

---

# Exeriese 2

---

Lan Zhang  
954517477

## 1 Problem 1

Solved on 09/17/2019. The seed I used for noise generation is 2612.

### 1.1 Loss Function

Given an input  $x$  with an output  $y$ , the prediction is  $\hat{y}$ . The performance of the squared loss function shown in Figure 1 (a).

The formulation of squared loss is:

$$squared\_loss = \frac{1}{N} \sum_0^N (y_i - \hat{y}_i)^2$$

Given the parameter  $\delta = 1.0$ . The formulation of huber loss is:

$$huber\_loss = \begin{cases} 0.5 * (y - \hat{y}) & \text{if } |y - \hat{y}| \leq \delta \\ \delta * |y - \hat{y}| - 0.5 * \delta^2 & \text{if } |y - \hat{y}| > \delta \end{cases}$$

The performance of the Huber loss function shown in Figure 1 (b).

I also implement the following hybrid loss function and the performance shown in Figure 1 (c).

$$hybrid\_loss = huber\_loss + |y - \hat{y}|$$

If I didn't use `tf.reduce_mean` in squared loss, the result of parameters  $W$  and  $b$  will turn out to be NaN. I found the issue will happen on TensorFlow 2.0. But on Tensorflow 1.14, sometimes it can output a correct answer. I'm not sure about the reason. But I think the loss is too large, which may result in NaN after several steps of gradients.

It's hard to tell which loss function works better(see Figure 1 (d)). I've run the code with those two loss functions multiple times and both of them work well. If I use the distance of the predicted parameters  $\hat{W}, \hat{b}$  and original parameters  $W, b$  as a measure of the performance( $distances = |W - \hat{W}| + |b - \hat{b}|$ ), the squared loss(.0133321) work better than Huber loss(.0238144).

### 1.2 Learning rate

I implemented patience scheduling. Then I tested the learning rate from 0.001 to 0.01 with the incrementation 0.002, Shown in Figure 2. Table 1 shows the parameters learned under different learning rate with patience scheduling. The results are quite similar. All of them approximate to the correct parameters.

### 1.3 Training steps

I've tested different training steps from 500 to 1500 with the incrementation 200. Shown in Figure 3. With a longer duration, the performance tends to be better. However, after "enough" steps, this effects start to stabilize (see Table 2).

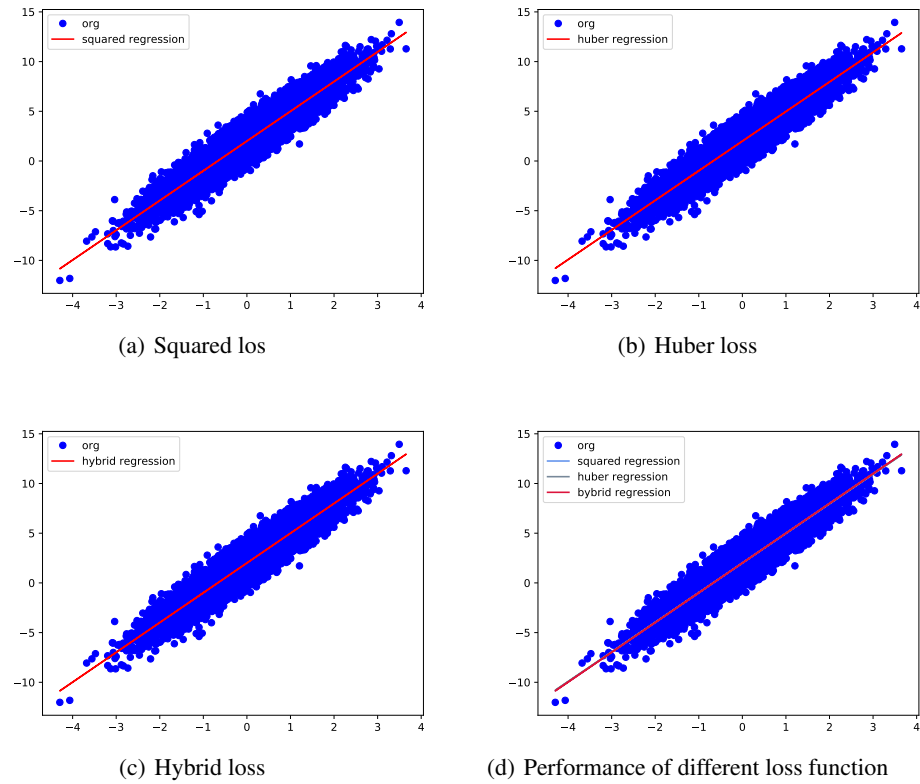


Figure 1: Loss function

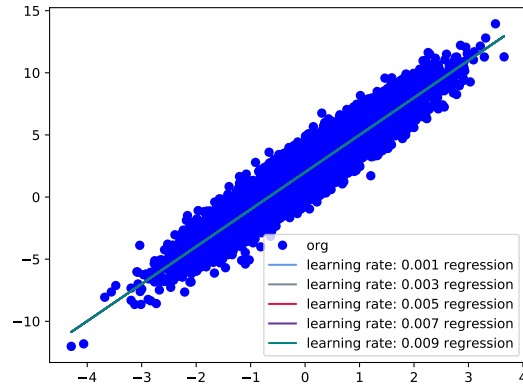


Figure 2: Learning rate

Table 1: The performance of different learning rate

| Learning rate | W         | b         |
|---------------|-----------|-----------|
| 0.001         | 2.9962702 | 2.0017185 |
| 0.003         | 2.9964287 | 2.001716  |
| 0.005         | 2.996503  | 2.0016577 |
| 0.007         | 2.996541  | 2.001543  |
| 0.009         | 2.9966085 | 2.0013998 |

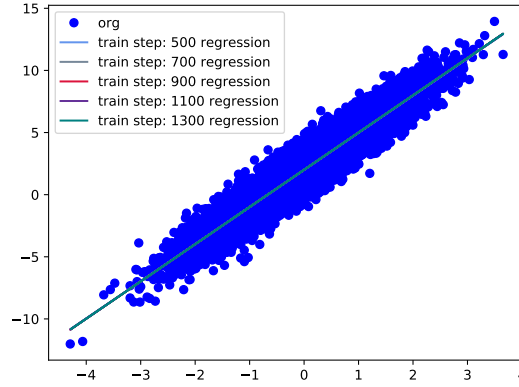


Figure 3: Training steps

Table 2: The performance of different training steps

| Training steps | W         | b         |
|----------------|-----------|-----------|
| 500            | 2.9877446 | 1.9952948 |
| 700            | 2.995941  | 2.001462  |
| 900            | 2.9962356 | 2.0016909 |
| 1100           | 2.9962835 | 2.001727  |
| 1300           | 2.9962983 | 2.0017405 |

#### 1.4 Initial Value of parameters

Table 3 and Table 4 respectively denote predicted parameters with different initial parameter  $W$  and  $b$  (both of them are from -1000 to 1000 with the incrementation 500).

If I initialize the parameter  $W$  with a small number like -1000 with hybrid Loss, the predicted parameter  $W$  will become -984.0698 after 1000 epochs. That is because the model has not yet converged after 1000 epochs. Figure 4 is the performance of different training epoch with the initial parameter  $W = -100$ . Despite the predicted parameters after 9000 epochs are far from the correct answer, it is gradually approaching to the correct answer.

Table 3: The performance of different  $W$  after 1000 epochs

| Initial W | W         | b           |
|-----------|-----------|-------------|
| -1000     | -984.0698 | -0.16686472 |
| -500      | -484.0088 | -0.12897438 |
| 0         | 2.9962702 | 2.0017185   |
| 500       | 484.0088  | 0.2609221   |
| 1000      | 984.0698  | 0.22580884  |

Table 4: The performance of different  $b$  after 1000 epochs

| Initial b | W           | b         |
|-----------|-------------|-----------|
| -1000     | -0.19055721 | -979.8584 |
| -500      | -0.19055721 | -479.8584 |
| 0         | 2.9962702   | 2.0017185 |
| 500       | 0.19055721  | 479.8584  |
| 1000      | 0.19055721  | 979.8584  |

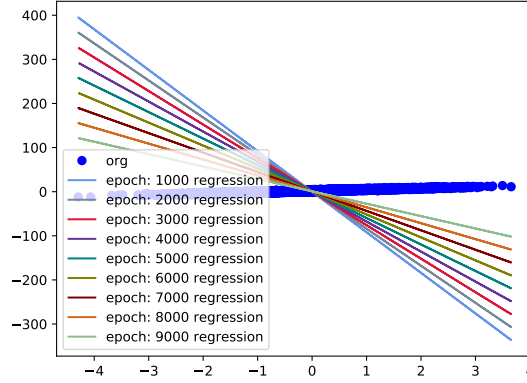


Figure 4: Training duration with a large initial parameter

### 1.5 Noise in data

Before this section, I generated noise in data using the normal distribution. Table 5 shows the results of the different noise distribution. The distribution of noise can affect the prediction. Table 6 shows the results of the normal distribution with different mean. Table 7 shows the results of the normal distribution with different std. The level of noise is higher, the model is harder to converge.

Table 5: The performance of different noise in data

| Noise distribution           | W         | b         |
|------------------------------|-----------|-----------|
| normal(mean=0.0, stddev=1.0) | 2.9794447 | 1.9967409 |
| gamma(alpha=1)               | 2.98501   | 2.8410964 |
| uniform(minval=0)            | 2.9989622 | 2.502729  |

Table 6: The performance of normal distribution with different mean

| mean | W         | b         |
|------|-----------|-----------|
| 0.0  | 2.9794447 | 1.9967409 |
| 1.0  | 2.964156  | 2.990973  |
| 2.0  | 2.9483728 | 3.9745786 |
| 3.0  | 2.924222  | 4.9402742 |

Table 7: The performance of normal distribution with different std

| mean | W         | b         |
|------|-----------|-----------|
| 1.0  | 2.9794447 | 1.9967409 |
| 2.0  | 2.8706882 | 1.9283952 |
| 3.0  | 2.6891909 | 1.8024324 |

### 1.6 Noise in data during training

First, I tested adding noises(normal distribution:  $mean = 0.$ ,  $stddev = 1.$ ) in data per epoch. The predicted parameters are  $\hat{W} = 1.4970689$ ,  $\hat{b} = 1.8999192$ .

Then, I added noises(normal distribution:  $mean = 0.1$ ,  $stddev = 0.1$ ) in weights per epoch. The predicted parameters are  $\hat{W} = -15.267623$ ,  $\hat{b} = 1.1710013$ .

Finally, I added noises(normal distribution:  $mean = 0.001$ ,  $stddev = 0.002$ ) in learning rate per epoch. The predicted parameters are  $\hat{W} = 2.9953077$ ,  $\hat{b} = 2.007313$ .

Noise in data and weight will change the performance of the model. If noise in learning rate is reasonable, it may affect training time. But it will not change the general performance of the model, except it reaches local minima. For example, it may lead to faster convergence or make the model hard to converge.

On other classification problems and mathematical models, noise in data and weight will also reduce the model's performance. But adding small noise in learning rate may make it possible to get rid of local minima.

## 1.7 Significance of seed

If I use the same seed to generate noise, I will get the same results every time I run the code without change any parameters. Because the pseudorandom number generating algorithms is designed for performing an operation on a seed. So the result of algorithms is determined by the seed. If I want to change the result every time, I can use the current time as the seed.

## 1.8 GPU VS CPU

I tested training time with GPU and CPU on Calab. I used each of them to run the same code for 1000 epochs. And the approximate running time of CPU per epoch is 0.029198 and of GPU is 0.05428.

The training time looks unreasonable. I think it may be because it's not the correct running time considering the network speed and the design of Calab. It may also because the linear regression model doesn't fit to calculate on GPU.

# 2 Problem 2

Solved on 09/21/2019.

## 2.1 Performance of three models

I implemented three models for fminist.

The first one is Vgg16(shown as Figure 5). Another one is multilayer CNNs with batch norm. The accuracy of Vgg16 on test dataset can achieve 92.45% after 15 epochs.

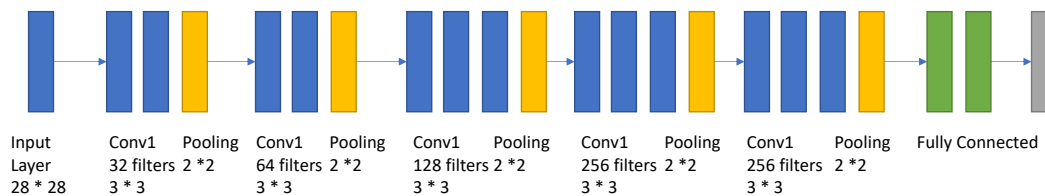


Figure 5: Vgg16

Another model achieves 90.67% after 120 epochs. Figure 6 is an example of images and CNNs model's prediction. Figure 7 is the weights of the first layer in CNNs model.

I also implement linear model as problem 1 in this dataset. It can achieve 82.6% after 120 epochs. Figure 8 is the weights of the first layer in CNNs model.

The following reports are based on multilayer CNNs model not Vgg16, because Vgg16 takes for a while per epoch.

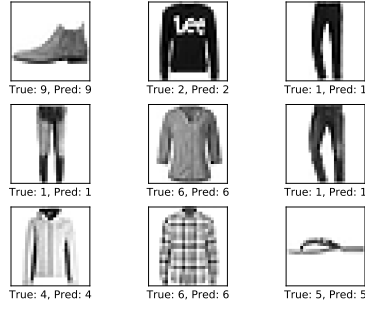


Figure 6: Samples

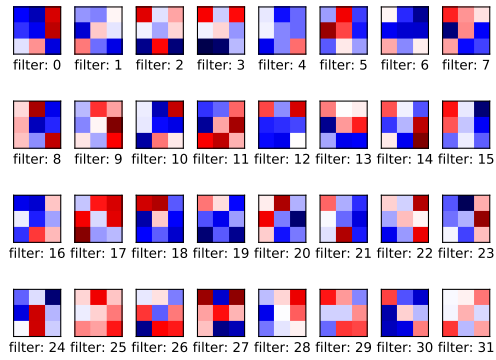


Figure 7: Filters' weights of the first layer in CNNs model

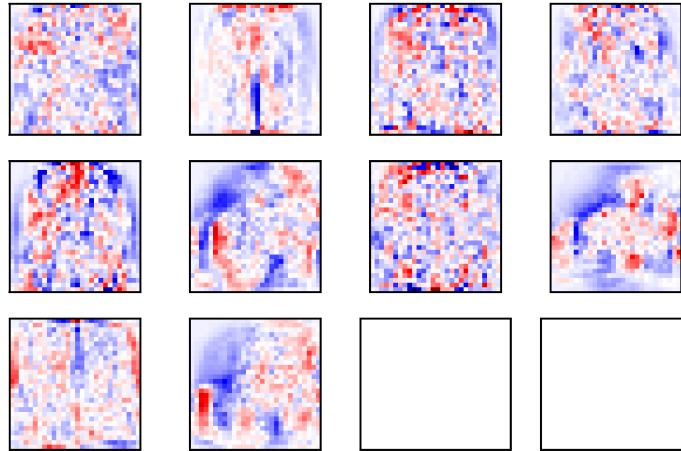


Figure 8: Weights of the first layer in linear model

## 2.2 Optimizer

Figure 9 is the loss and accuracy of train dataset and test dataset using Adam optimizer. Figure 10 is the performance of Adagrad optimizer. Figure 11 is the performance of Nadam optimizer. Adam optimizer could converge faster and generalize better. I tested SGD Optimizer( $momentum = 0.9$ ).

After 8 iterations, the loss will turn to  $NaN$ . Maybe it is the same reason I listed in problem 1. The exploding gradients problem lead to overflow at some point during gradient descent and cause NaNs.

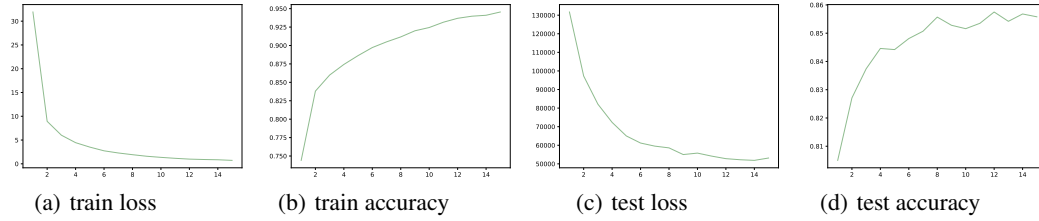


Figure 9: Adam Optimizer

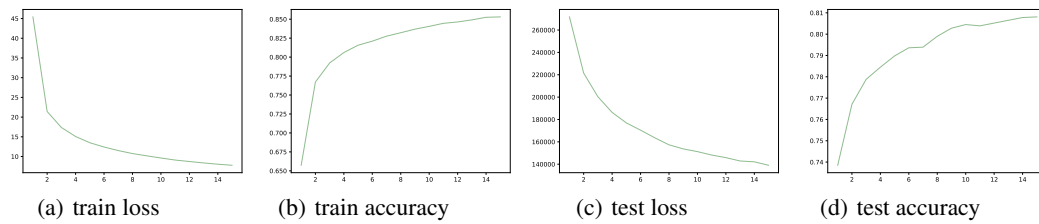


Figure 10: Adagrad Optimizer

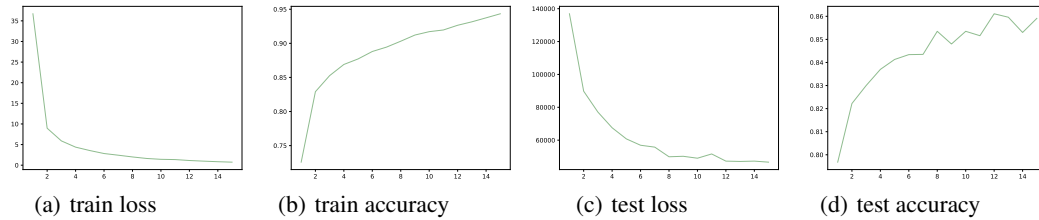


Figure 11: NAdam Optimizer

### 2.3 Train for longer epochs

Figure 12 is the result if I train the model for more epochs. At first, I thought the decreasing accuracy on test dataset means the model is already overfitting. But with longer epochs, it turns out the model just fall into local minima. And the accuracy can reach 90.67% after 120 epochs.

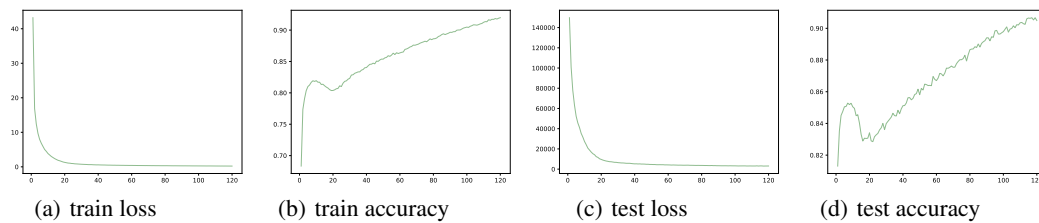


Figure 12: More Training time

## 2.4 Train/Val Split

Figure 13 is the result (best accuracy is 86.57% on test dataset) that the proportion of train dataset to test dataset is 9:1. Before this section, the ratio of train and test dataset is 6:1, and the result (best accuracy on the test dataset is 85.87%) shows in Figure 9. Figure 14 is the result (best accuracy is 86.10% on test dataset) that the proportion of train dataset to test dataset is 7:3. Figure 15 is the result (best accuracy is 84.72% on test dataset) that the proportion of train dataset to test dataset is 6:4.

The model needs enough data to train, so less train dataset has greater variance. But more train dataset can not guarantee a robust measure of error.

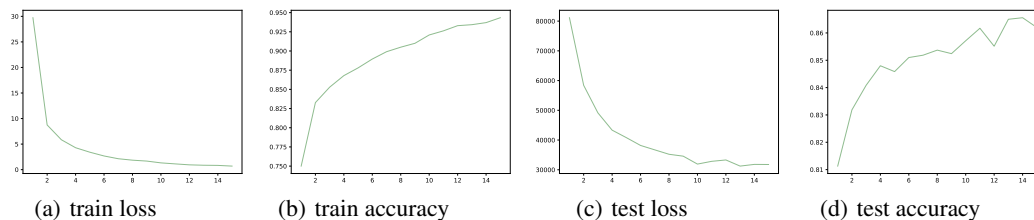


Figure 13: train dataset : test dataset = 9 : 1

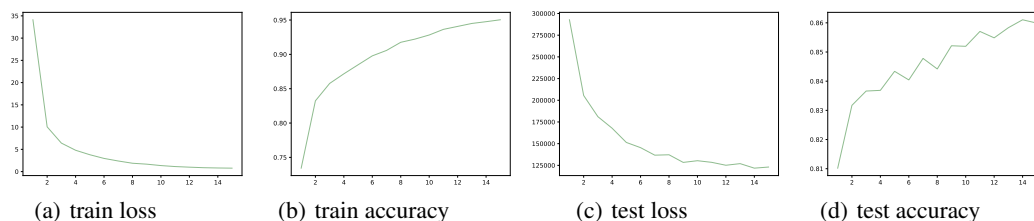


Figure 14: train dataset : test dataset = 7 : 3

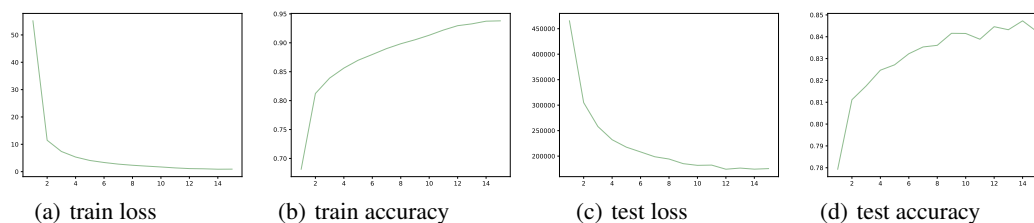


Figure 15: train dataset : test dataset = 6 : 4

## 2.5 Batch size

Figure 16 is the result of batch size 500, which can reach 86.84% on accuracy within 15 epochs. Figure 9 is the result of batch size 1000, whose maximum of accuracy is 85.87%. Figure 17 is the result of batch size 1500, whose maximum of accuracy is 85.20%. Batch size will affect the performance of the model. But a smaller batch size also leads to longer training time. So we need to balance those two factors based on different requirements.

## 2.6 GPU VS CPU

I tested training time with GPU and CPU on Calab. The approximate running time of CPU per epoch is 114.79436 and of GPU is 6.03704.



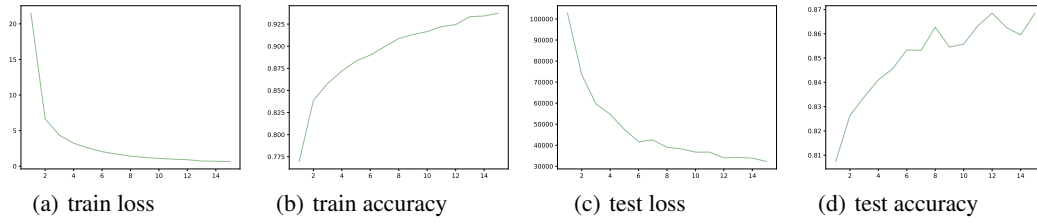


Figure 16: Batch size: 500

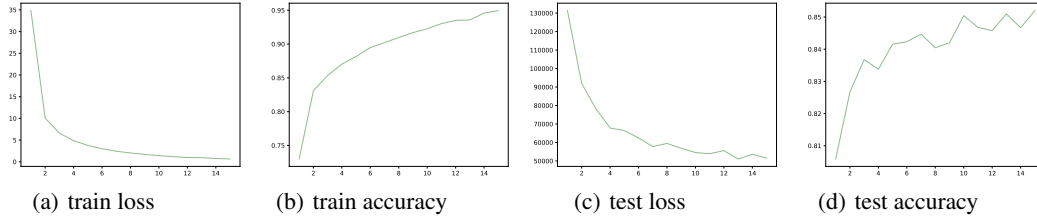


Figure 17: Batch size: 1500

Unlike problem 1, with the model CNNs, the training time on CPU per epoch is almost 20 times than training time on GPU.

## 2.7 Overfit

The CNNs model is not overfitting according to Figure 12. Normally, I'll run a model for a few epochs. If the accuracy(loss) on train dataset is very desired while those measures on the test dataset are still unreasonable, I assume the model overfits. For example, it would be a red signal if the model has 99% accuracy on the training dataset but only 50% accuracy on the test dataset. I'll use dropout and regularization to prevent the model overfit.

## 2.8 Performance compared with other algorithms

I tried random forest classifier( $n\_estimators=100$ ,  $max\_depth=30$ ), which can achieve 87.57% on accuracy. Decision tree algorithm could achieve 79.02%. KNeighbors algorithm could achieve 85.54%.

## 2.9 Cluster on weights

Because there are many weights in CNNs model, I clustered the weights from the linear model.

I tried kmeans and ksne with the constrain that there are two Figure 18 ??

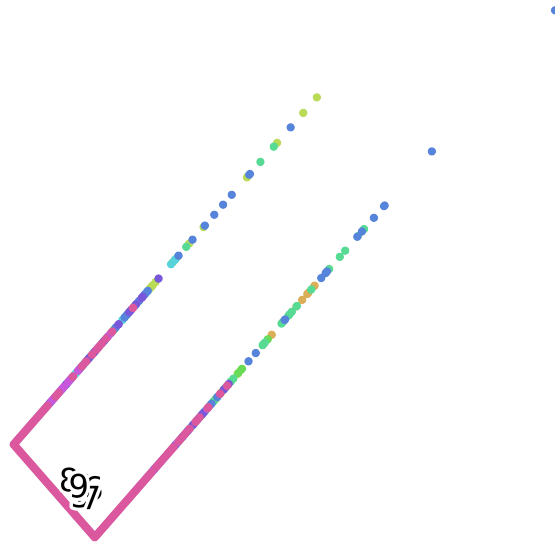


Figure 18: KMEANS

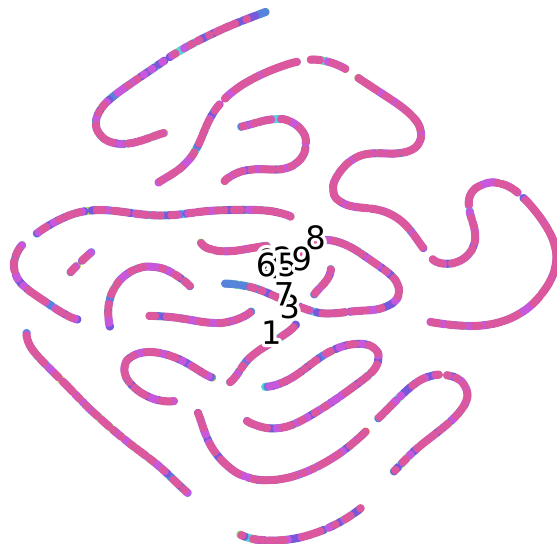


Figure 19: KSNE