

Written Part

Problem 1

(15 points)

Swish is an activation function that can be viewed as a “modification” of the logistic function. Swish has been shown to be better than ReLU in some deep learning models. Recall that the logistic function is defined as

$$\theta(s) = \frac{1}{1 + \exp(-s)}$$

Now Swish is defined as

$$\varphi(s) = s \cdot \theta(s)$$

What is $\varphi'(s)$?

$$\varphi'(s) = \theta(s) + s \cdot \theta'(s) = \frac{1}{1 + e^{-s}} + s \cdot \frac{e^{-s}}{(1 + e^{-s})^2} = \frac{1 + e^{-s} + se^{-s}}{(1 + e^{-s})^2}$$

Problem 2

(20 points)

In the class, we briefly talked about the famous **PageRank** algorithm, which assigns a rank to each web page. Given a set of pages: $P = \{P_1, P_2, \dots, P_n\}$ and their incoming links: $in(P_i) = \{P_j \mid P_j \text{ has a link to } P_i, i \neq j\}$. The basic idea is that a web page should be ranked higher if it has 1) more incoming hyperlinks and 2) incoming links from high-ranked pages. The PageRank can be computed through the following iterative algorithm:

1. Initialize $\text{PageRank}_0(P_i) = \frac{1}{n}$ for each page.
2. Compute the updated PageRank for each page:

$$\text{PageRank}_1(P_i) = \sum_{P_j \in in(P_i)} \frac{\text{PageRank}_0(P_j)}{|out(P_j)|},$$

where $|out(P_j)|$ denotes the out-degree of P_j .

3. Repeat step 2 until the ranks converge.

Now, consider the following example with 3 pages:

$$\begin{aligned} P &= \{P_1, P_2, P_3\} \\ in(P_1) &= \{P_2, P_3\} \\ in(P_2) &= \{P_3\} \\ in(P_3) &= \{P_1\} \end{aligned}$$

We define the following variables:

$$\begin{aligned} \mathbf{v}_t &= [\text{PageRank}_t(P_1), \text{PageRank}_t(P_2), \text{PageRank}_t(P_3)]^T \\ \mathbf{P} &= \begin{bmatrix} 0 & 1 & 0.5 \\ 0 & 0 & 0.5 \\ 1 & 0 & 0 \end{bmatrix}, \end{aligned}$$

where \mathbf{v}_t is the ranks after the t -th iteration, and \mathbf{P} is the transition matrix between pages. We can rewrite the update rule above as $\mathbf{v}_t = \mathbf{P}\mathbf{v}_{t-1}$. Answer the following questions:

- (A) Based on the iterative algorithm above, please compute the values of $\mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3, \mathbf{v}_4$, and \mathbf{v}_5 .
- (B) Recall that in the class, we mentioned that the algorithm converges when $\mathbf{v}^* = \mathbf{P}\mathbf{v}^*$. Solve this equation to derive \mathbf{v}^* . Note that you should provide a normalized \mathbf{v}^* , i.e., $\sum \mathbf{v}^* = 1$, as the answer.

cb.) Let $\mathbf{v}^* = [a, b, c]^T$. Then $\begin{bmatrix} a \\ b \\ c \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0.5 \\ 0 & 0 & 0.5 \\ 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} a \\ b \\ c \end{bmatrix} \Leftrightarrow \begin{cases} a = b + c/2 \\ b = c/2 \\ c = a \end{cases}$

which gives us $a = 2b = c$ after simplification. Since \mathbf{v}^* is normalized, we have $\mathbf{v}^* = [\frac{2}{5}, \frac{1}{5}, \frac{2}{5}]^T$

(a.) At first, we have $\mathbf{v}_0 = [\frac{1}{3}, \frac{1}{3}, \frac{1}{3}]^T$, Therefore

$$\mathbf{v}_1 = \mathbf{P}\mathbf{v}_0 = \begin{bmatrix} 0 & 1 & 0.5 \\ 0 & 0 & 0.5 \\ 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} \frac{1}{3} \\ \frac{1}{3} \\ \frac{1}{3} \end{bmatrix} = \begin{bmatrix} \frac{1}{2} \\ \frac{1}{6} \\ \frac{1}{3} \end{bmatrix},$$

$$\mathbf{v}_2 = \mathbf{P}\mathbf{v}_1 = \begin{bmatrix} 0 & 1 & 0.5 \\ 0 & 0 & 0.5 \\ 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} \frac{1}{2} \\ \frac{1}{6} \\ \frac{1}{3} \end{bmatrix} = \begin{bmatrix} \frac{1}{3} \\ \frac{1}{6} \\ \frac{1}{2} \end{bmatrix},$$

$$\mathbf{v}_3 = \mathbf{P}\mathbf{v}_2 = \begin{bmatrix} 0 & 1 & 0.5 \\ 0 & 0 & 0.5 \\ 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} \frac{1}{3} \\ \frac{1}{6} \\ \frac{1}{2} \end{bmatrix} = \begin{bmatrix} \frac{5}{12} \\ \frac{1}{4} \\ \frac{1}{3} \end{bmatrix},$$

$$\mathbf{v}_4 = \mathbf{P}\mathbf{v}_3 = \begin{bmatrix} 0 & 1 & 0.5 \\ 0 & 0 & 0.5 \\ 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} \frac{5}{12} \\ \frac{1}{4} \\ \frac{1}{3} \end{bmatrix} = \begin{bmatrix} \frac{9}{12} \\ \frac{1}{6} \\ \frac{5}{12} \end{bmatrix},$$

$$\mathbf{v}_5 = \mathbf{P}\mathbf{v}_4 = \begin{bmatrix} 0 & 1 & 0.5 \\ 0 & 0 & 0.5 \\ 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} \frac{9}{12} \\ \frac{1}{6} \\ \frac{5}{12} \end{bmatrix} = \begin{bmatrix} \frac{7}{8} \\ \frac{5}{24} \\ \frac{9}{12} \end{bmatrix} *$$

Problem 3

(20 points)

We talked about the K-means clustering algorithm in class, which iteratively does partition optimization and prototype optimization until convergence. Given a set of 2-dimensional data points $\mathbf{X} = \{(1, 2), (3, 4), (7, 0), (10, 2)\}$, we would like to perform K-means clustering with $K = 2$. Answer the following questions. Note that you should write down the resulting cluster centroids and partitions of each iteration.

- (A) Perform the K-means algorithm with $K = 2$ and initial centroids $\{\mu_1, \mu_2\} = \{(1, 2), (3, 4)\}$ until convergence.
- (B) Perform the K-means algorithm with $K = 2$ and initial centroids $\{\mu_1, \mu_2\} = \{(1, 2), (7, 0)\}$ until convergence. Are the results different from (A)?
- (C) Now consider adding a data point $(5, 6)$ to \mathbf{X} . Show that the K-means algorithm converges to a local minimum with certain initial centroids. You can demonstrate this by showing that there exists a set of initial centroids that does not converge to the global minimum.

a.) Cluster 1: $\{(1, 2)\}$ Cluster 2: $\{(3, 4), (7, 0), (10, 2)\} \rightarrow \{\mu_1, \mu_2\} = \{(1, 2), (\frac{20}{3}, 2)\}$

Cluster 1: $\{(1, 2), (3, 4)\}$ Cluster 2: $\{(7, 0), (10, 2)\} \rightarrow \{\mu_1, \mu_2\} = \{(2, 3), (8.5, 1)\}$

Cluster 1: $\{(1, 2), (3, 4)\}$ Cluster 2: $\{(7, 0), (10, 2)\} \rightarrow \{\mu_1, \mu_2\} = \{(2, 3), (8.5, 1)\} \rightarrow \text{Converge.}$

b.) Cluster 1: $\{(1, 2), (3, 4)\}$ Cluster 2: $\{(7, 0), (10, 2)\} \rightarrow \{\mu_1, \mu_2\} = \{(2, 3), (8.5, 1)\}$

Cluster 1: $\{(1, 2), (3, 4)\}$ Cluster 2: $\{(7, 0), (10, 2)\} \rightarrow \{\mu_1, \mu_2\} = \{(2, 3), (8.5, 1)\} \rightarrow \text{Converge. , No different.}$

c.) Consider initial centroids: $\{\mu_1, \mu_2\} = \{(1, 2), (7, 0)\}$

Cluster 1: $\{(1, 2), (3, 4), (5, 6)\}$ Cluster 2: $\{(7, 0), (10, 2)\} \rightarrow \{\mu_1, \mu_2\} = \{(3, 4), (8.5, 1)\}$

Cluster 1: $\{(1, 2), (3, 4), (5, 6)\}$ Cluster 2: $\{(7, 0), (10, 2)\} \rightarrow \{\mu_1, \mu_2\} = \{(3, 4), (8.5, 1)\} \rightarrow \text{Converge.}$

And consider initial centroids: $\{\mu_1, \mu_2\} = \{(1, 2), (5, 6)\}$ * We know that these two has different answers, not matter who is optimal, it do show that k-means alg. could converge to a local min. with certain centroids.

Cluster 1: $\{(1, 2)\}$ Cluster 2: $\{(3, 4), (5, 6), (7, 0), (10, 2)\} \rightarrow \{\mu_1, \mu_2\} = \{(1, 2), (6.25, 3)\}$

Cluster 1: $\{(1, 2), (3, 4)\}$ Cluster 2: $\{(5, 6), (7, 0), (10, 2)\} \rightarrow \{\mu_1, \mu_2\} = \{(2, 3), (7.33, 2.66)\}$

Cluster 1: $\{(1, 2), (3, 4)\}$ Cluster 2: $\{(5, 6), (7, 0), (10, 2)\} \rightarrow \{\mu_1, \mu_2\} = \{(2, 3), (7.33, 2.66)\} \rightarrow \text{Converge.}$

Programming Part

(a)

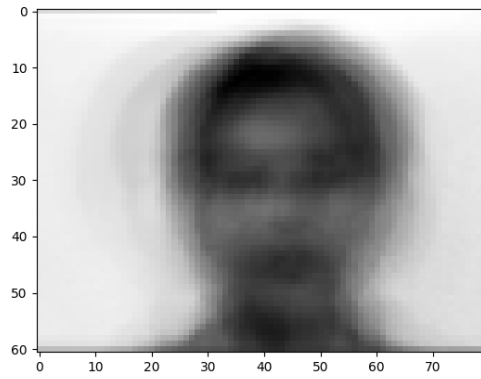


Fig1. Mean vector

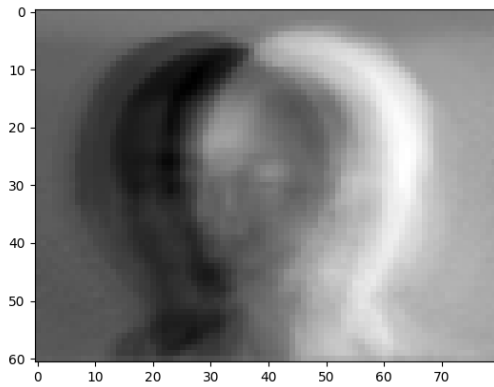


Fig.2 1st eigenvector

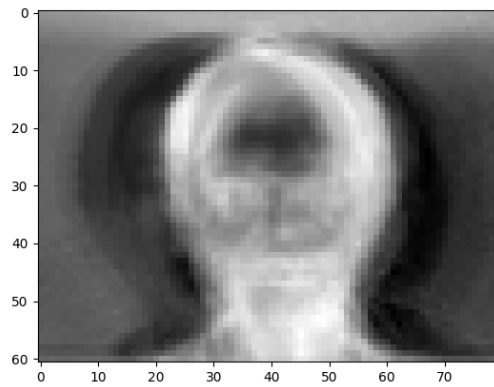


Fig3. 2nd eigenvector

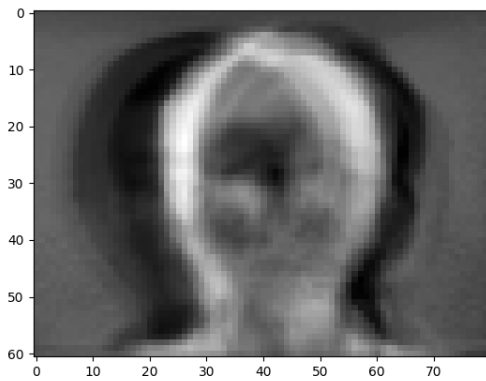


Fig.4 3rd eigenvector

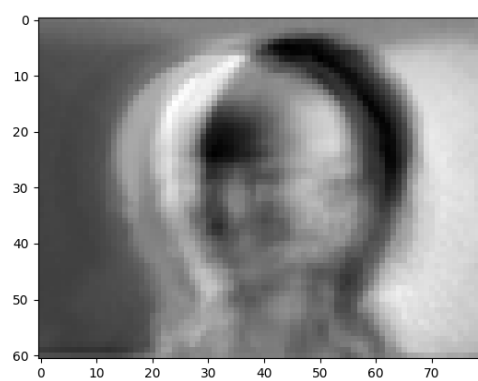


Fig5. 4th eigenvector

(b)

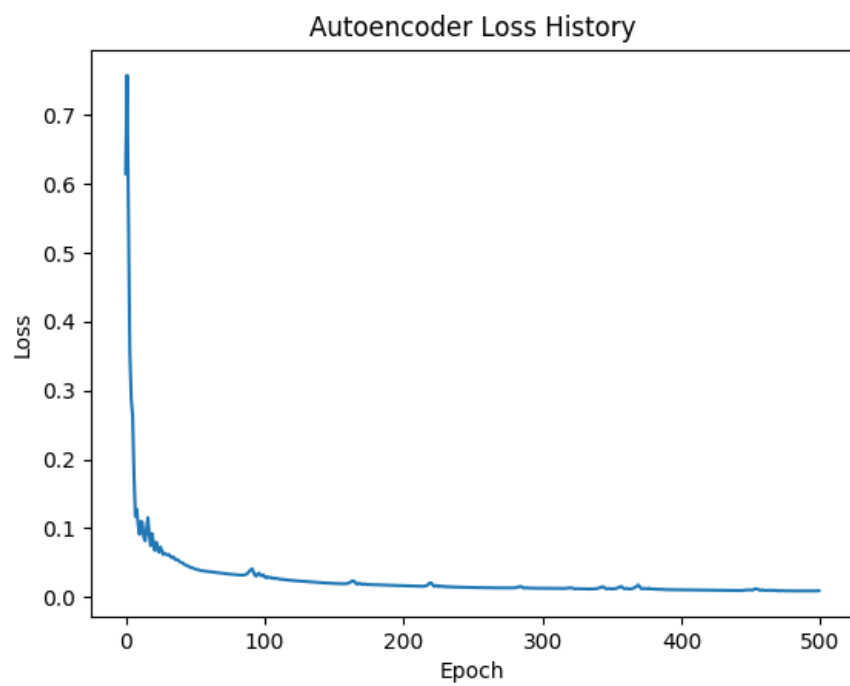


Fig6. Loss curve of Autoencoder

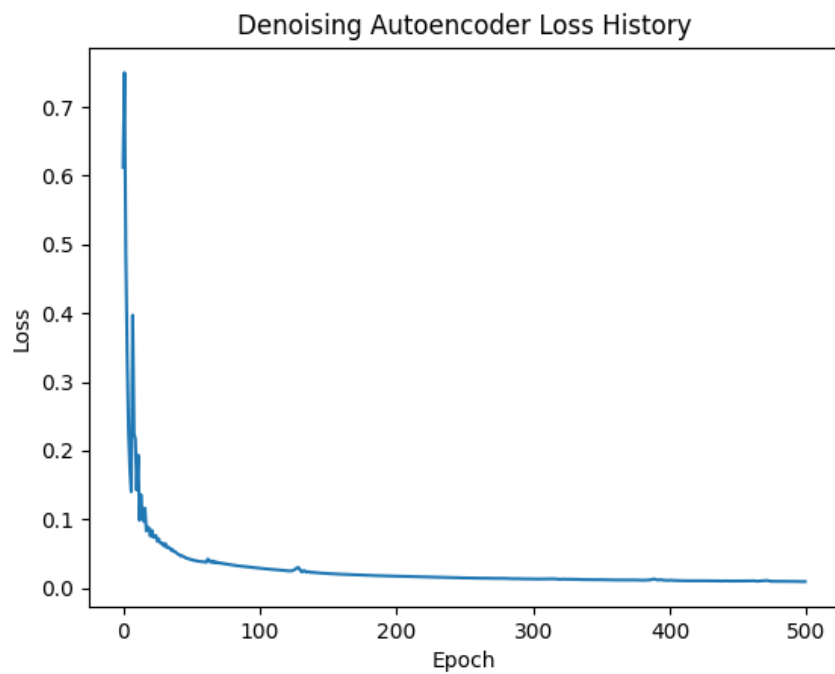


Fig7. Loss curve of DenoisingAutoencoder

(c)

Original Image



Fig8. Original image

PCA Reconstructed Image



Fig9. Reconstructed with PCA (MSE= 0.010710469688056314)

Autoencoder Reconstructed Image



Fig10. Reconstructed with Autoencoder(MSE= 0.015646440532965353)

DenoisingAutoencoder Reconstructed Image



Fig11. Reconstructed with DenosingAutoencoder(MSE= 0.012897145575758865)

(d)

I've tried the deeper and shallower architecture with all the settings remains the same as the sample.

Original one: **Acc:** 0.9333333333333333 and **Reconstruction Loss:** 0.013306048900814367

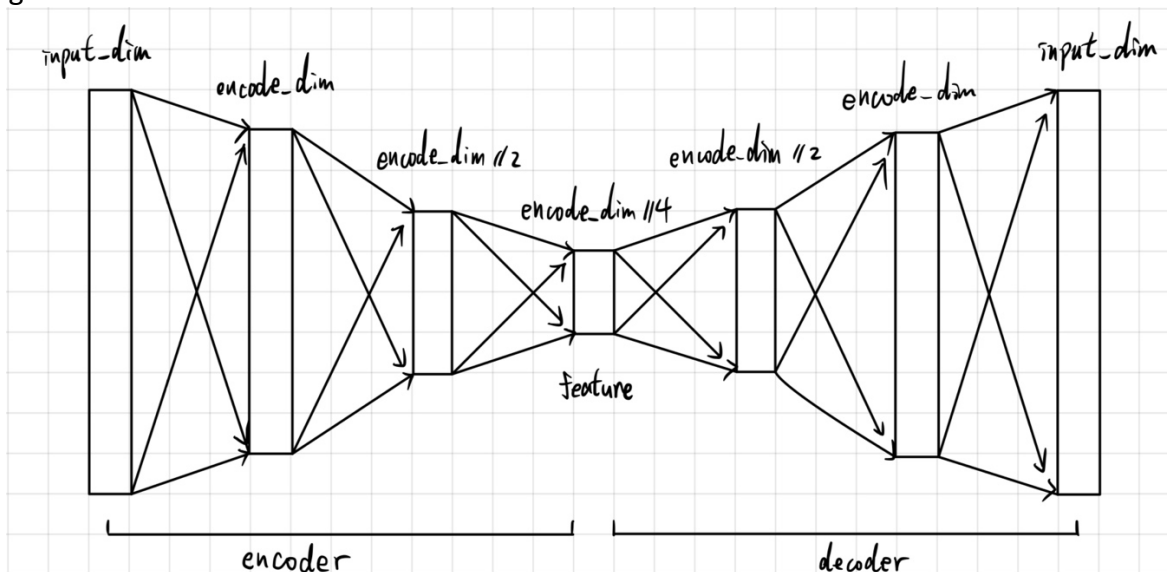


Fig12. The deeper architecture

The deeper one has the

Acc: 0.06666666666666667 and **Reconstruction Loss:** 0.042025545732162176

We can see that its performance is worse than the original one. I think it might because of the decoding would be harder for deeper model.

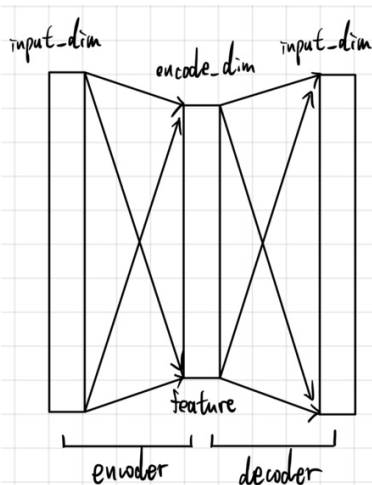


Fig13. The shallower architecture

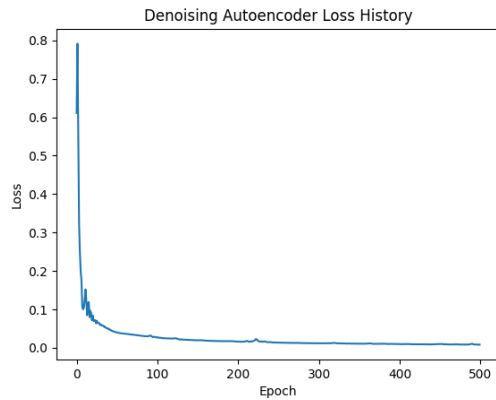
The shallower one has the

Acc: 0.8666666666666667 and **Reconstruction Loss:** 0.019420283759560463

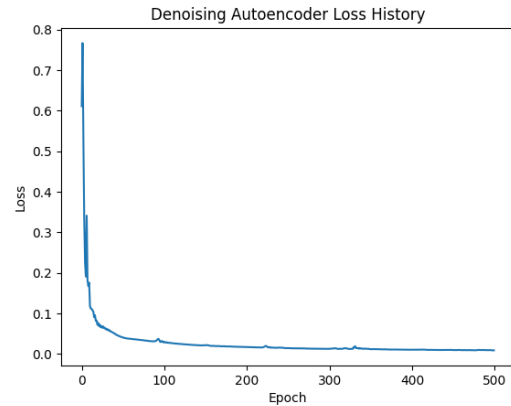
We can see that the shallow one performs worse than the original one too. But has a slightly difference, this might suggest that the task is not that complex so the easy one could handle but still need the suitable settings to get the best performance. Which means that the architecture is truly important for the task.

(e)

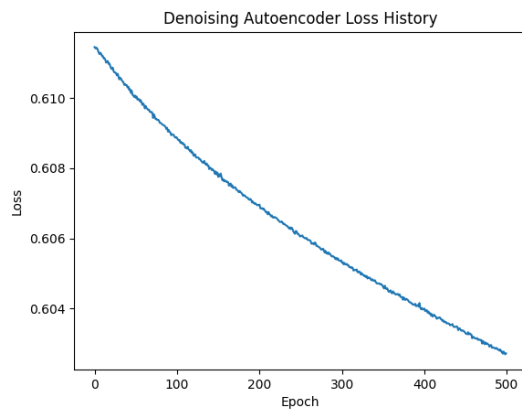
I've tested 5 different optimizers: Adam, AdamW, Adadelata, PMSpropk, and SGD.



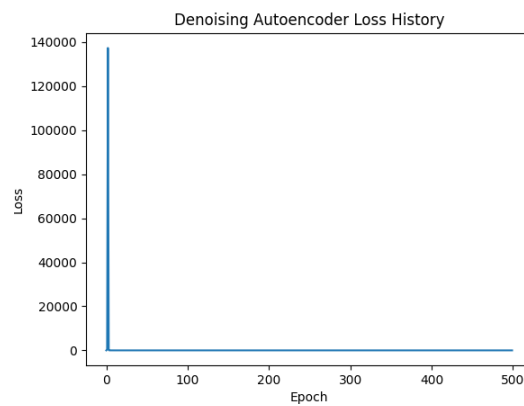
Adam



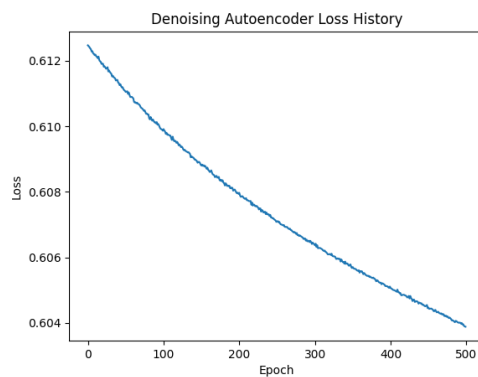
AdamW



Adadelata



RMSprop



SGD

Optimizer	Accuracy	Reconstruct Loss
Adam	0.9333	0.0133
AdamW	0.9333	0.0132
Adadelata	0.8667	0.6919
RMSprop	0.0667	0.1657
SGD	0.9000	0.6934

We can see that the Adam and AdamW have the similar results, which is pretty reasonable. And Adadelat's curve is converge so slow and has quite good accuract but bloody terrible loss. In contrast, the RMSprop has bloody terrible accuracy but converge so fast with quite a large loss. And lastly, SGD has pretty good accuracy bt loss terrible as bad as Adadelat, which would oberved by the loss history of these 2 models, I think the trend indeed influence the loss.