

Q₁

$$\varphi'(s) = s \theta'(s) + \theta(s)$$

$$\begin{aligned} \frac{d\theta(s)}{ds} &= [(1+e^{-s})^{-1}]' = -e^{-s} \cdot (-1) \cdot (1+e^{-s})^{-2} \\ &= \frac{e^{-s}}{(1+e^{-s})^2} = (1-\theta(s))\theta(s) \end{aligned}$$

$$\begin{aligned} \Rightarrow \varphi'(s) &= (s+1)\theta(s) - s(\theta(s))^2 \\ &= \frac{1+s}{1+e^{-s}} - \frac{s}{(1+e^{-s})^2} \end{aligned}$$

Q₂

$$(A) V_0 = \left[\frac{1}{3}, \frac{1}{3}, \frac{1}{3} \right]^T$$

$$V_1 = \begin{bmatrix} 0 & 1 & 1/2 \\ 0 & 0 & 1/2 \\ 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} 1/3 \\ 1/3 \\ 1/3 \end{bmatrix} = \begin{bmatrix} 1/2 \\ 1/6 \\ 1/3 \end{bmatrix}$$

$$V_2 = \begin{bmatrix} 0 & 1 & 1/2 \\ 0 & 0 & 1/2 \\ 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} 1/2 \\ 1/6 \\ 1/3 \end{bmatrix} = \begin{bmatrix} 1/3 \\ 1/6 \\ 1/2 \end{bmatrix}$$

$$V_3 = \begin{bmatrix} 0 & 1 & 1/2 \\ 0 & 0 & 1/2 \\ 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} 1/3 \\ 1/6 \\ 1/2 \end{bmatrix} = \begin{bmatrix} 5/12 \\ 1/4 \\ 1/3 \end{bmatrix}$$

$$V_4 = \begin{bmatrix} 0 & 1 & 1/2 \\ 0 & 0 & 1/2 \\ 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} 5/12 \\ 1/4 \\ 1/3 \end{bmatrix} = \begin{bmatrix} 5/12 \\ 1/6 \\ 5/12 \end{bmatrix}$$

$$V_5 = \begin{bmatrix} 0 & 1 & 1/2 \\ 0 & 0 & 1/2 \\ 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} 5/12 \\ 1/6 \\ 5/12 \end{bmatrix} = \begin{bmatrix} 3/8 \\ 5/24 \\ 5/12 \end{bmatrix}$$

(B)

Suppose $V^* = (a, b, c)$, $a+b+c=1$

$$\Rightarrow \begin{bmatrix} a \\ b \\ c \end{bmatrix} = \begin{bmatrix} 0 & 1 & 1/2 \\ 0 & 0 & 1/2 \\ 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} a \\ b \\ c \end{bmatrix} \Rightarrow \begin{array}{l} a = b + \frac{1}{2}c \\ b = \frac{1}{2}c \\ c = a \end{array} \Rightarrow \frac{5}{2}a = 1$$

$$\Rightarrow V^* = \left[\frac{2}{5}, \frac{1}{5}, \frac{2}{5} \right]^T \#$$

3. Let $d^{(l)} = d^{(l)}$

If $L=1$

$\Rightarrow d_0 + d_1 = 100, d_1 = 90 \Rightarrow$ the only possible number of weights is 900

If $L=2$

$\Rightarrow d_0 + d_1 + d_2 = 100 \Rightarrow d_1 + d_2 = 90$ $\left. \begin{array}{l} d_1 = 50 \\ 2d_1 = 100 \end{array} \right\}$

$$W = 10d_1 + d_1d_2 = 10d_1 + d_1(90 - d_1) = -d_1^2 + 100d_1$$

$\Rightarrow (d_1, d_2) = (50, 40)$ has max. number of weights

$\Rightarrow (d_1, d_2) = (1, 89)$ has min. number of weights $\begin{array}{l} = 2500 \\ = 99 \end{array}$

If $L=3$

$$\Rightarrow d_1 + d_2 + d_3 = 90 = g(d_1, d_2, d_3)$$

$$W(d_1, d_2, d_3) = 10d_1 + d_1d_2 + d_2d_3$$

$$= 10(90 - d_2 - d_3) + (90 - d_2 - d_3)d_2 + d_2d_3$$

$$= 900 - 10d_2 - 10d_3 + 90d_2 - d_2^2$$

$$= 900 + 80d_2 - d_2^2 - 10d_3$$

$$= -(d_2 - 40)^2 + 2500 - 10d_3$$

$\Rightarrow (d_1, d_2, d_3) = (49, 40, 1)$ has max. number of weights $= 2490$

$\Rightarrow (d_1, d_2, d_3) = (1, 1, 88)$ has min. number of weights $= 99$

Hence, the maximum number of weights is 2500,

$$L=2, (d_0, d_1, d_2) = (10, 50, 40)$$

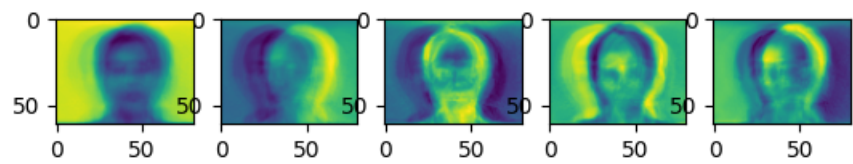
the minimum number of weights is 99,

$$L=2, (d_0, d_1, d_2) = (10, 1, 89) \text{ or}$$

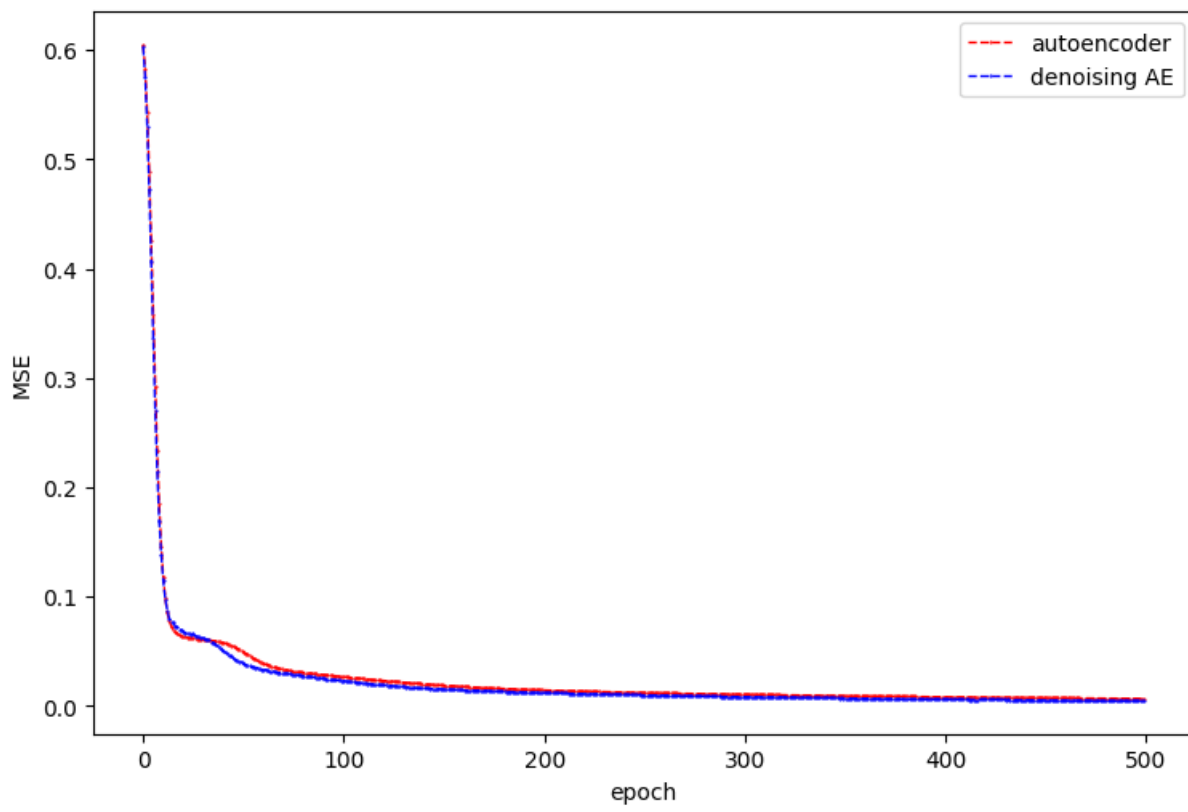
$$L=3, (d_0, d_1, d_2, d_3) = (10, 1, 1, 88) \#$$

Programming Part Report

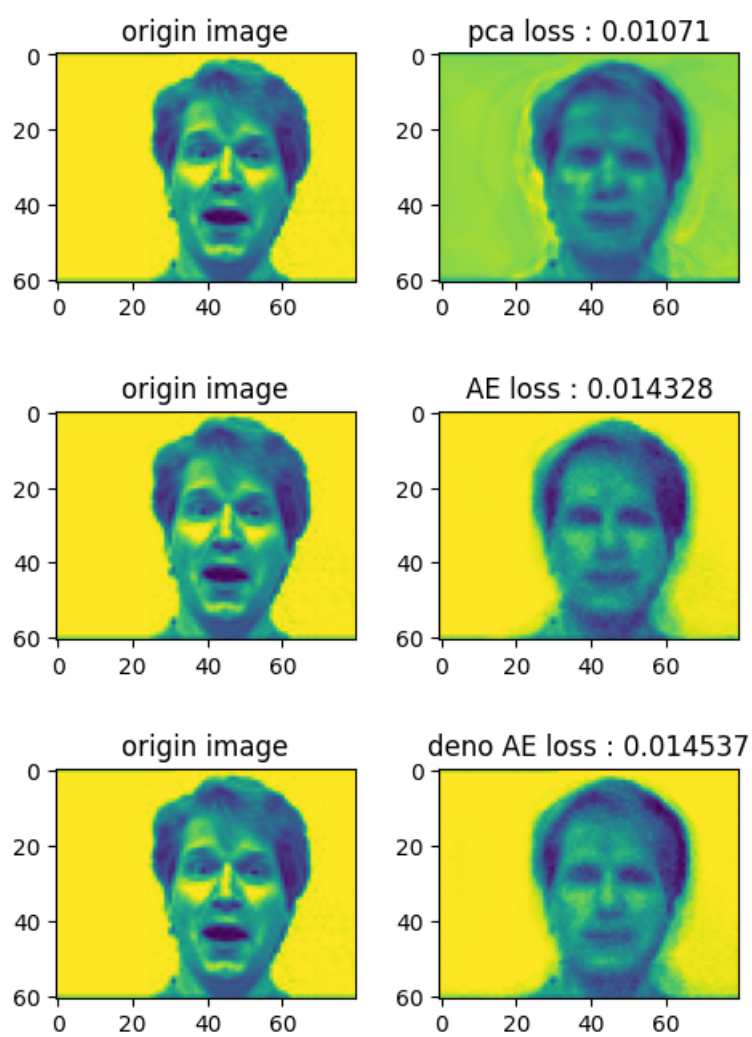
(a)



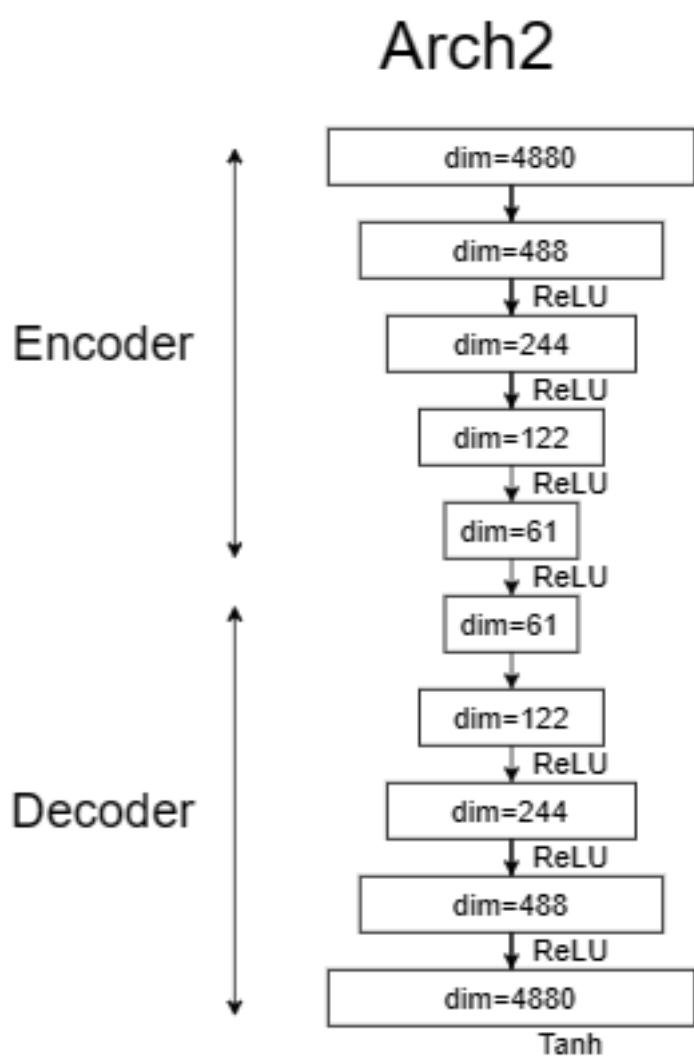
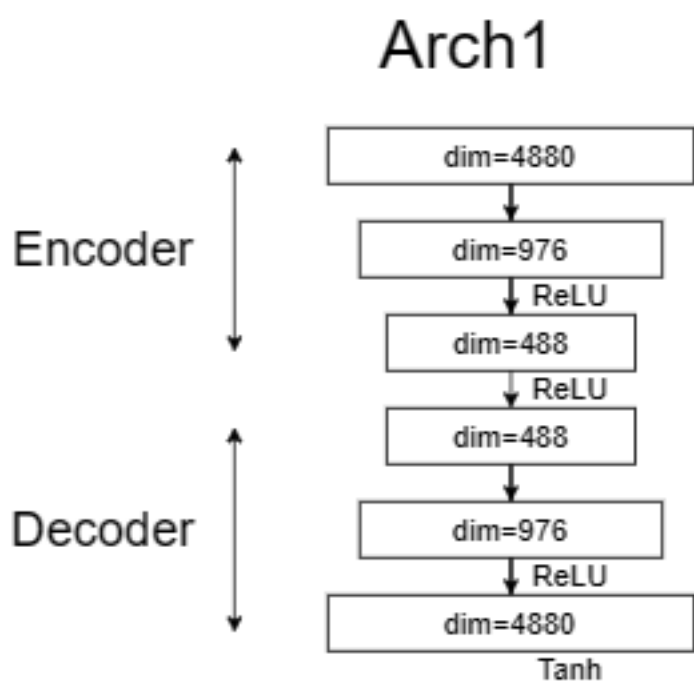
(b)

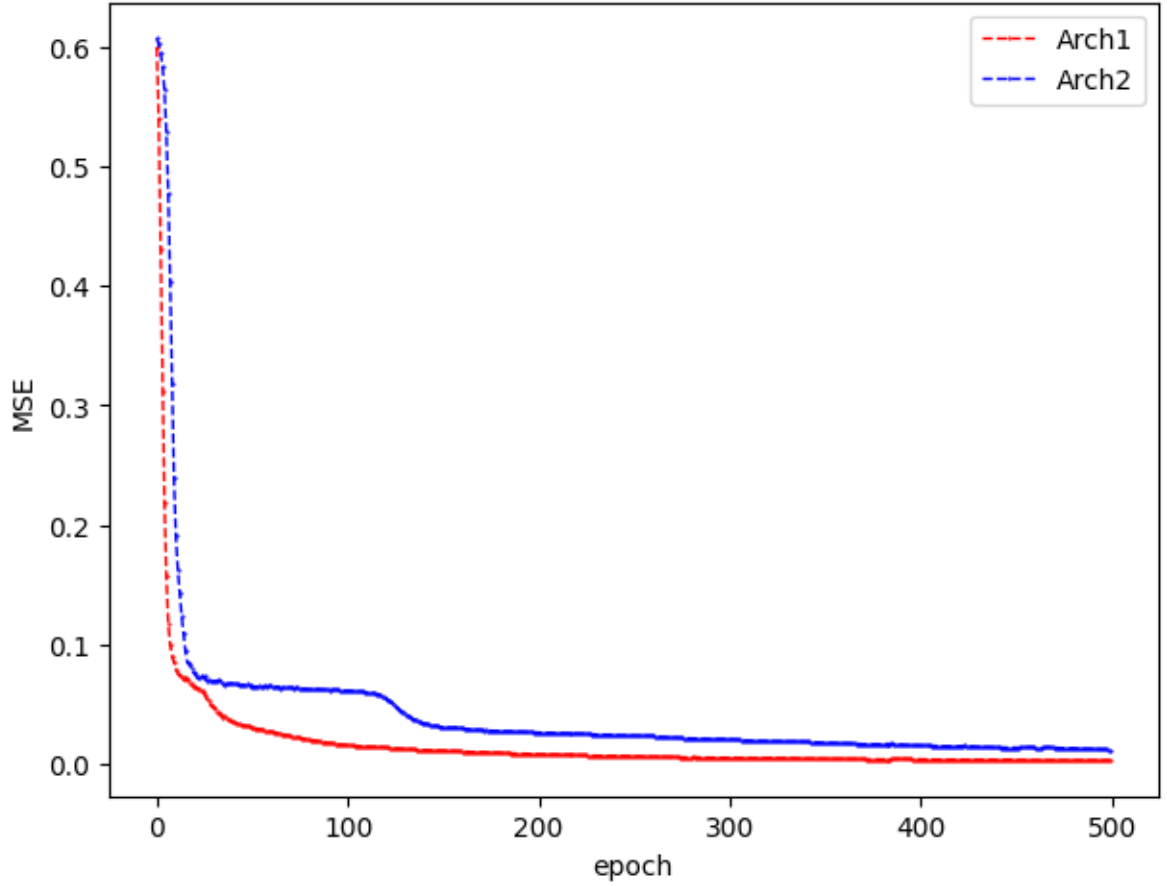


(c)



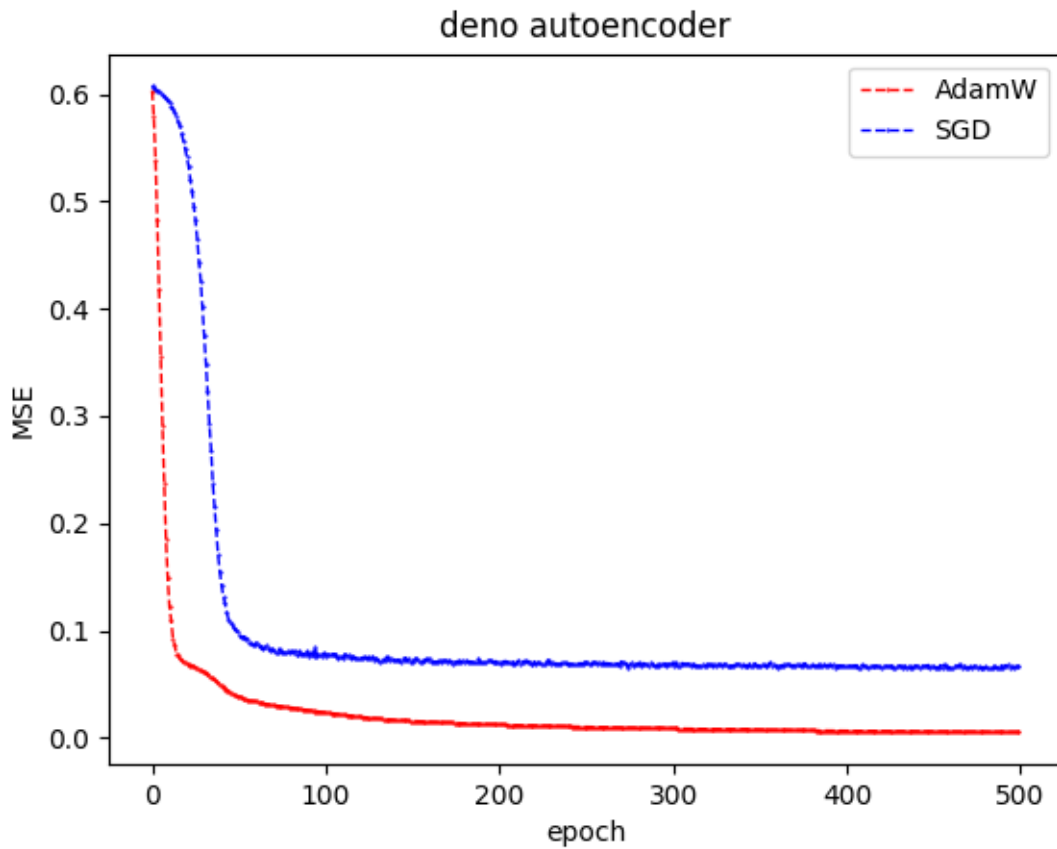
(d)





The two architectures are trained under epoch:500, batchSize:45, learning rates:1e-4 and optimizer:AdamW. Observing the training curve, I find the curve of Arch1 is generally smaller than the curve of Arch2, that is, the MSE loss of Arch1 is smaller than the one of Arch2. Also, when it comes to reconstruction loss, the loss of Arch1 is about 0.01413 while the loss of Arch2 is about 0.01783. Based on these results, I think a shallow but fatter denoising autoencoder may outperform a deeper but thinner one. This might because a fatter network can preserve more information when decoding and a deeper network may be more difficult to optimize.

(e)



I tried two optimizers, one is vanilla SGD and another one is AdamW. Based on the training curve, I find AdamW can converge slightly faster than SGD. However, the overall performance of AdamW defeat the one of SGD totally. The result may because AdamW can dynamically adjust learning rates and use momentum to improve training while vanilla SGD simply train the model.