

# CSCE636 Homework 2

## Solution Keys

October 14, 2019

### 1

Let  $\mathbf{H} = \mathbf{U}\mathbf{\Sigma}\mathbf{U}^T$ , where  $\mathbf{U}\mathbf{U}^T = \mathbf{U}^T\mathbf{U} = \mathbf{I}$  and  $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_n]$ ,  $\mathbf{\Sigma} = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$ . Then we know that  $\mathbf{u}_1, \dots, \mathbf{u}_n$  is a group of orthogonal vectors and  $\mathbf{u}_t^T \mathbf{u}_t = 1, 1 \leq t \leq n$ . Let  $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_k]$ , and  $\mathbf{A}^T \mathbf{U} = [\mathbf{b}_1, \dots, \mathbf{b}_n]$ ,  $\mathbf{b}_t = \mathbf{A}^T \mathbf{u}_t, 1 \leq t \leq n$ , we have

$$\begin{aligned} \text{trace}(\mathbf{A}^T \mathbf{H} \mathbf{A}) &= \text{trace}(\mathbf{A}^T \mathbf{U} \mathbf{\Sigma} \mathbf{U}^T \mathbf{A}) \\ &= \text{trace}\left(\sum_{t=1}^n \lambda_t \mathbf{b}_t \mathbf{b}_t^T\right) \\ &= \sum_{t=1}^n \text{trace}(\lambda_t \mathbf{b}_t \mathbf{b}_t^T) \\ &= \sum_{t=1}^n \lambda_t \text{trace}(\mathbf{b}_t \mathbf{b}_t^T) \\ &= \sum_{t=1}^n \lambda_t \text{trace}(\mathbf{b}_t^T \mathbf{b}_t) \\ &= \sum_{t=1}^n \lambda_t \mathbf{b}_t^T \mathbf{b}_t \\ &= \sum_{t=1}^n \lambda_t \mathbf{u}_t^T \mathbf{A} \mathbf{A}^T \mathbf{u}_t \\ &= \sum_{t=1}^n \lambda_t \|\mathbf{A}^T \mathbf{u}_t\|_2^2 \end{aligned}$$

8 points up to here.

Since  $\mathbf{A}^T \mathbf{A} = \mathbf{I}$  and  $\mathbf{A} \in R^{n \times k}, k < n$ , we know that  $\mathbf{a}_1, \dots, \mathbf{a}_k$  must be a group of orthogonal vectors. Let  $\mathbf{A}' = [\mathbf{a}_{k+1}, \dots, \mathbf{a}_n]$  and  $\mathbf{a}_1, \dots, \mathbf{a}_k, \mathbf{a}_{k+1}, \dots, \mathbf{a}_n$  can form a group of standard orthogonal basis of  $R^n$ . Then  $[\mathbf{A}, \mathbf{A}'] [\mathbf{A}, \mathbf{A}']^T = [\mathbf{A}, \mathbf{A}']^T [\mathbf{A}, \mathbf{A}'] = \mathbf{I}$ . So we have

$$\|\mathbf{A}^T \mathbf{u}_t\|_2^2 \leq \|\mathbf{A}^T \mathbf{u}_t\|_2^2 + \|(\mathbf{A}')^T \mathbf{u}_t\|_2^2 = \|[\mathbf{A}, \mathbf{A}']^T \mathbf{u}_t\|_2^2 = \mathbf{u}_t^T [\mathbf{A}, \mathbf{A}'] [\mathbf{A}, \mathbf{A}']^T \mathbf{u}_t = \mathbf{u}_t^T \mathbf{u}_t = 1$$

The above inequality becomes equality only when  $(\mathbf{A}')^T \mathbf{u}_t = \mathbf{0}$ , thus  $\mathbf{u}_t$  is in the orthogonal complement space of the linear space formed by  $\mathbf{a}_{k+1}, \dots, \mathbf{a}_n$ , and this orthogonal complement space is exactly the column space of  $\mathbf{A}$ . So in this case there must exist a vector  $\mathbf{q}_t$  such that  $\mathbf{u}_t = \mathbf{A} \mathbf{q}_t, \mathbf{q}_t \in R^k$ . In addition,  $\mathbf{u}_t$  is a unitary vector so we have  $\mathbf{1} = \mathbf{u}_t^T \mathbf{u}_t = \mathbf{q}_t^T \mathbf{A}^T \mathbf{A} \mathbf{q}_t = \mathbf{q}_t^T \mathbf{q}_t$ . In conclusion, we always have  $\|\mathbf{A}^T \mathbf{u}_t\|_2^2 \leq 1$  and only when  $\mathbf{u}_t = \mathbf{A} \mathbf{q}_t, \mathbf{q}_t$  is a unit vector, then we will have  $\|\mathbf{A}^T \mathbf{u}_t\|_2^2 = 1$ . This is the complete proof of the first sub-question, 4 points.

Besides, we know that  $\mathbf{A}^T \mathbf{U} \mathbf{U}^T \mathbf{A} = \mathbf{A}^T \mathbf{A} = \mathbf{I}$ ,  $\mathbf{I} \in \mathbf{R}^{k \times k}$  so we know that

$$\begin{aligned} k &= \text{trace}(\mathbf{A}^T \mathbf{U} \mathbf{U}^T \mathbf{A}) \\ &= \text{trace}(\mathbf{U}^T \mathbf{A} \mathbf{A}^T \mathbf{U}) \\ &= \sum_{t=1}^n \mathbf{u}_t^T \mathbf{A} \mathbf{A}^T \mathbf{u}_t \\ &= \sum_{t=1}^n \|\mathbf{A}^T \mathbf{u}_t\|_2^2 \end{aligned}$$

This is the complete proof of the second sub-question, 4 points.

Thus we have

$$\begin{aligned} \text{trace}(\mathbf{A}^T \mathbf{H} \mathbf{A}) &= \sum_{t=1}^n \lambda_t \|\mathbf{A}^T \mathbf{u}_t\|_2^2 \\ &\leq \sum_{t=1}^k \lambda_t \|\mathbf{A}^T \mathbf{u}_t\|_2^2 + \lambda_{k+1} \sum_{t=k+1}^n \|\mathbf{A}^T \mathbf{u}_t\|_2^2 \\ &= \sum_{t=1}^k \lambda_t \|\mathbf{A}^T \mathbf{u}_t\|_2^2 + \lambda_{k+1} (k - \sum_{t=1}^k \|\mathbf{A}^T \mathbf{u}_t\|_2^2) \\ &= \sum_{t=1}^k (\lambda_t - \lambda_{k+1}) \|\mathbf{A}^T \mathbf{u}_t\|_2^2 + k \lambda_{k+1} \\ &\leq \sum_{t=1}^k (\lambda_t - \lambda_{k+1}) + k \lambda_{k+1} \\ &= \sum_{t=1}^k \lambda_t \end{aligned}$$

4 points.

All inequalities become equalities when  $\|\mathbf{A}^T \mathbf{u}_t\|_2 = 1$  for  $1 \leq t \leq k$  and  $\|\mathbf{A}^T \mathbf{u}_t\|_2 = 0$  for  $t > k$ . So we can conclude that  $\mathbf{u}_t = \mathbf{A} \mathbf{q}_t$  for  $1 \leq t \leq k$  and apparently  $\mathbf{q}_1, \dots, \mathbf{q}_k$  are an arbitrary group of unit orthogonal vectors. Thus  $\mathbf{A} = [\mathbf{u}_1, \dots, \mathbf{u}_k] \mathbf{Q}$ ,  $\mathbf{Q} = [\mathbf{q}_1, \dots, \mathbf{q}_k]^T \in \mathbf{R}^{k \times k}$ . So far we can claim that  $\max(\text{trace}(\mathbf{A}^T \mathbf{H} \mathbf{A})) = \sum_{t=1}^k \lambda_t$ , and  $\mathbf{A}^* = [\mathbf{u}_1, \dots, \mathbf{u}_k] \mathbf{Q}$ , where  $\mathbf{Q}$  is an arbitrary orthogonal matrix.

2

Let  $h(z) = \tanh z = \frac{e^z - e^{-z}}{e^z + e^{-z}}$ , we have

$$\frac{dh}{dz} = \frac{(e^z + e^{-z})^2 - (e^z - e^{-z})^2}{(e^z + e^{-z})^2} = 1 - h(z)^2$$

Thus we can get

$$\begin{aligned}
\nabla E_{in}(\mathbf{w}) &= \nabla \frac{1}{N} \sum_{n=1}^N (h(\mathbf{w}^T \mathbf{x}_n) - \mathbf{y}_n)^2 \\
&= \frac{2}{N} \sum_{n=1}^N (h(\mathbf{w}^T \mathbf{x}_n) - \mathbf{y}_n) \nabla (h(\mathbf{w}^T \mathbf{x}_n)) \\
&= \frac{2}{N} \sum_{n=1}^N (h(\mathbf{w}^T \mathbf{x}_n) - \mathbf{y}_n) (1 - h(\mathbf{w}^T \mathbf{x}_n)^2) \nabla (\mathbf{w}^T \mathbf{x}_n) \\
&= \frac{2}{N} \sum_{n=1}^N (h(\mathbf{w}^T \mathbf{x}_n) - \mathbf{y}_n) (1 - h(\mathbf{w}^T \mathbf{x}_n)^2) \mathbf{x}_n \\
&= \frac{2}{N} \sum_{n=1}^N (\tanh(\mathbf{w}^T \mathbf{x}_n) - \mathbf{y}_n) (1 - \tanh^2(\mathbf{w}^T \mathbf{x}_n)) \mathbf{x}_n
\end{aligned}$$

If  $\mathbf{w} \rightarrow \infty$ , then  $\tanh^2(\mathbf{w}^T \mathbf{x}_n)$  will be very close to 1 and  $\nabla E_{in}(\mathbf{w})$  will approach to 0. When this happened, the back propagation will change a very little amount in  $\mathbf{w}$ , which essentially causes gradient vanishing.

### 3

$$\begin{aligned}
s^{(1)} &= W^{(1)} x^{(0)} \begin{bmatrix} 0.1 & 0.3 \\ 0.2 & 0.4 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \end{bmatrix} = \begin{bmatrix} 0.7 \\ 1 \end{bmatrix} \\
x(1) &= \begin{bmatrix} 1 \\ \tanh(s^{(1)}) \end{bmatrix} = \begin{bmatrix} 1 \\ \tanh(0.7) \\ \tanh(1) \end{bmatrix} = \begin{bmatrix} 1 \\ 0.6 \\ 0.76 \end{bmatrix} \\
s^{(2)} &= W^{(2)T} x^{(1)} = [-1.48] \\
x^{(2)} &= \begin{bmatrix} 1 \\ \tanh(s^{(2)}) \end{bmatrix} = \begin{bmatrix} 1 \\ \tanh(-1.48) \end{bmatrix} = \begin{bmatrix} 1 \\ -0.90 \end{bmatrix} \\
x^{(3)} &= s^{(3)} = W^{(3)T} x^{(2)} = [-0.8] \\
\delta^{(3)} &= 2(x^{(3)} - y) = [-3.6] \\
\frac{\partial e}{\partial W^{(3)}} &= x^{(2)} \delta^{(3)} = \begin{bmatrix} 1 \\ -0.9 \end{bmatrix} [-3.6] = \begin{bmatrix} -3.6 \\ 3.24 \end{bmatrix} \\
\delta^{(2)} &= 2 \times (1 - 0.9^2) \times \delta^{(3)} = -1.368 \\
\frac{\partial e}{\partial W^{(2)}} &= x^{(1)} \delta^{(2)} = \begin{bmatrix} 1 \\ 0.6 \\ 0.76 \end{bmatrix} [-1.368] = \begin{bmatrix} -1.368 \\ -0.82 \\ -1.04 \end{bmatrix} \\
\delta^{(1)} &= \begin{bmatrix} 1 - 0.6^2 \\ -3 \times (1 - 0.76^2) \end{bmatrix} \times \delta^{(2)} = \begin{bmatrix} -0.88 \\ 1.73 \end{bmatrix} \\
\frac{\partial e}{\partial W^{(1)}} &= x^{(0)} \delta^{(1)T} = \begin{bmatrix} 1 \\ 2 \end{bmatrix} [-0.88 \ 1.73] = \begin{bmatrix} -0.88 & 1.73 \\ -1.76 & 3.46 \end{bmatrix}
\end{aligned}$$

## 4

(e) For  $d = 32$ , the reconstruction error of PCA and AE are 129.4 and 129.9; for  $d = 64$ , the error are 85.8 and 86.2; for  $d = 128$ , the error are 45.6 and 46.5, respectively.

The reconstruction error of PCA is always very close, but a little bit smaller than that of AE. This is because in this settings, AE is a linear model, and PCA always gives optimal solutions (minimal reconstruction errors here) among all linear models. In addition, AE is a convex problem when weights are shared, so it can achieve minimal reconstruction error in theory. In practice, however, results obtained by gradient-descent usually reach a value near the minimal unless training infinite epochs.

(f) For  $p = d = 32, 64, 128$ , the Frobenius norm of  $\mathbf{W} - \mathbf{G}$  are  $\|\mathbf{W} - \mathbf{G}\|_F = 7.9, 11.3, 16.1$  respectively. Therefore,  $\mathbf{W}$  and  $\mathbf{G}$  are not the same (5 points).

Exploring the relations between  $\mathbf{W}$  and  $\mathbf{G}$  is an open-ended question. Your solutions will be considered as correct if you give either of the below 2:

Sol1: According to note,  $\mathbf{G}$  is a column orthogonal matrix because  $\mathbf{G}$  is component of orthogonal eigenvectors of centered data matrix. Therefore,  $\mathbf{G}^T \mathbf{G} = \mathbf{I}$ . Also, suppose  $\mathbf{W}$  obtained by AE is as optimal as PCA, then  $\mathbf{W}$  can be denotes as  $\mathbf{W} = \mathbf{G}\mathbf{Q}$ , where  $\mathbf{Q}$  is an orthogonal matrix. Therefore,  $\mathbf{W}$  is column orthogonal and we have  $\mathbf{W}^T \mathbf{W} = \mathbf{I}$ . Frobenius norm of  $\mathbf{G}^T \mathbf{G} - \mathbf{W}^T \mathbf{W}$  could be computed experimentally and  $\|\mathbf{G}^T \mathbf{G} - \mathbf{W}^T \mathbf{W}\|_F = 0.007, 0.01, 0.02$  for  $p = d = 32, 64, 128$  respectively, which is in line with our analysis the matrix  $\mathbf{G}$  and  $\mathbf{W}$  all have relation that  $\mathbf{W}^T \mathbf{W} \approx \mathbf{G}^T \mathbf{G} = \mathbf{I}$ .

Sol2: Experimentally show that  $\mathbf{W}\mathbf{W}^T$  and  $\mathbf{G}\mathbf{G}^T$  are close to each other.

Other reasonable solutions will also be accepted as long as experimental justification is provided.

(g) For  $d = 32, 64, 128$ , the autoencoder without shared weights have the reconstruction errors of 129.7, 86.0 and 46.1, respectively. They are very close to the reconstruction errors of the autoencoder with shared weights. The autoencoder with shared weights is just the special case of the autoencoder without shared weights, and the gradient descent process may push the encoder and decoder to have same weights so as to approach the optimal network.

(h) The original data is not linearly separable. By building a nonlinear model using multiple layers and nonlinear activations, you are expected to achieve much lower reconstruction error compared with PCA and above two autoencoders.

This is one network setting.  $batchsize = 256$ ,  $epochs = 2000$ , a 2-layer encoder and a dimension-symmetrical decoder with  $\tanh$  activation functions for the first layers and  $relu$  for the last layer. Under this configuration, the reconstruction error is 47.7, which is much lower than 85.8 obtained by PCA.

## 5

You are expected to achieve the test accuracy better than 97%.

Here is one network setting:  $epochs = 10$ ,  $batchsize = 128$ ,  $num\_hid\_units = 512$ ,  $num\_hid\_units = 3$ ,  $relu$  as activation function. Testing accuracy is 97.7%.