

① eigen-decomposition  $\Rightarrow U \Lambda U^T$   
 eigenvalues  
 eigenvectors

In SVD,  $A = U \Sigma V^T$ . The values of eigenvalue in  $\Sigma$  should be non-negative ①  
 and diagonal element should satisfy "decreasing order" ②  
 (eigenvalues)

$$\text{ex) } \sigma_1 \geq \sigma_2 \geq \sigma_3 \geq \dots \geq 0$$

We should change  $\Lambda$  to  $\Sigma$  and satisfy the above condition: ①, ②

We also need to modify  $U$  and  $V^T$  slightly. Then we get.

$$A = \begin{bmatrix} U_{11} & -U_{12} & U_{13} \\ U_{21} & -U_{23} & U_{22} \\ U_{31} & -U_{32} & U_{33} \end{bmatrix} \begin{bmatrix} 3 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} U_{11} & U_{21} & U_{31} \\ U_{13} & U_{23} & U_{33} \\ U_{12} & U_{22} & U_{32} \end{bmatrix}$$

②.

If  $H$  is a symmetric matrix then

$$\text{svd}(H) = \underbrace{U \Sigma V^T}_{\text{nxn}} = \underbrace{U \Sigma U^T}_{\text{nxn}}$$

because,  
it is symmetric

$$\max_{A \in \mathbb{R}^{n \times k}} \text{trace}(A^T H A)$$

$$= \max \text{trace}(A^T U \Sigma U^T A)$$

$$= \max \text{trace}((U^T A)^T \Sigma U^T A)$$

$$= \max \text{trace}(Q^T \Sigma Q) \text{ where } Q = U^T A \text{ and } Q \in \mathbb{R}^{n \times k}: Q^T Q = I_k$$

This simplifies into

$$\max_{Q \in \mathbb{R}^{n \times k}: Q^T Q = I} \text{trace}(Q^T \Sigma Q)$$

Here we use the property of 'trace'

for example,  $\text{trace}(ABC) = \text{trace}(BCA)$

We can rewrite the previous equation by the following

$$= \max \text{trace}(\underbrace{I_n Q Q^T}_{\text{rearrange}})$$

$\text{trace}(A)$  is just the summation of diagonal components.

$$= \max \left( \sum_{i=1}^n \lambda_i b_i b_i^T \right) \quad \text{where } b_i \text{ is the diagonal element of } Q$$

Let  $Q^T = A^T U$ . Remember that  $Q = U^T A$  previously.  $A = Q \in \mathbb{R}^{n \times K}$

$$= \max \left( \sum_{i=1}^n \lambda_i a_i^T U U^T a_i \right)$$

$$= \max \left( \sum_{i=1}^n \lambda_i a_i^T a_i \right) \quad \text{because } U U^T = I \quad (\because U \text{ is orthogonal and } H \text{ is symmetric})$$

$$a_i^T a_i = \begin{cases} 0 & i > K \\ 1 & i \leq K \end{cases}$$

we get rid of  
- mat here

$$\textcircled{1} \quad = \sum_{i=1}^K \lambda_i \quad \left( \because \max \left( \sum_{i=1}^K \lambda_i b_i b_i^T + \sum_{i=k+1}^n \lambda_i b_i b_i^T \right) = 0 \right)$$

when  $1 \leq i \leq K$ ,  $M_i = A g_i$ ,  $g_i = g_1 \dots g_K$  ( $g$  is a orthogonal vector)

$$\|A^T U_i\|_2 = 1$$

when  $K < i \leq n$

$$\|A^T U_i\|_2 = 0 \quad \text{this is important}$$

$$\therefore A^* = [U_1, U_2, \dots, U_K] Q \quad \left( \text{and } Q = [g_1, g_2, \dots, g_K]^T \in \mathbb{R}^{K \times K} \right)$$

(①, ②, ③ lead to the fact that)

$$\lambda_1 + \dots + \lambda_K = \max_{A \in \mathbb{R}^{n \times K}, A^T A = I_K} \text{trace}(A^T H A)$$

$Q$  is an arbitrary orthogonal matrix.

## Result and Analysis

Jin Hyun Park, UIN: 633001823

3-a

- In `class Kernel_LR`, bias of `nn.Linear` is set to `False`.
- After hyperparameter tuning with different learning rate, epoch, batch size, and sigma, the best hyperparameters that I have found is the following
  - o Learning rate = 0.01, epoch = 100, batch size = 32, and sigma = 1
- Test accuracy = 0.987 (98.7%)

3-b

- In `class RBF`, bias of `nn.Linear` is set to `False`
- After hyperparameter tuning with different hidden dimension size, learning rate, epoch, batch size, and sigma, the best hyperparameters that I have found is the following
  - o Hidden dim size = 32, learning rate = 0.01, epoch = 100, batch size = 128, and sigma = 1
- Test accuracy = 0.985 (98.5%)

3-c

- In `class FFN`, bias of `nn.Linear` is set to `True`
- After hyperparameter tuning with different hidden dimension size, learning rate, epoch, and batch size, the best hyperparameters that I have is the following:
  - o Hidden dim size = 32, learning rate = 0.01, epoch = 100, batch size = 128
- Test accuracy = 0.989 (98.9%)

Summary of Question 3

	Hidden dimension size	Test accuracy
Kernel LR	# of training samples	98.7%
RBF	32	98.5%
FFN	32	98.9%

- For comparison, the hidden layer dimension size is set to 32 for both RBF and FFN.
- RBF achieved similar test accuracy to Kernel LR (only 0.2% difference). This implies that RBF can perform similarly to Kernel LR with fewer parameters.
- FFN performs slightly better than RBF (only 0.4% difference). This implies that RBF can perform similarly to FFN with fewer parameters. This is possible because we used the K-means algorithm to find centroids.

4-b

# of principal component (p)	Reconstruction Error (RE)
32	129.37
64	85.82
128	45.63

4-d

# of hidden representation (d)	Reconstruction Error (RE)
32	132.70
64	86.63
128	46.60

4-e

p=d	RE of PCA	RE of AE
32	129.37	132.70
64	85.82	86.63
128	45.63	46.60

- The reconstruction error of PCA is slightly better than AE (Auto Encoder) but the differences are very small. This is because we did not use a non-linear activation function. With a linear activation function, PCA and AE basically do the same thing.

4-f

p=d	frobeniu_norm_error(G, W)	frobeniu_norm_error(G <sup>T</sup> .G, W <sup>T</sup> .W)
32	8.11	0.17
64	11.44	0.03
128	15.92	0.02

Note 1: It is represented as  $G \cdot G^T$  and  $W \cdot W^T$  instead of  $G^T \cdot G$  and  $W^T \cdot W$  in the following URL: [1]  
<http://people.tamu.edu/~sji/classes/PCA.pdf>.

Note 2:  $G^T \cdot G$  and  $W^T \cdot W$  are the same as  $G \cdot G^T$  and  $W \cdot W^T$  in my code implementation.

According to [1], a comparison of frobeniu\_norm\_error between  $G^T \cdot G$  and  $W^T \cdot W$  is a better approach than a comparison of frobeniu\_norm\_error between  $G$  and  $W$ . Unsurprisingly, it turned out that the errors were significantly reduced.

4-g

p=d	Shared weights	Non-shared weights
32	132.70	130.28
64	86.63	86.43
128	46.60	46.76

The result of shared weights and non-shared weights is very similar. The decoding part of the network essentially learns to mirror the encoder part for the non-shared weights model.

#### 4-h

- After hyperparameter tuning with different epochs and batch sizes, the best hyperparameters that I have found is the following ( $d$  is set to 64).
  - o Batch size = 64 and epoch = 3000
- I used the following network structure:
  - o (Encoder) layer 1 → ReLU → (Encoder) layer 2 → [Hidden Representation] → (Decoder) layer 3 → ReLU → (Decoder) layer 4
  - o For more details, please see the following two figures

```

W_encoder1 = torch.empty(self.n_features, 2 * self.d_hidden_rep)
self.w_encoder1 = nn.Parameter(nn.init.kaiming_normal_(W_encoder1))

W_encoder2 = torch.empty(2 * self.d_hidden_rep, self.d_hidden_rep)
self.w = nn.Parameter(nn.init.kaiming_normal_(W_encoder2)) # w == w_encoder2

W_decoder1 = torch.empty(self.d_hidden_rep, 2 * self.d_hidden_rep)
self.w_decoder1 = nn.Parameter(nn.init.kaiming_normal_(W_decoder1))

W_decoder2 = torch.empty(2 * self.d_hidden_rep, self.n_features)
self.w_decoder2 = nn.Parameter(nn.init.kaiming_normal_(W_decoder2))

```

```

encoder1 = torch.mm(self.w_encoder1.T, X)
encoder1 = torch.relu(encoder1)
encoder2 = torch.mm(self.w.T, encoder1)

decoder1 = torch.mm(self.w_decoder1.T, encoder2)
decoder1 = torch.relu(decoder1)
decoder2 = torch.mm(self.w_decoder2.T, decoder1)

return decoder2

```

	Reconstruction error ( $d=64$ )
PCA	85.82
Auto Encoder with shared weights	86.63
Auto Encoder with non-shared weights	86.43
Auto encoder with non-linearity and more layers	70.67

The Auto Encoder with non-linearity (= non-linear activation function) and more layers outperformed the others with respect to reconstruction error. This is because it (= the last method, 4-h) becomes **manifold learning** which is an approach to non-linear dimensionality reduction. I strongly believe that with more layers and trying out different combinations of hyperparameters (such as activation functions and batch size), it will be possible to get a lower reconstruction error.

5. Bonus.

$K_{\text{train}}$  : training kernel

$K_{\text{test}}$  : testing kernel

let  $A \in \mathbb{R}^{n \times K}$ .  $A$  = training data,  $K$  = the number of features

"  $B \in \mathbb{R}^{m \times K}$ .  $B$  = testing data,  $K$  = "

$$K_{\text{train}} = A \cdot A^T \in \mathbb{R}^{n \times n}$$

$$\text{Svd}(K_{\text{train}}) = U \Sigma V^T$$

it is symmetric

$$\begin{aligned} &= U \Sigma U^T \quad \text{we can divide } \Sigma \\ &\quad \text{because } \Sigma \text{ is symmetric} \\ &= (U \Sigma^{1/2}) (U \Sigma^{1/2})^T \quad \text{and p.s.d.} \end{aligned}$$

set  $\Rightarrow A = U \Sigma^{1/2}$  --- ① all of diagonal element is  $> 0$  because it is compact SVD

$$\text{In one line, } K_{\text{train}} = A \cdot A^T = U \Sigma^{1/2} (\Sigma^{1/2})^T$$

(NOTE)

" ① and ② implies that the training and the testing data can be recovered or obtained by the kernel."

This implies that we can use  $U$  and  $\Sigma$  for training data. And,  $U$  and  $\Sigma$  can be obtained by svd of  $A^T A$

For testing,

$$K_{\text{test}} = A \cdot B^T$$

$$= U \Sigma^{1/2} B^T$$

$$U^T K_{\text{test}} = U^T U \Sigma^{1/2} B^T$$

$$= \Sigma^{1/2} B^T$$

$$B^T = (\Sigma^{1/2})^{-1} U^T K_{\text{test}}$$

$$B = K_{\text{test}}^T \cdot U \cdot (\Sigma^{1/2})^T \quad \dots \textcircled{2}$$

This implies that we can use  $U$  and  $\Sigma$  for testing data. And,  $U$  and  $\Sigma$  can be obtained by svd of  $A^T A$

∴ Therefore, we can use only the solver and data given above to train a logistic classifier and perform testing.