learning rate: 10^{-2} .

From Figures 3.19 and 3.20, we can see that using 16-dimension hidden layers makes the model converge faster in this case.

```
(d1) root@5b48b3907782:/workspace/lab1# python main.py --is_writer_logging True --is_image_saving True --lr 1e-2 --optim adam --activation leaky_relu --dataset xor
0% | 0/30000 [00:00<?, 7it/s]
[INFO]: Accuracy reaches 1.0. Early stop at epoch 100. Final prediction:
0% | 99/30000 [00:00<02:04, 240.68it/s]
(d1) root@5b48b3907782:/workspace/lab1# python main.py --is_writer_logging True --is_image_saving True --lr 1e-2 --optim adam --activation leaky_relu --dataset linear
0% | 144/30000 [00:00<00:20, 1454.93it/s]
[INFO]: Accuracy reaches 1.0. Early stop at epoch 200. Final prediction:
1% | 199/30000 [00:00<01:53, 263.64it/s]
```

Figure 3.19: Results with 4-dimension hidden layers. In this case, the model converges in less than 100 and 200 epochs, respectively.

```
(d1) root@5b48b3907782:/workspace/lab1# python main.py --is_writer_logging True --is_image_saving True --lr 1e-2 --optim adam --activation leaky_relu --dataset xor --dim_hidden 16
0% | 0/30000 [00:00<?, 7it/s]
[INFO]: Accuracy reaches 1.0. Early stop at epoch 100. Final prediction:
0% | 99/30000 [00:00<02:23, 207.72it/s]
(d1) root@5b48b3907782:/workspace/lab1# python main.py --is_writer_logging True --is_image_saving True --lr 1e-2 --optim adam --activation leaky_relu --dataset linear --dim_hidden 16
0% | 0/30000 [00:00<?, 7it/s]
[INFO]: Accuracy reaches 1.0. Early stop at epoch 100. Final prediction:
0% | 99/30000 [00:00<02:49, 176.10it/s]
```

Figure 3.20: Results with 16-dimension hidden layers. In this case, the model converges in less than 100 and 100 epochs, respectively.

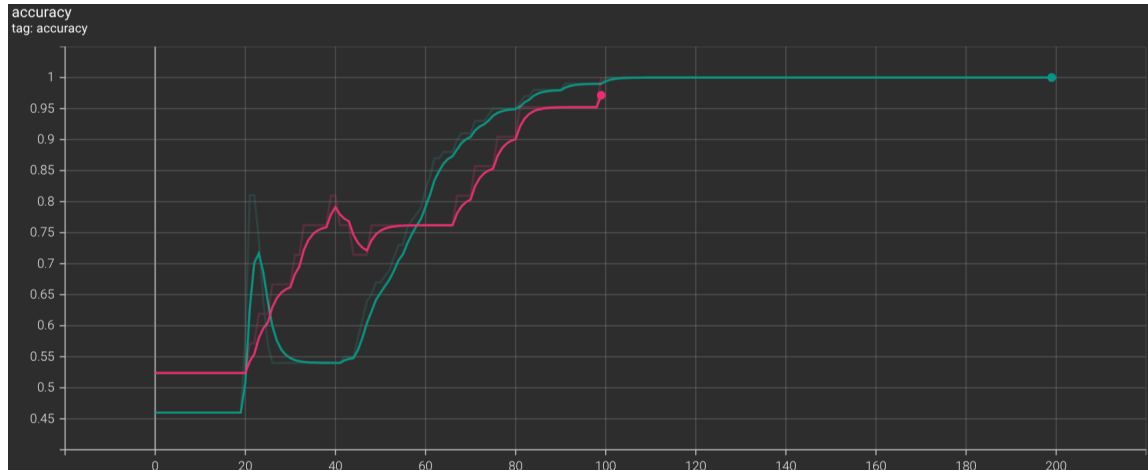


Figure 3.21: Accuracy with 4-dimension hidden layers.

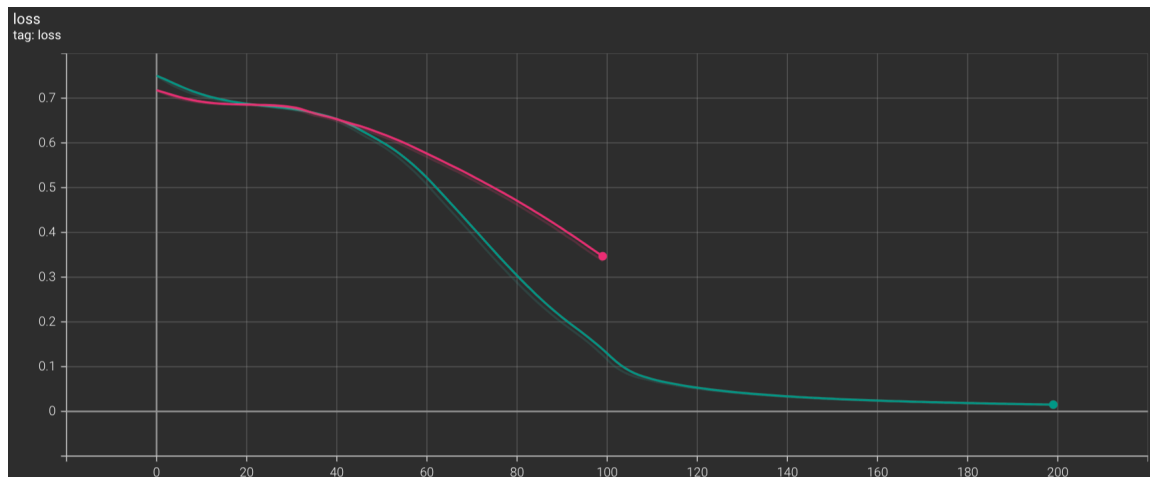


Figure 3.22: Loss with 4-dimension hidden layers.

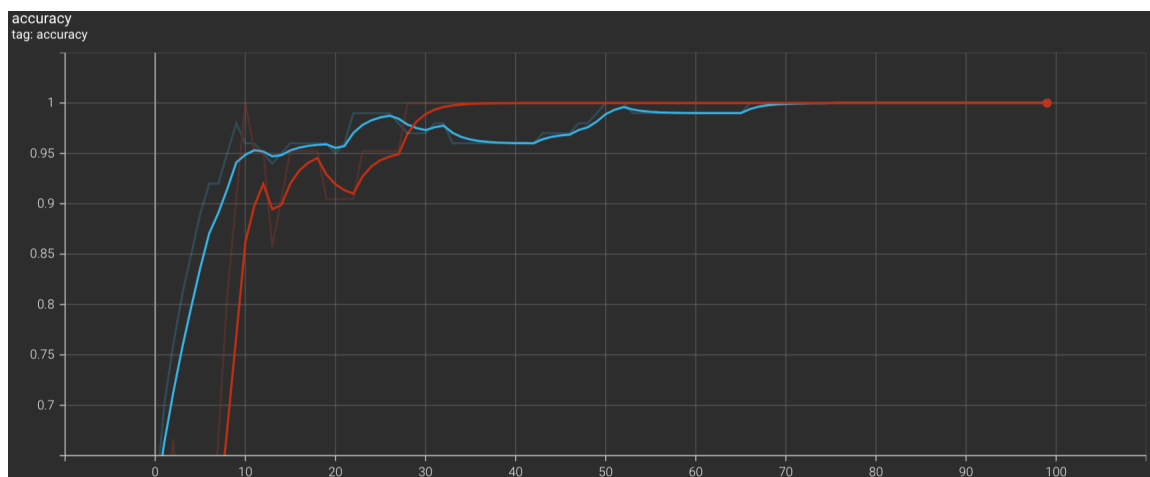


Figure 3.23: Accuracy with 16-dimension hidden layers.

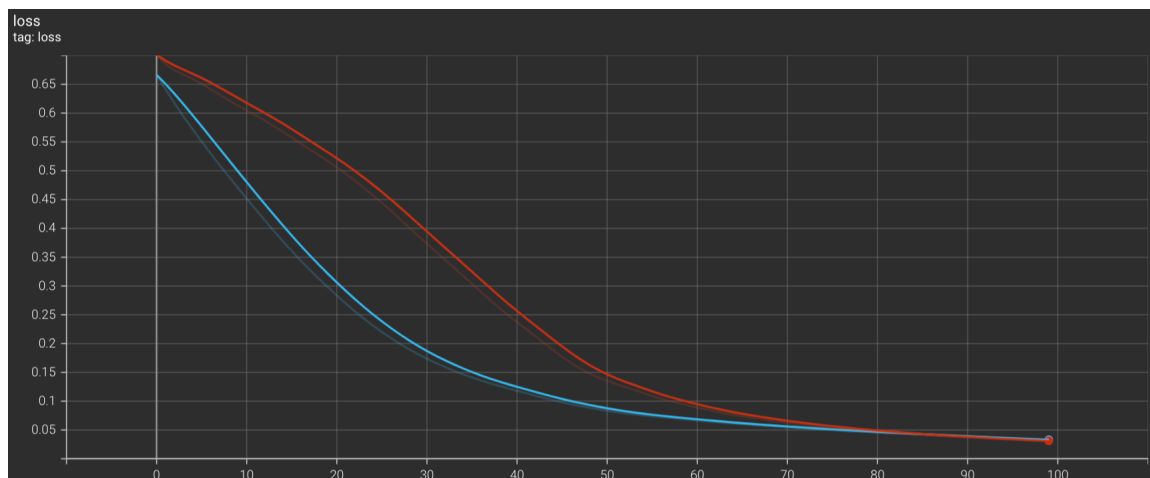


Figure 3.24: Loss with 16-dimension hidden layers.

Chapter 4

Conclusion

After trying so many experiments, I found some patterns for this task:

- [Adam](#) optimizer is **much better** than [SGD](#) optimizer.
- [LeakyReLU](#) is **a little better** than [ReLU](#).
- Learning rate 10^{-2} is great, too large or too small is hard to converge.
- **Larger** dimension of hidden layers converge faster.