

Hand-written Part

学号: +11902210 姓名: 张凡 (交换生)

code question report 在后面

Problem 1

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

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

Now Swish is defined as

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

What is $\varphi'(s)$?

解:

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

$$\begin{aligned} \varphi'(s) &= \theta(s) + s \cdot \frac{d(\theta(s))}{ds} \quad \Leftarrow \frac{d(\theta(s))}{ds} = \theta'(s) = \left(\frac{1}{1 + e^{-s}} \right)' \\ &= \frac{1}{1 + e^{-s}} + \frac{s \cdot e^{-s}}{(1 + e^{-s})^2} &= -\frac{e^{-s} \cdot (-1)}{(1 + e^{-s})^2} = \frac{e^{-s}}{(1 + e^{-s})^2} \\ &= \frac{1 + e^{-s} + s \cdot e^{-s}}{(1 + e^{-s})^2} \\ &= \frac{1 + (s+1)e^{-s}}{(1 + e^{-s})^2} & \varphi'(s) = \frac{1 + (s+1)e^{-s}}{(1 + e^{-s})^2} \end{aligned}$$

Problem 2

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: $\text{in}(P_j)$ (P_j has a link to P_i , if $i \neq j$). 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: $\text{PageRank}_t(P_i) = \frac{\sum_{P_j \in \text{in}(P_i)} \text{PageRank}_t(P_j)}{|\text{out}(P_i)|}$
where $|\text{out}(P_i)|$ denotes the out-degree of P_i .
3. Repeat step 2 until the ranks converge.

Now, consider the following example with 3 pages:

$$P = \{P_1, P_2, P_3\} \quad \text{in}(P_1) = \{P_2, P_3\} \quad \text{in}(P_2) = \{P_3\} \quad \text{in}(P_3) = \{P_1\}$$

We define the following variables:

$$V_t = (\text{PageRank}_t(P_1), \text{PageRank}_t(P_2), \text{PageRank}_t(P_3))^T$$

$$P = \begin{bmatrix} 0 & 1 & 0.5 \\ 0 & 0 & 0.5 \\ 1 & 0 & 0 \end{bmatrix}.$$

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

- (A) Based on the iterative algorithm above, please compute the values of v_1, v_2, v_3, v_4 and v_5 .

Step 1: according to the algorithm.

$$\text{Step 1: } P = \{P_1, P_2, P_3\} \Rightarrow n=3 \Rightarrow V_0 = \left(\frac{1}{3}, \frac{1}{3}, \frac{1}{3} \right)^T$$

initialize PageRank_i(P_i) = $\frac{1}{3}$
($i=1, 2, 3$) $\Rightarrow V_0 = \left(\frac{1}{3}, \frac{1}{3}, \frac{1}{3} \right)$

$$\text{Step 2-1: } V_1 = P \cdot V_0 = P V_0 = \left(\begin{array}{ccc|c} 0 & 1 & \frac{1}{2} & \frac{1}{3} \\ 0 & 0 & \frac{1}{2} & \frac{1}{3} \\ 1 & 0 & 0 & \frac{1}{3} \end{array} \right) \left(\begin{array}{c} \frac{1}{3} \\ \frac{1}{3} \\ \frac{1}{3} \end{array} \right) = \left(\begin{array}{c} \frac{1}{2} \\ \frac{1}{6} \\ \frac{1}{3} \end{array} \right)$$

$$\left(\begin{array}{cc|c} 0 & 1 & \frac{1}{2} \\ 0 & 0 & \frac{1}{2} \\ 1 & 0 & 0 \end{array} \right) \left(\begin{array}{c} \frac{1}{3} \\ \frac{1}{3} \\ \frac{1}{3} \end{array} \right) = \left(\begin{array}{c} \frac{1}{2} \\ \frac{1}{6} \\ \frac{1}{3} \end{array} \right)$$

$$\text{Step 2-2: } V_2 = P V_1 = \left(\begin{array}{cc|c} 0 & 1 & \frac{1}{2} \\ 0 & 0 & \frac{1}{2} \\ 1 & 0 & 0 \end{array} \right) \left(\begin{array}{c} \frac{1}{2} \\ \frac{1}{6} \\ \frac{1}{3} \end{array} \right) = \left(\begin{array}{c} \frac{1}{3} \\ \frac{1}{6} \\ \frac{1}{2} \end{array} \right)$$

$$\text{Step 2-3: } V_3 = P V_2 = \left(\begin{array}{cc|c} 0 & 1 & \frac{1}{2} \\ 0 & 0 & \frac{1}{2} \\ 1 & 0 & 0 \end{array} \right) \left(\begin{array}{c} \frac{1}{3} \\ \frac{1}{6} \\ \frac{1}{2} \end{array} \right) = \left(\begin{array}{c} \frac{5}{12} \\ \frac{1}{4} \\ \frac{1}{3} \end{array} \right)$$

$$\text{Step 2-4: } V_4 = P V_3 = \left(\begin{array}{cc|c} 0 & 1 & \frac{1}{2} \\ 0 & 0 & \frac{1}{2} \\ 1 & 0 & 0 \end{array} \right) \left(\begin{array}{c} \frac{5}{12} \\ \frac{1}{4} \\ \frac{1}{3} \end{array} \right) = \left(\begin{array}{c} \frac{5}{12} \\ \frac{1}{6} \\ \frac{5}{12} \end{array} \right)$$

$$\text{Step 2-5: } V_5 = P V_4 = \left(\begin{array}{cc|c} 0 & 1 & \frac{1}{2} \\ 0 & 0 & \frac{1}{2} \\ 1 & 0 & 0 \end{array} \right) \left(\begin{array}{c} \frac{5}{12} \\ \frac{1}{6} \\ \frac{5}{12} \end{array} \right) = \left(\begin{array}{c} \frac{3}{8} \\ \frac{5}{24} \\ \frac{5}{12} \end{array} \right)$$

In conclusion

$$\left\{ \begin{array}{l} V_0 = \left(\frac{1}{3}, \frac{1}{3}, \frac{1}{3} \right)^T \\ V_1 = \left(\frac{1}{2}, \frac{1}{6}, \frac{1}{3} \right)^T \\ V_2 = \left(\frac{1}{3}, \frac{1}{6}, \frac{1}{2} \right)^T \\ V_3 = \left(\frac{5}{12}, \frac{1}{4}, \frac{1}{3} \right)^T \\ V_4 = \left(\frac{5}{12}, \frac{1}{6}, \frac{5}{12} \right)^T \\ V_5 = \left(\frac{3}{8}, \frac{5}{24}, \frac{5}{12} \right)^T \end{array} \right.$$

- (B) Recall that in the class, we mentioned that the algorithm converges when $V^* = PV^*$ solve this equation to derive V^* . Note that you should provide a normalized V^* as answer.

解：设 $V = \alpha \begin{pmatrix} v_1 \\ v_2 \\ v_3 \end{pmatrix}$ (\Leftrightarrow) $PV^* = V^* \Rightarrow (P - E)V^* = 0$ (E为单位矩阵)

$$P = \begin{pmatrix} 0 & 1 & \frac{1}{2} \\ 0 & 0 & \frac{1}{2} \\ 1 & 0 & 0 \end{pmatrix} \quad (P - E) = \begin{pmatrix} -1 & 1 & \frac{1}{2} & \frac{1}{2} \\ 0 & -1 & \frac{1}{2} & 0 \\ 1 & 0 & -1 & 0 \end{pmatrix}$$

$$(P - E)V^* = \alpha \begin{pmatrix} -v_1 + v_2 + \frac{1}{2}v_3 \\ -v_2 + \frac{1}{2}v_3 \\ v_1 - v_3 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}$$

$$\Rightarrow \begin{cases} -v_1 + v_2 + \frac{1}{2}v_3 = 0 \\ -v_2 + \frac{1}{2}v_3 = 0 \\ v_1 - v_3 = 0 \end{cases} \Rightarrow \begin{cases} 2v_2 = v_3 \\ v_1 = v_3 \end{cases}$$

$$V^* = \alpha \begin{pmatrix} 2v_2 \\ v_2 \\ 2v_2 \end{pmatrix} = \alpha \cdot v_2 \begin{pmatrix} 2 \\ 1 \\ 2 \end{pmatrix}$$

$$V^* \Rightarrow \text{normalized} \quad \sqrt{2^2 + 1^2 + 2^2} = 3$$

$$\text{解得 } V^* = \begin{pmatrix} \frac{2}{3} \\ \frac{1}{3} \\ \frac{2}{3} \end{pmatrix}$$

Consider $L \in \{1, 2, 3\}$
and

Problem 3

For a $d^{(0)} - d^{(1)} - d^{(2)} - \dots - d^{(L)}$ fully-connected neural network, the total number of neurons is

$$D = \sum_{l=1}^L d^{(l)} \quad \text{and the total number of weight is } \sum_{l=1}^L d^{(l-1)}d^{(l)}$$

which ignore the so-called 'bias' terms in the network. Let the number of inputs $d^{(0)} = 10$ and the total number of neurons $D = 100$. Over all possible choices of L , $d^{(1)}, d^{(2)}, \dots, d^{(L)}$. What is the maximum total number of weights (and what neural network architecture does it correspond to)? What is the minimum total number of weights (and what neural network architecture does it correspond to)?

解：maximum: 选用最大的 $d^{(l)}$ 使得 L 为 ~~尽可能大~~ $d^{(1)} = d^{(2)} = \dots = d^{(L)} = 10$

$$\sum_{l=1}^L d^{(l-1)}d^{(l)} = L \cdot 10 \cdot 10 = 100L$$

$$d^{(0)}d^{(1)} + d^{(1)}d^{(2)} + \dots + d^{(L-1)}d^{(L)}$$

$$\text{类似地, } a+b = 100 \quad (a>0, b>0)$$

$$6000 = (a+b)^2 \geq a^2 + 2ab + b^2 \geq a^2 + b^2.$$

$L \in \{1, 2, 3\}$ 因此在总 D 等于的情况下， L 越大， $\sum_{l=1}^L d^{(l-1)}d^{(l)}$ 越大。

因此 L 越大，

最大值的 L 的设计方法为 $d^{(l)} = 10$

但更好的 $d^{(l)}$ 是

我们设计过几次 ($L=1, 2, 3$)

$$\textcircled{1} \quad L=1 \quad \sum_{l=1}^L d^{(l-1)}d^{(l)} = 10 \times 100 = 1000$$

$$\textcircled{2} L=2 \text{ 时 } \sum_{i=1}^2 d(i-1)d(i) = d(0)d(1) + d(1)d(2)$$

$$\begin{cases} d(1)=x \\ d(2)=y \end{cases} \quad W = 10x+xy = 10x+x(100-x) = -x^2+110x = -(x-55)^2+3025 \leq 3025$$

$$\textcircled{3} L=3 \text{ 时 } W = 10d(1) + d(1)d(2) + d(2)d(3)$$

$$\begin{cases} x,y,z=d_1,d_2,d_3 \\ x+y+z=100 \end{cases}$$

$$d(1) > 0 \quad W = 10x+xy+yz = 10x+xy+(100-x-y)y = 10x+100y-y^2 = 10 \cdot (100y-2)$$

$$\frac{\partial W}{\partial x} = 10 \rightarrow 0 \quad \frac{\partial W}{\partial y} = -2y+100 = -y^2+100y-y^2$$

$$\text{当 } z \text{ 取最小值 } z \geq 0, \quad y = 45, \quad z = 1$$

$$x = 54$$

$$1000 < \frac{3025}{3015} < 3025$$

$$3015 <$$

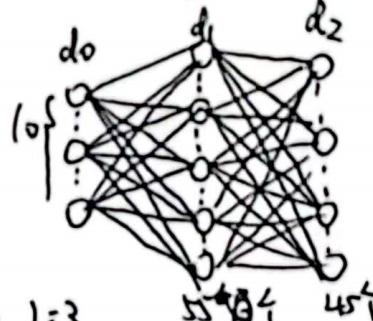
$$= -(\cancel{y^2}) + \cancel{2500} + 10x + 100y - y^2 = -(\cancel{y^2}) + 1000 - y^2 - 10z + 90y$$

$$= \frac{3000}{3025} - (y-45)^2 - 10z$$

$$= 3025 - 0 - 10 = 3015$$

因此在 $L=2$ 时 $\exists d(1)=55, d(2)=45$ 取最大值

神经网络架构:



$$\sum_{i=1}^L d(i-1)d(i) = 3025 \quad (\text{最大值})$$

minimize (最小值): 也就是说 $L=1, L=2, L=3$

$$\textcircled{1} L=1 \text{ 时, } \text{寻求最大值} \quad \sum_{i=1}^L d(i-1)d(i) = 10 \times 100 = 1000$$

$$\textcircled{2} L=2 \text{ 时, 同样设 } d(1)=x, d(2)=y$$

$$W = -(x-55)^2 + 3025 \quad \nabla |x-55| \text{ 最大}$$

$$x > 0 \quad x=1 \quad W = -(1-55)^2 + 3025 = -1 + 110 \times 1 = -1 + 110 = 109$$

$$\textcircled{3} L=3 \text{ 时, 同样设 } d(1)=x, d(2)=y, d(3)=z$$

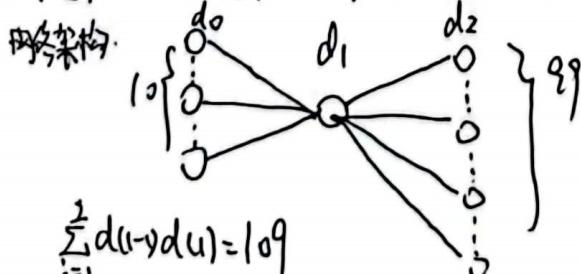
$$W = -(y-45)^2 - 10z + 3025 \quad x, y, z > 0$$

$$z \text{ 的约束 } 10 < 0, |y-45| \text{ 最大值, } z \text{ 最大值 } \quad x=1, y=1$$

$$W = 10x + 1 \times 1 + 1 \times 98 = 109 \quad z = 100 - 2 = 98$$

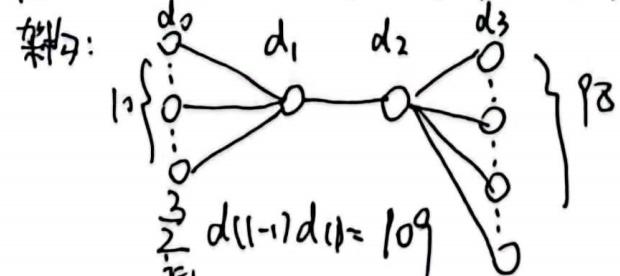
$$1000 > 109 \Rightarrow 109 = 109 \quad \text{最小值有多种情况}$$

$$\text{情况1: } L=2, d(1)=1, d(2)=99$$



$$\sum_{i=1}^L d(i-1)d(i) = 109$$

$$\text{情况2: } L=3, d(1)=d(2)=1, d(3)=98$$



$$\sum_{i=1}^L d(i-1)d(i) = 109$$

Code Part Problem Report

學號: t11902210 姓名: 張一凡

- (a) Plot the mean vector as an image as well as the top 4 eigenvectors, each as an image, by calling the given plot component routine. Each of those top eigenvectors is usually called an eigenface that physically means an informative “face ingredient.”

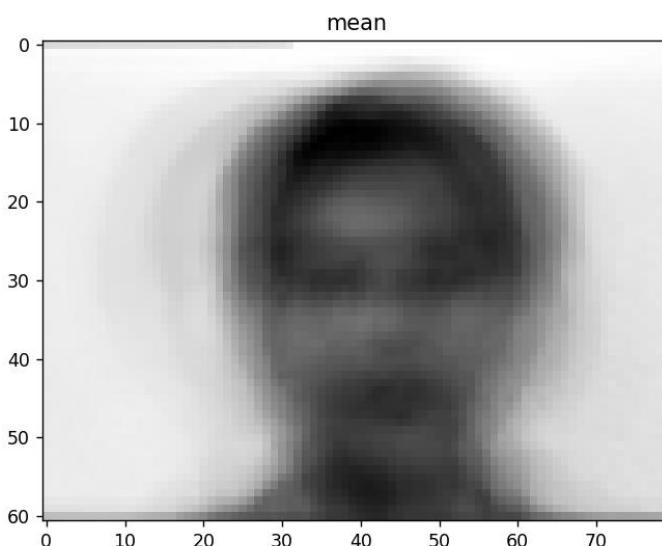
Answer:

First, what I am showing is the code that generates the average vector and the first four feature vectors. I use the imshow function to display the output, and because it is a grayscale image, I have added the parameter cmap='gray'. However, since the image set is trained by PCA to form a 4880-dimensional vector, there are certain requirements for the image's reshaping. Later, I used the prompts and load in the question_. The data function attempted to output the image specification and found that the image is 61 * 80 in size, which was successfully output. In addition, when outputting feature vectors, due to the convenience of performing transformations in the PCA definition, the specification is 4880 * 40. Therefore, I need to perform some transposing before outputting in order. Here is my code:

```
# plot Q1
img = (pca.mean).reshape(61, 80) # image size is 61*80 = 4880
plt.imshow(img, cmap='gray')
plt.title("mean")
plt.show()

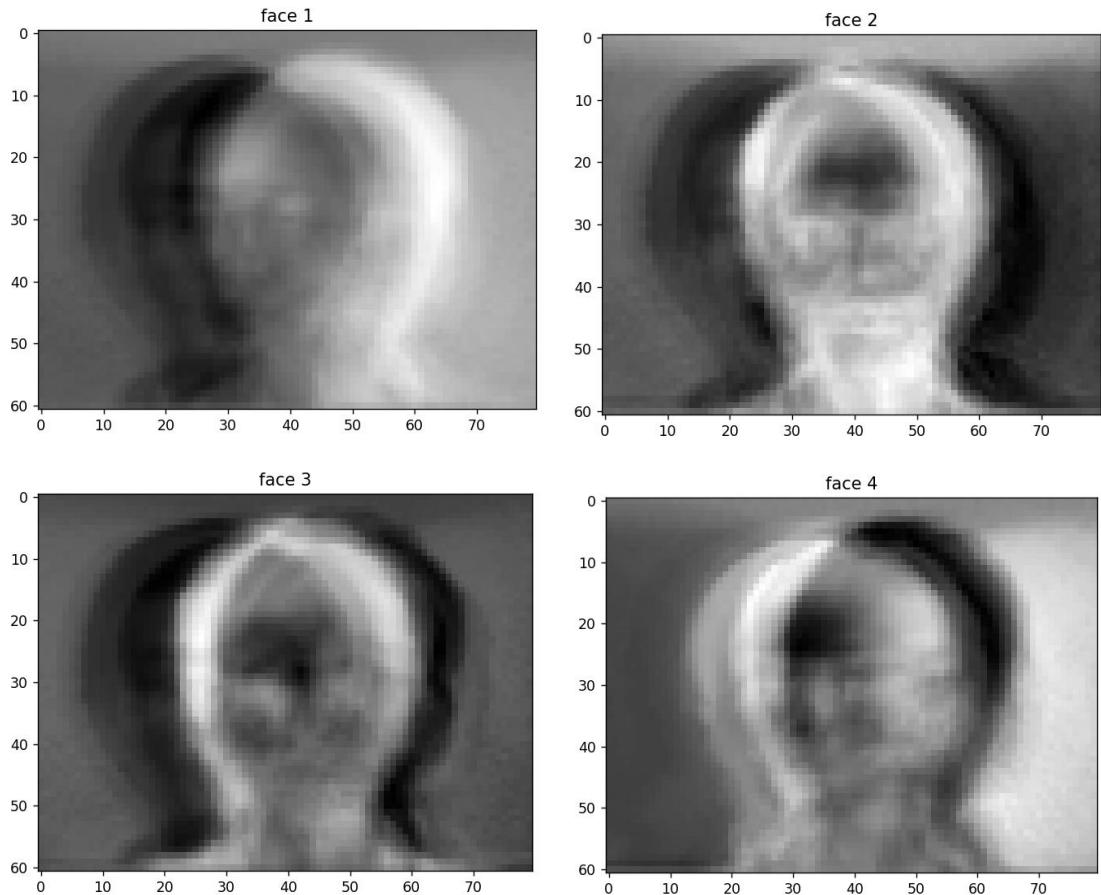
for i in range(4):
    eigenvector = pca.components.T[i]
    img = (eigenvector).reshape(61, 80)
    plt.imshow(img, cmap='gray')
    plt.title(f'face {i+1}')
    plt.show()
```

The first thing to display is the image drawn through the average vector, which displays the basic elements of the portrait, while also fusing the common features of multiple facial images as the average face.



The following are the images drawn by the first four feature vectors that I have selected.

I have marked the result of which vector is displayed on each image. From the image display, it can be seen that different feature vectors correspond to different faces, which are somewhat different from the average face in the whiter part in the middle. This may be due to the difference between each vector and the principal component, resulting in the display of different features, such as lighting and facial expressions, Gender characteristics, etc. (personal analysis)



(b) Plot the training curve of Autoencoder and DenoisingAutoencoder

Answer:

Firstly, I also demonstrate the process of calculating the loss rate according to the requirements of the question in Autoencoder and Denoising Autoencoder, and outputting the curve based on the number of iterations or epochs. I have set a variable loss in each epoch_ Num and divide it by the number to calculate the average loss for each iteration process and place it in a list loss_ In the record, finally output the curve through plot. In the DenoisingAutoencoder, I decided to use Gaussian noise after searching for data. Since the process of recording losses and generating curves is basically the same for both the Autoencoder and DenoisingAutoencoder, I chose to display the code of the DenoisingAutoencoder, as well as the code that generates noise.

```

loss_record = []
for e in tqdm(range(epochs), desc="Denoising-fitting"):
    loss_num = 0
    k = 0
    for item in data_load:
        x = item[0]
        noise = self.add_noise(x) + x

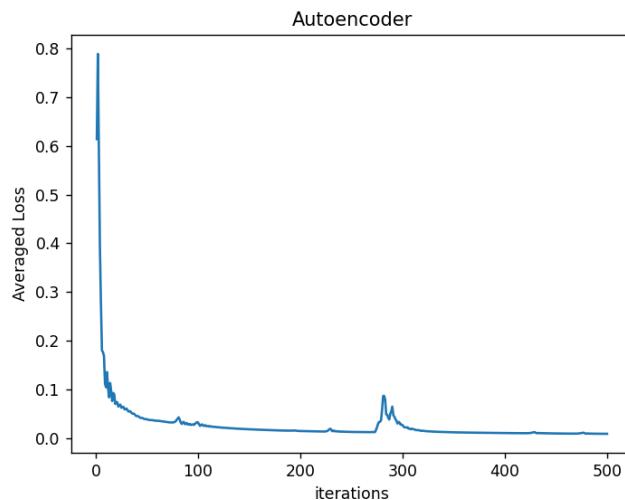
        optimizer.zero_grad()
        gx = self.forward(noise)
        loss = loss_func(gx, x)
        loss.backward()
        optimizer.step()

    def add_noise(self, x):
        #TODO: 2%
        # Gaussian noise
        noise = self.noise_factor * torch.randn(*x.shape)
        result = torch.clamp(noise, -1, 1)
        return noise
        #raise NotImplementedError

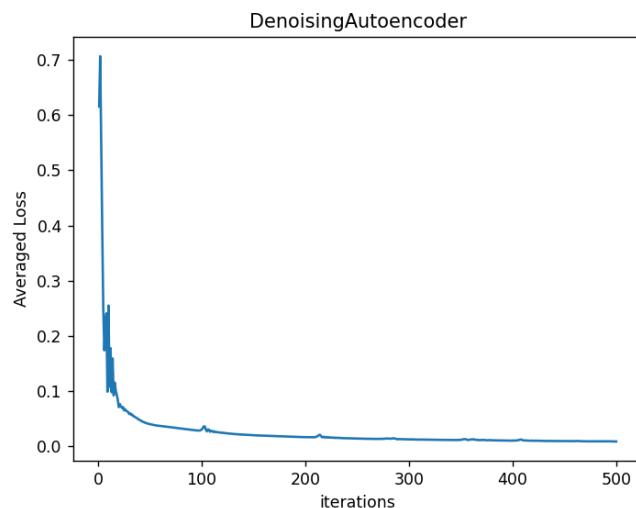
```

Here is the training curve of Autoencoder and Denoising Autoencoder.

Autoencoder average loss curve



DenoisingAutoencoder average loss curve



From the results presented by the two curves, it can be seen that both Autoencoder and DenoisingAutoencoder can handle this problem well, with fast convergence speed. Therefore, they will gradually learn the feature representation of the input data, thereby gradually reducing the error. However, due to the addition of noise, DenoisingAutoencoder has greater advantages over Autoencoder. In the observation of this problem, it is mainly reflected in two aspects: 1. Although not obvious, from the subtle differences in the curve, it can be seen that the curve DenoisingAutoencoder exhibits faster convergence speed and ultimately lower error, indicating that adding noise will make the final loss rate smaller; 2. Although the training curve of DenoisingAutoencoder exhibits more fluctuations and fluctuations in the initial stage, it is extremely stable after convergence. On the contrary, Autoencoder exhibits significant fluctuations and instability in the later stage, indicating that DenoisingAutoencoder has stronger robustness and better performance.

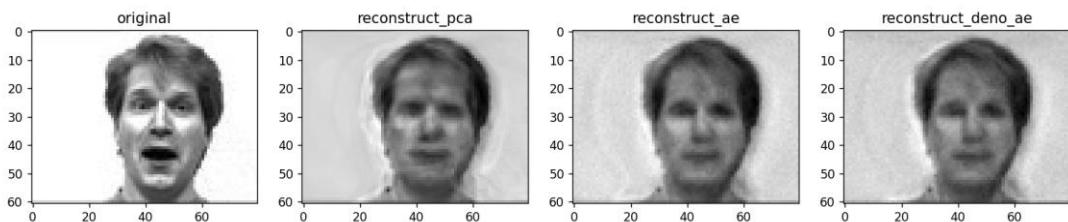
(c) Plot the original image and the images reconstructed with PCA, Autoencoder, and DenoisingAutoencoder side by side, ideally as large as possible. Then, list the mean squared error between the original image and each reconstructed image

Answer:

Firstly, I also presented the code to solve this problem. As the requirement in the question is to draw four images side by side, I called the subplots function to generate sub images and reshape each of the four images and display them on one image. The following is my implementation code.

```
# plot Q3
fig, axs = plt.subplots(1, 4, figsize=(20, 5))
img_original = (img_vec).reshape(61, 80) # image size is 61*80 = 4880
axs[0].imshow(img_original, cmap='gray')
axs[0].set_title("original")
img_pca = (img_reconstruct_pca).reshape(61, 80) # image size is 61*80 = 4880
axs[1].imshow(img_pca, cmap='gray')
axs[1].set_title("reconstruct_pca")
img_ae = (img_reconstruct_ae).reshape(61, 80) # image size is 61*80 = 4880
axs[2].imshow(img_ae, cmap='gray')
axs[2].set_title("reconstruct_ae")
img_deno_ae = (img_reconstruct_deno_ae).reshape(61, 80) # image size is 61*80 = 4880
axs[3].imshow(img_deno_ae, cmap='gray')
axs[3].set_title("reconstruct_deno_ae")
plt.show()
```

The following are the image results I generated, from left to right are the original image, reconstructed with PCA image, reconstructed with Autoencoder image, and reconstructed with DenoisingAutoencoder image.



Next, I will list the mean square error between the original image and each reconstructed image, which is the output result of the given code.

```
Reconstruction Loss with PCA: 0.010710469688056319
Reconstruction Loss with Autoencoder: 0.012720367232229133
Reconstruction Loss with DenoisingAutoencoder: 0.01331009228793101
```

Firstly, among the square error reconstruction errors, PAC has the lowest reconstruction error of 0.011, while AE and denoised AE have lower reconstruction errors of 0.0127 and 0.0133, respectively. This result may indicate that compared to PCA, the ability of automatic encoders and denoising automatic encoders to reconstruct the original image in this task is relatively weak. This may be because PCA preserves the information of the original data as much as possible during the reconstruction process, so it can perform well in terms of square error reconstruction error. At the same time, the automatic encoder and denoising automatic encoder are based on neural networks and can capture nonlinear structures in the data. Although they are not as good as PAC in this problem, they also perform well in terms of absolute results in the reconstructed image. Finally, when observing the results of image display, AE and denoising AE cannot visually generate images as close to the original image as possible compared to PCA, making the PCA method clearer. After searching for information, my conclusion is that they may be trying to learn more "smooth" or "average" data representations, which may cause the reconstructed image to appear a bit blurry, but this may indicate that AE and denoising AE are learning the internal structure of the data, rather than just remembering the input data. Overall, low reconstruction errors do not always mean better performance.

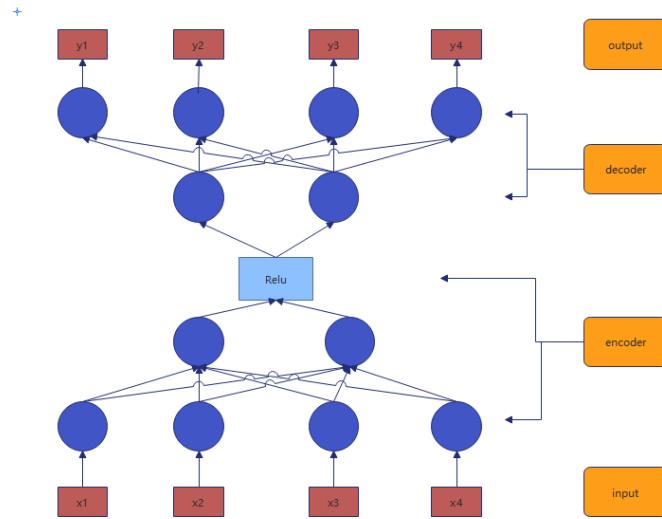
(d) Modify the architecture in Autoencoder in its constructor. Try at least two different network architectures for the denoising autoencoder. You can consider trying a deeper or shallower or fatter or thinner network. You can also consider adding convolutional layers and/or other activation functions. Draw the architecture that you have tried and discuss your findings, particularly in terms of the reconstruction error that the architecture can achieve after decent optimization.

Answer:

I have made three modifications to this issue and combined them with the original results to display and output. I will summarize my observations and analysis later on. At the same time, I also used concept maps as required to complete the structural diagrams of these self encoders. Due to space limitations, I will set the dimensions of all inputs and outputs to 4, which is convenient for comparison and observation. I will place concept maps of each variant and original structure at the end of each code for inspection and explanation.

Original design

```
self.encoder = nn.Sequential(
    nn.Linear(input_dim, encoding_dim),
    nn.Linear(encoding_dim, encoding_dim//2),
    nn.ReLU()
)
self.decoder = nn.Sequential(
    nn.Linear(encoding_dim//2, encoding_dim),
    nn.Linear(encoding_dim, input_dim),
)
```



Acc from Autoencoder: 0.8666666666666667

Acc from DenoisingAutoencoder: 0.9333333333333333

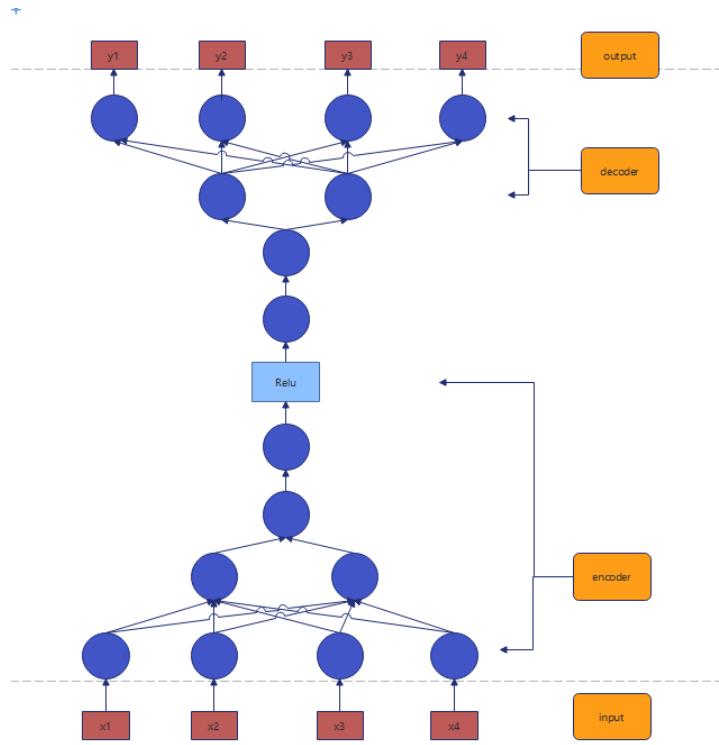
Reconstruction Loss with Autoencoder: 0.015411890219807605

Reconstruction Loss with DenoisingAutoencoder: 0.012947307032105934

This is the accuracy and Reconstruction Loss of the original data, and each variant below will be compared with this and a conclusion will be drawn.

Deeper design

```
# deeper
self.encoder = nn.Sequential(
    nn.Linear(input_dim, encoding_dim),
    nn.Linear(encoding_dim, encoding_dim // 2),
    nn.Linear(encoding_dim // 2, encoding_dim // 3),
    nn.Linear(encoding_dim // 3, encoding_dim // 4),
    nn.ReLU()
)
self.decoder = nn.Sequential(
    nn.Linear(encoding_dim // 4, encoding_dim // 3),
    nn.Linear(encoding_dim // 3, encoding_dim // 2),
    nn.Linear(encoding_dim // 2, encoding_dim),
    nn.Linear(encoding_dim, input_dim),
)
```



Acc from Autoencoder: 0.8666666666666667

Acc from DenoisingAutoencoder: 0.8666666666666667

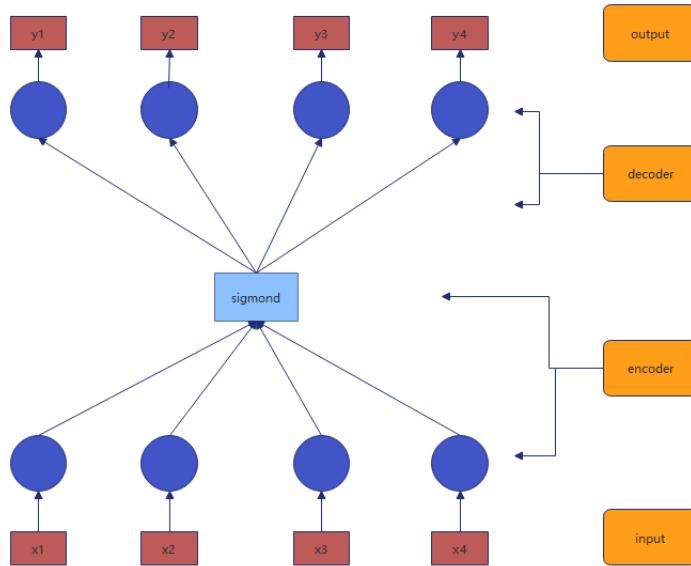
Reconstruction Loss with Autoencoder: 0.014501936661465207

Reconstruction Loss with DenoisingAutoencoder: 0.014558046191633872

The first type is transformed into a deeper network, with an accuracy AE of 0.87 for both DenoiseAE and DenoiseAE, which is basically the same as the original accuracy AE, but Dnoisae is slightly lower; In terms of Reconstruction Loss, both AE and DenoiseAE are approximately 0.0145. Compared with the original design, AE may be due to error issues, but the increase in DenoiseAE is more significant. I think the primary reason is the occurrence of overfitting. In the case of limited data, the deeper network may be overfitting on the training data, which makes its performance on the test data worse. At the same time, the deeper network may be more difficult to train when the optimizer is certain.

Shallower + sigmoid design

```
# shallower + sigmoid
self.encoder = nn.Sequential(
    nn.Linear(input_dim, encoding_dim),
    nn.Sigmoid()
)
self.decoder = nn.Sequential(
    nn.Linear(encoding_dim, input_dim),
)
```



Acc from Autoencoder: 0.5333333333333333

Acc from DenoisingAutoencoder: 0.1666666666666666

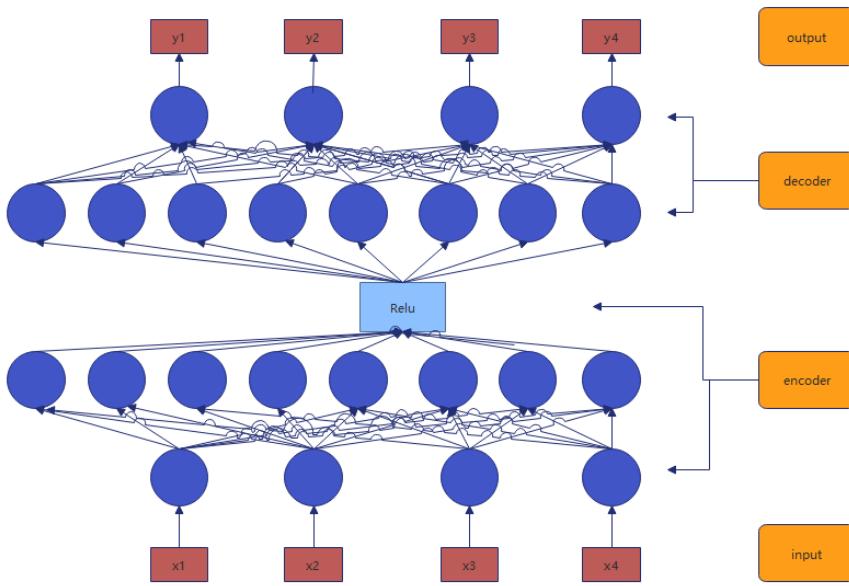
Reconstruction Loss with Autoencoder: 0.028569871536348768

Reconstruction Loss with DenoisingAutoencoder: 0.03531376575089117

The second is a shallower network with sigmoid activation function instead of relu, whose accuracy rate AE (0.53) and Denoise AE (0.17) are much lower than the original accuracy rate; In terms of Reconstruction Loss, AE and DenoiseAE were 0.029 and 0.035, respectively, indicating a significant increase compared to the original design. I think the reason is that underfitting occurs, and shallow networks may not have sufficient ability to learn complex data representations. This leads to significant errors when attempting to reconstruct input data, resulting in increased reconstruction losses and a significant decrease in accuracy; At the same time, the sigmoid activation function may limit the representation ability of the network. The output of the sigmoid function is between 0 and 1, which may not fully show all the characteristics of the data. The gradient of the original relu function in the positive part is 1, which will not have the problem of gradient disappearance. The result is more stable, reducing reconstruction loss and improving accuracy.

Fatter design

```
# fatter
self.encoder = nn.Sequential(
    nn.Linear(input_dim, encoding_dim),
    nn.Linear(encoding_dim, encoding_dim*2),
    nn.ReLU()
)
self.decoder = nn.Sequential(
    nn.Linear(encoding_dim*2, encoding_dim),
    nn.Linear(encoding_dim, input_dim),
)
```



Acc from Autoencoder: 0.9333333333333333

Acc from DenoisingAutoencoder: 0.9333333333333333

Reconstruction Loss with Autoencoder: 0.013918745855689495

Reconstruction Loss with DenoisingAutoencoder: 0.012682685163480024

The above are the three variations and original forms that I have shown that should be designed for this issue. The third type is transformed into a fatter network, with accuracy AE and DenoiseAE both reaching 0.93, which is relatively stable compared to the original accuracy; In terms of Reconstruction Loss, AE and DenoiseAE are 0.013 and 0.012 respectively, which are relatively stable and slightly improved compared to the original design. I think the reason is that it increases the width of the network with more hidden units, which can learn and store more information. Such a network may be more able to capture complex patterns in the data and generate lower errors when reconstructing inputs. This deeper network comparison with the first variant may be less likely to cause overfitting in the same situation, so it is more stable or even slightly improved compared with the original result.

The above are the three forms of variation and original design that I have presented for this issue, and the results and analysis are presented, with a focus on the Reconstruction Loss aspect. In fact, in the analysis, I believe that the impact and results of changing neural networks are similar to ordinary neural network models. However, in AE and DenoiseAE, it is mainly more convenient and advanced in image processing, which also shows me the accuracy of image processing and the relationship between reconstruction loss and layer number and layer width design.

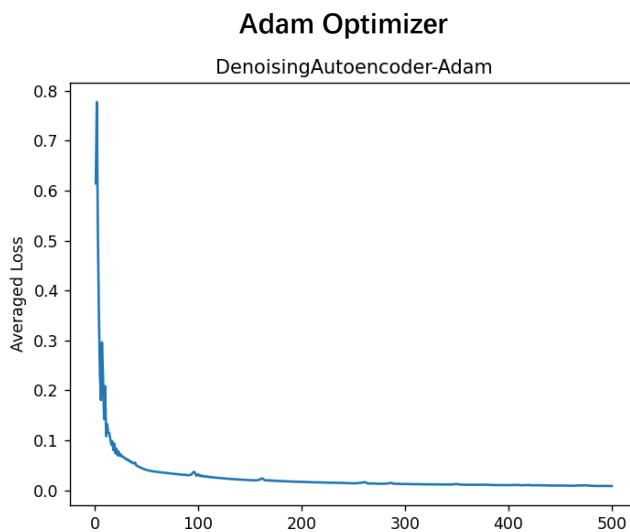
(e) Test at least 2 different optimizers, compare the training curve of DenoisingAutoencoder and discuss what you have found in terms of the convergence speed and overall performance of the model.

Answer:

An optimizer is an algorithm or method used in machine learning and deep learning to improve model performance and reduce prediction errors. The main objective of the optimizer is to minimize (or maximize) the loss function or objective function. In the neural network, the loss function measures the difference between the predicted value of the model and the actual value. The task of the optimizer is to find the parameters that can minimize the loss function (i.e. the weight and deviation of the model). I have selected two optimizers based on the requirements of the topic - Adam and SGD. Below, I will display their loss curves in DnoiseAE and explain them separately. Finally, I will make a comparative analysis. The generated image code is basically similar to Q1, so the image code part is not displayed.

```
# Adam
optimizer = optim.Adam(self.parameters(), lr=learning_rate)
# SGD
# optimizer = optim.SGD(self.parameters(), lr=learning_rate, momentum=0.7)
X_fit = torch.tensor(X).clone().detach().float()

data_pre = torch.utils.data.TensorDataset(X_fit)
```

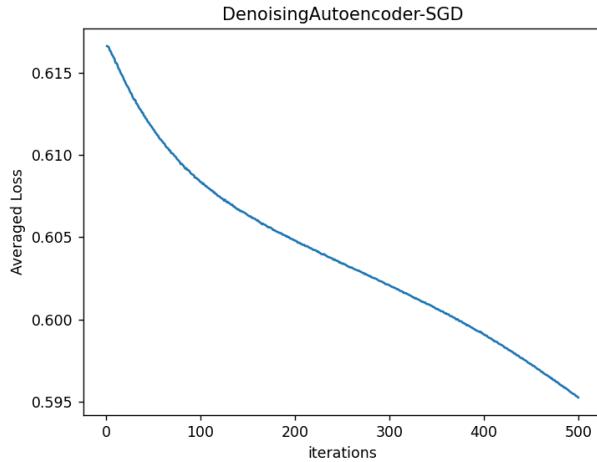


Acc from DenoisingAutoencoder: 0.9333333333333333

Reconstruction Loss with DenoisingAutoencoder: 0.013592173395249749

Adam (Adaptive Moment Estimation) is an optimization algorithm for deep learning models, particularly suitable for processing large-scale data and parameters. Adam combines the advantages of two other optimization algorithms: Adaptive Gradient Algorithm and Root Mean Square Propagation. Adam adjusts the learning rate separately for each parameter, which makes it particularly suitable for sparse data and non-uniform objective functions. At the same time, Adam also integrates the concept of momentum, which can help the optimizer skip local minima and accelerate convergence. Finally, Adam uses a deviation correction mechanism to obtain more accurate estimates faster from the initial stage of initialization and low gradient regions.

SGD Optimizer



Acc from DenoisingAutoencoder: 0.8666666666666667

Reconstruction Loss with DenoisingAutoencoder: 0.6826610005493821

SGD (Stochastic Gradient Descent) is a commonly used optimization algorithm in deep learning, used to train various types of models, especially neural networks. Random gradient descent only uses one training sample or a small batch of training samples during each update, which enables SGD to train faster and effectively handle large datasets. A common improvement method for SGD is to use momentum, which can help SGD accelerate in relevant directions and suppress oscillations, enabling it to converge to the optimal solution faster and more accurately.

The Adam optimizer has a significantly faster convergence speed in the results I have shown, with lower accuracy (0.93) and reconstruction loss (0.13). The final result and performance are both good. Compared to Adam, SGD has lower performance, only achieving an accuracy of 0.87 and a reconstruction loss of 0.68, while the convergence speed is slower. By comparing the results of Adam and SGD optimizers and searching for relevant information, I have made the following discussion: In terms of convergence speed, the Adam optimizer accelerates the learning process by calculating first-order moment estimation and second-order moment estimation of gradients. This adaptive learning rate adjustment method allows Adam to converge faster, so in the early stages of training, Adam usually performs better than SGD; In terms of reconstruction loss: The Adam optimizer can adjust the learning rate based on the gradient of each parameter, which can more effectively optimize complex functions, potentially resulting in lower reconstruction loss and higher accuracy; Overall performance: Although SGD can achieve better accuracy in certain situations, in many tasks, the Adam optimizer is usually more stable and does not require manual adjustment of learning rates. This makes Adam the preferred optimizer for many deep learning applications. But this does not mean that it is the best optimizer in all situations. Choosing the best optimizer often depends on the specific task, model architecture, and dataset. In some cases, such as when the dataset is very large or noisy, other optimizers such as SGD or RMSProp may perform better. In summary, this also provides us with advice and guidance on selecting an optimizer.