
Exercise 11 of Machine Learning [IN 2064]

Name: Yiman Li
Matr-Nr: 03724352
cooperate with Kejia Chen(03729686)

Problem 1

After adding the data of new user Leslie, and using the SVD decomposition function in MATLAB, we get that:

$$M' = U\Sigma V^T \quad (1)$$

namely:

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 3 & 3 & 3 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 0 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 0 & 0 & 2 & 2 \\ 0 & 3 & 0 & 0 & 4 \end{bmatrix} = \begin{bmatrix} 0.1375 & -0.0170 \\ 0.4126 & -0.0511 \\ 0.5502 & -0.0681 \\ 0.6877 & -0.0852 \\ 0.0417 & 0.5627 \\ 0.0521 & 0.7033 \\ 0.0208 & 0.2813 \\ 0.1739 & 0.3079 \end{bmatrix} \begin{bmatrix} 12.5217 & 0 \\ 0 & 9.9377 \end{bmatrix} \quad (2)$$
$$\begin{bmatrix} 0.5602 & -0.0874 & -0.2886 & -0.3028 & 0.7096 \\ 0.6019 & 0.0056 & 0.5067 & 0.6172 & -0.0050 \end{bmatrix}$$

Here, U connects people to concepts, and $U(8, 1) = 0.1739$ is smaller than some of the other entries in the first column, because Leslie doesn't rate science fiction very high compared to John and Jack. In contrast, Leslie may more likely to see romance movie since $U(8, 2) = 0.3079$ is larger than $U(8, 1)$.

Problem 2

We already know that

$$p(z) = \mathcal{N}(z|\mathbf{0}, I) \quad (3)$$

$$p(\mathbf{x}|z) = \mathcal{N}(\mathbf{W}z + \boldsymbol{\mu}, \Phi) \quad (4)$$

For the original data, we have the corresponding log likelihood as below:

$$\begin{aligned} \ln p(\mathbf{X}|\boldsymbol{\mu}, \mathbf{W}, \Phi) &= \sum_{n=1}^N \ln p(\mathbf{x}_n|\mathbf{W}, \boldsymbol{\mu}, \sigma^2) \\ &= -\frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln|C| - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^T C^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) \end{aligned} \quad (5)$$

where C is defined by

$$C = \mathbf{W}\mathbf{W}^T + \Phi \quad (6)$$

- Setting the derivative of the log likelihood with respect to $\boldsymbol{\mu}$ equal to zero gives the expected result $\boldsymbol{\mu} = \bar{\mathbf{x}}$ where $\bar{\mathbf{x}}$ is the mean value of the data.

- Maximization with respect to \mathbf{W} according to [1], we have

$$\mathbf{W}_{ML} = \mathbf{U}_M(\mathbf{L}_M - \Phi)^{\frac{1}{2}} \mathbf{R} \quad (7)$$

where \mathbf{U}_M is a $D \times M$ matrix whose columns are given by any subset of the eigenvectors, \mathbf{L}_M has elements given by the corresponding eigenvalues λ_i , and \mathbf{R} is an arbitrary $M \times M$ orthogonal matrix.

- Maximization with respect to Φ according to [1], we have

$$\phi_i = \frac{1}{D - M} \sum_{i=M+1}^D \lambda_i \quad (8)$$

where D is the original data dimension and M is the number of principle eigenvectors.

After doing a linear transformation $\mathbf{y} = \mathbf{A}\mathbf{x}$, the mean value and the covariance is now $\mathbf{A}\boldsymbol{\mu}_{ML}$ and $\mathbf{A}\Phi_{ML}\mathbf{A}^T$, which is exactly the corresponding maximum likelihood solution for the transformed data set. And when set the derivative with respect to \mathbf{W} in the transformed space, since the elements in matrix \mathbf{L}_M is scaled by $\mathbf{A}\mathbf{A}^T$ and then take the power of $\frac{1}{2}$, so the new solution is now $\mathbf{A}\mathbf{W}_{ML}$.

And since the matrix \mathbf{A} is orthogonal, so the quantity $\mathbf{W}\mathbf{W}^T$ appears in the covariance matrix \mathbf{C} takes the form

$$\tilde{\mathbf{W}}\tilde{\mathbf{W}}^T = \mathbf{W}\mathbf{A}\mathbf{A}^T\mathbf{W}^T = \mathbf{W}\mathbf{W}^T \quad (9)$$

and hence is independent of \mathbf{A} . We can also get this solution by simple matching patterns the MLE solutions.

Now, if \mathbf{A} is orthogonal and Φ a scaled identity matrix, the model characteristics are also preserved since

$$\mathbf{A}\Phi_x\mathbf{A}^T = \sigma^2\mathbf{I}\mathbf{A}\mathbf{A}^T = \sigma^2\mathbf{I}^2 = \sigma^2\mathbf{I} \quad (10)$$

Problem 3

Suppose we make a matrix transformation formulated by $\mathbf{Y} = \mathbf{X}\mathbf{W} + \mathbf{b}$, where $\mathbf{W} \in \mathbb{R}^{D \times D}$, and $\mathbf{b} \in \mathbb{R}^{N \times D}$, then we have

$$\begin{aligned} \text{Var}(\mathbf{Y}) &= \mathbf{W}^T \left(\frac{1}{N} \mathbf{X}^T \mathbf{X} - \bar{\mathbf{x}}\bar{\mathbf{x}}^T \right) \mathbf{W} \\ &= \mathbf{W}^T (\mathbf{\Gamma} \mathbf{\Lambda} \mathbf{\Gamma}^T) \mathbf{W} \end{aligned} \quad (11)$$

- a) With \mathbf{S} being the identity matrix \mathbf{I} , which means that \mathbf{Y} remains the same as \mathbf{X} after the transformation, so there are still 70% of the variance preserved.
- b) With \mathbf{R} being the row orthogonal matrix, we have $\text{Var}(\mathbf{Y}) = \mathbf{R}^T (\mathbf{\Gamma} \mathbf{\Lambda} \mathbf{\Gamma}^T) \mathbf{R}$, with $\mathbf{\Gamma}' = \mathbf{R}^T \mathbf{\Gamma}$. Here $\mathbf{R}^T \mathbf{\Gamma} \mathbf{\Gamma}^T \mathbf{R} = \mathbf{I}$, and $\mathbf{\Gamma}^T \mathbf{R} \mathbf{R}^T \mathbf{\Gamma} = \mathbf{I}$, which means that $\mathbf{\Gamma}'$ is still a orthonormal matrix, so there are 70% of the variance being preserved.
- c) With $\mathbf{P} = \text{diag}(5, -5, \dots, 5, -5)$, then the new covariance matrix is now scaled by a factor 25, which will not change what the matrix looks like, so it still preserve 70% of the variance.
- d) With $\mathbf{Q} = \text{diag}(1, 2, 3, \dots, D-1, D)$, which means the new covariance matrix is now $\text{diag}(\lambda_1, 4\lambda_2, 9\lambda_3, \dots)$, so we have no idea what the order of principle components in new matrix will be, which means the preserved variance can not be told without additional information.
- e) Since adding the same value to the original matrix will not change the covatriance matrix, so we still preserve 70% of original variance.
- f) Since $\text{rank}(\mathbf{A}) = 5$, which means that the data lies in a 5-dimensional space or subspace, so the top 5 principle components captures all the variance, that is to say 100%.

Problem 4

a) The mean of the N points is shown as below:

$$\bar{x} = \frac{1}{N} \cdot \mathbf{X}^T \cdot \mathbf{1}_N = \begin{bmatrix} 2 & 1 & 1 \end{bmatrix} \quad (12)$$

Shift the points by their mean \bar{x} , we get

$$\tilde{\mathbf{X}} = \mathbf{X} - \bar{x} = \begin{bmatrix} 2 & 2 & 1 \\ 0 & 0 & -3 \\ 2 & -2 & 1 \\ -4 & 0 & 1 \end{bmatrix} \quad (13)$$

So the variance is

$$\Sigma_{\mathbf{X}} = \text{Var}(\tilde{\mathbf{X}}) = \begin{bmatrix} 6 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 3 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \times \begin{bmatrix} 6 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 3 \end{bmatrix} \times \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad (14)$$

- The first principle component is $[1 \ 0 \ 0]^T$, with variance value 6;
- The second principle component is $[0 \ 0 \ 1]^T$, with variance value 3;
- The third principle component is $[0 \ 1 \ 0]^T$, with variance value 2.

b) According to the function 14, we can see that the top-2 principle components are 6 and 3, so we obtain the projected data \mathbf{Y} as below:

$$\mathbf{Y} = \begin{bmatrix} 2 & 2 & 1 \\ 0 & 0 & -3 \\ 2 & -2 & 1 \\ -4 & 0 & 1 \end{bmatrix} \times \begin{bmatrix} 1 & 0 \\ 0 & 0 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 2 & 1 \\ 0 & -3 \\ 2 & 1 \\ -4 & 1 \end{bmatrix} \quad (15)$$

c) When new data point have exactly the same PCA performance, which means that the covariance remains the same after adding this new data point (scaling factor will not change what the covariance matrix looks like), so we can say that $x_5 = [2 \ 1 \ 1]$

Problem 5

To be seen in the end.

References

- [1] Christopher M. Bishop. *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Springer, 1 edition, 2007.