

Homework 0 - Introduction to Probabilistic Graphical Models

kipngeno koech - bkoech

April 5, 2025

1 Distributed SGD [20 points]

Consider that we have a system of m worker nodes and a parameter server performing distributed SGD (stochastic gradient descent). In each iteration, every worker node receives the model from the parameter server, computes one gradient step of the objective function locally using its local data, and sends the gradient to the parameter server. The parameter server does the aggregation of gradients using either synchronous SGD or asynchronous SGD.

The gradient calculation time X_i taken by each node i follows the exponential distribution with rate $\lambda = 2$, which has the following probability density function (PDF):

$$f_X(x) = \begin{cases} \lambda e^{-\lambda x} & \text{if } x \geq 0 \\ 0 & \text{else} \end{cases} \quad (1)$$

Answer the following questions and make sure to explain your answers:

1. What is the cumulative distribution function (CDF) of $f_X(x)$, i.e., $F_X(x)$?

We find the CDF by integrating the PDF:

$$F_X(x) = \int_{-\infty}^x f_X(t) dt \quad (2)$$

$$= \int_0^x \lambda e^{-\lambda t} dt \quad (3)$$

$$= [-e^{-\lambda t}]_0^x \quad (4)$$

$$= -e^{-\lambda x} + 1 \quad (5)$$

$$= 1 - e^{-\lambda x} \quad (6)$$

integrating from 0 to x since the PDF is 0 for $x < 0$.

$$F_X(x) = \begin{cases} 0 & \text{if } x < 0 \\ 1 - e^{-\lambda x} & \text{if } x \geq 0 \end{cases} \quad (7)$$

2. Define $X_{m:m}$ as the maximum of m i.i.d. (independently and identically distributed) instances X_1, \dots, X_m following the distribution X . What is the CDF of $X_{m:m}$, and what is the expected value $E[X_{m:m}]$?
3. Define $X_{1:m}$ as the minimum of m i.i.d instances X_1, \dots, X_m following the distribution X . What is the CDF of $X_{1:m}$, and what is the expected value $E[X_{1:m}]$?
4. In this sub-problem, we will simulate and compare the expected runtime per iteration of synchronous SGD and asynchronous SGD for different values of m . The time for each worker node to finish one gradient computation is exponentially distributed as given in part (a) with $\lambda = 2$, and it is i.i.d. across workers and iterations. Assume there is no communication delay.

Simulate 5000 iterations of training using Python for different values of m ranging from 1 to 20, and obtain the average runtime per iteration. Make a comparative plot of the average runtimes per iteration of synchronous and asynchronous SGD versus m . Explain the trends observed in the plot in 1-2 sentences. You may use packages inside `numpy.random` to draw random samples from the exponential distribution. Attach your plot and code in PDF format to the end of your homework.

5. Write down the theoretical expressions for the expected runtimes per iteration of synchronous and asynchronous SGD in terms of m and λ (Hint: You can use the expressions derived in parts (b) and (c)). On the figure generated in part (d), also plot the theoretical expected runtimes versus m . Check whether the theoretical and simulated values align.

2 K-means [20 points]

Given a set of data points $\{\mathbf{x}_n\}_{n=1}^N$, k-means clustering minimizes the following distortion measure (also called the “objective” or “clustering cost”):

$$D = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|_2^2 \quad (8)$$

where $\boldsymbol{\mu}_k$ is the prototype of the k -th cluster and r_{nk} is a binary indicator variable. If \mathbf{x}_n is assigned to the cluster k , r_{nk} is 1, and otherwise r_{nk} is 0. For each cluster, $\boldsymbol{\mu}_k$ is the prototype representative for all the data points assigned to that cluster.

1. In lecture, we stated but did not prove that $\boldsymbol{\mu}_k$ is the mean of all points associated with the k th cluster, thus motivating the name of the algorithm. You will now prove this statement. Assuming all r_{nk} are known (i.e., assuming you know the cluster assignments of all N data points), show that the objective D is minimized when each $\boldsymbol{\mu}_k$ is chosen as the mean of all data points assigned to cluster k , for any k . This justifies the iterative procedure of k-means¹.

Let us denote the set of data points assigned to cluster k as $C_k = \{n : r_{nk} = 1\}$. The objective function can be rewritten as:

$$D = \sum_{k=1}^K \sum_{n \in C_k} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|_2^2 \quad (9)$$

$$= \sum_{k=1}^K \sum_{n \in C_k} (\|\mathbf{x}_n\|_2^2 - 2\mathbf{x}_n^T \boldsymbol{\mu}_k + \|\boldsymbol{\mu}_k\|_2^2) \quad (10)$$

$$= \sum_{k=1}^K \left(\sum_{n \in C_k} \|\mathbf{x}_n\|_2^2 - 2\boldsymbol{\mu}_k^T \sum_{n \in C_k} \mathbf{x}_n + |C_k| \|\boldsymbol{\mu}_k\|_2^2 \right) \quad (11)$$

where $|C_k|$ is the number of data points assigned to cluster k . Now, we can differentiate D with respect to $\boldsymbol{\mu}_k$ and set it to zero to find the optimal $\boldsymbol{\mu}_k$:

$$\frac{\partial D}{\partial \boldsymbol{\mu}_k} = -2 \sum_{n \in C_k} \mathbf{x}_n + 2|C_k| \boldsymbol{\mu}_k = 0 \quad (12)$$

$$\sum_{n \in C_k} \mathbf{x}_n = |C_k| \boldsymbol{\mu}_k \quad (13)$$

$$\boldsymbol{\mu}_k = \frac{1}{|C_k|} \sum_{n \in C_k} \mathbf{x}_n \quad (14)$$

Thus, the optimal $\boldsymbol{\mu}_k$ is the mean of all data points assigned to cluster k .

2. As discussed in lecture, sometimes we wish to scale each feature in order to ensure that “larger” features do not dominate the clustering. Suppose that each data point \mathbf{x}_n is a d -dimensional feature vector and that we scale the j th feature by a factor $w_j > 0$. Letting \mathbf{W} denote a $d \times d$ diagonal matrix with the j ’th diagonal entry being w_j , $j = 1, 2, \dots, d$, we can write our transformed features as $\mathbf{x}' = \mathbf{W}\mathbf{x}$.

Suppose we fix the r_{nk} , i.e., we take the assignment of data points \mathbf{x}_n to clusters k as given. Our goal is then to find the cluster centers $\boldsymbol{\mu}_k$ that minimize the distortion measure

$$D = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|\mathbf{W}\mathbf{x}_n - \boldsymbol{\mu}_k\|_2^2. \quad (15)$$

Show that the cluster centers $\{\boldsymbol{\mu}_k\}$ that do so are given by $\boldsymbol{\mu}_k = \frac{1}{\sum_{n=1}^N r_{nk}} \mathbf{W} \sum_{n=1}^N r_{nk} \mathbf{x}_n$.

Let us denote the set of data points assigned to cluster k as $C_k = \{n : r_{nk} = 1\}$. The objective function can be rewritten as:

$$D = \sum_{k=1}^K \sum_{n \in C_k} \|\mathbf{W}\mathbf{x}_n - \boldsymbol{\mu}_k\|_2^2 \quad (16)$$

$$= \sum_{k=1}^K \sum_{n \in C_k} (\|\mathbf{W}\mathbf{x}_n\|_2^2 - 2(\mathbf{W}\mathbf{x}_n)^T \boldsymbol{\mu}_k + \|\boldsymbol{\mu}_k\|_2^2) \quad (17)$$

$$= \sum_{k=1}^K \left(\sum_{n \in C_k} \|\mathbf{W}\mathbf{x}_n\|_2^2 - 2\boldsymbol{\mu}_k^T \sum_{n \in C_k} \mathbf{W}\mathbf{x}_n + |C_k| \|\boldsymbol{\mu}_k\|_2^2 \right) \quad (18)$$

where $|C_k|$ is the number of data points assigned to cluster k . Now, we can differentiate D with respect to $\boldsymbol{\mu}_k$ and set it to zero to find the optimal $\boldsymbol{\mu}_k$:

$$\frac{\partial D}{\partial \boldsymbol{\mu}_k} = -2 \sum_{n \in C_k} \mathbf{W}\mathbf{x}_n + 2|C_k| \boldsymbol{\mu}_k = 0 \quad (19)$$

$$\sum_{n \in C_k} \mathbf{W}\mathbf{x}_n = |C_k| \boldsymbol{\mu}_k \quad (20)$$

$$\boldsymbol{\mu}_k = \frac{1}{|C_k|} \sum_{n \in C_k} \mathbf{W}\mathbf{x}_n \quad (21)$$

¹More rigorously, one would also need to show that if all $\boldsymbol{\mu}_k$ are known, then r_{nk} can be computed by assigning \mathbf{x}_n to the nearest $\boldsymbol{\mu}_k$. You are not required to do so.

We can rewrite the sum over C_k as a sum over all points using the indicator variables r_{nk} , and note that $|C_k| = \sum_{n=1}^N r_{nk}$:

$$\boldsymbol{\mu}_k = \frac{1}{\sum_{n=1}^N r_{nk}} \sum_{n=1}^N r_{nk} \mathbf{W} \mathbf{x}_n \quad (22)$$

Since \mathbf{W} is a diagonal matrix that doesn't depend on the summation index n , we can factor it out of the summation:

$$\boldsymbol{\mu}_k = \frac{1}{\sum_{n=1}^N r_{nk}} \mathbf{W} \sum_{n=1}^N r_{nk} \mathbf{x}_n \quad (23)$$

Thus, the cluster centers $\boldsymbol{\mu}_k$ that minimize the distortion measure are given by:

$$\boldsymbol{\mu}_k = \frac{1}{\sum_{n=1}^N r_{nk}} \mathbf{W} \sum_{n=1}^N r_{nk} \mathbf{x}_n \quad (24)$$

3 3-Dimensional Principal Component Analysis [20 points]

In this problem, we will perform PCA on 3-dimensional data step by step. We are given three data points:

$$\mathbf{x}_1 = [0, -1, -2] \quad (25)$$

$$\mathbf{x}_2 = [1, 1, 1] \quad (26)$$

$$\mathbf{x}_3 = [2, 0, 1] \quad (27)$$

and we want to find 2 principal components of the given data.

1. First, find the covariance matrix $\mathbf{C}_X = \mathbf{X}^T \mathbf{X}$ where $\mathbf{X} = \begin{bmatrix} \mathbf{x}_1 - \bar{\mathbf{x}} \\ \mathbf{x}_2 - \bar{\mathbf{x}} \\ \mathbf{x}_3 - \bar{\mathbf{x}} \end{bmatrix}$, where $\bar{\mathbf{x}} = \frac{1}{3}(\mathbf{x}_1 + \mathbf{x}_2 + \mathbf{x}_3)$ is the mean of the data samples.

Then, find the eigenvalues and the corresponding eigenvectors of \mathbf{C}_X . Feel free to use any numerical analysis program such as numpy, e.g., `numpy.linalg.eig` can be useful. However, you should explain what you inputted into this program.

Finding the mean $\bar{\mathbf{x}}$:

$$\bar{\mathbf{x}} = \frac{1}{3}(\mathbf{x}_1 + \mathbf{x}_2 + \mathbf{x}_3) \quad (28)$$

$$= \frac{1}{3}([0, -1, -2] + [1, 1, 1] + [2, 0, 1]) \quad (29)$$

$$= \frac{1}{3}([3, 0, 0]) \quad (30)$$

$$= [1, 0, 0] \quad (31)$$

let us find \mathbf{X} :

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1 - \bar{\mathbf{x}} \\ \mathbf{x}_2 - \bar{\mathbf{x}} \\ \mathbf{x}_3 - \bar{\mathbf{x}} \end{bmatrix} \quad (32)$$

$$= \begin{bmatrix} [0, -1, -2] - [1, 0, 0] \\ [1, 1, 1] - [1, 0, 0] \\ [2, 0, 1] - [1, 0, 0] \end{bmatrix} \quad (33)$$

$$= \begin{bmatrix} [-1, -1, -2] \\ [0, 1, 1] \\ [1, 0, 1] \end{bmatrix} \quad (34)$$

Now, we can find the covariance matrix \mathbf{C}_X :

$$\mathbf{C}_X = \mathbf{X}^T \mathbf{X} \quad (35)$$

$$= \begin{bmatrix} [-1, 0, 1] \\ [-1, 1, 0] \\ [-2, 1, 1] \end{bmatrix} \begin{bmatrix} [-1, -1, -2] \\ [0, 1, 1] \\ [1, 0, 1] \end{bmatrix} \quad (36)$$

$$= \begin{bmatrix} 2 & 1 & 3 \\ 1 & 2 & 3 \\ 3 & 3 & 6 \end{bmatrix} \quad (37)$$

Now, we can find the eigenvalues and eigenvectors of \mathbf{C}_X using numpy:

```
import numpy as np
C_X = np.array([[2, 1, 3], [1, 2, 3], [3, 3, 6]])
eigenvalues, eigenvectors = np.linalg.eig(C_X)
print("Eigenvalues:", eigenvalues)
print("Eigenvectors:", eigenvectors)
```

The output will give us the eigenvalues and eigenvectors of the covariance matrix. The eigenvalues are:

$$\lambda_1 = 9.00 \quad (38)$$

$$\lambda_2 = 1.00 \quad (39)$$

$$\lambda_3 = 0.00 \quad (40)$$

The corresponding eigenvectors are:

$$\mathbf{u}_1 = \begin{pmatrix} 1 \\ 1 \\ 2 \end{pmatrix} \quad \mathbf{u}_2 = \begin{pmatrix} -1 \\ 1 \\ 0 \end{pmatrix} \quad \mathbf{u}_3 = \begin{pmatrix} -1 \\ -1 \\ 1 \end{pmatrix} \quad (41)$$

I used the numpy function `numpy.linalg.eig` to compute the eigenvalues and eigenvectors of the covariance matrix. The input to this function was the covariance matrix \mathbf{C}_X that we computed above.

2. Using the result above, find the first two principal components of the given data.

The first two principal components are the eigenvectors corresponding to the two largest eigenvalues. In this case, the first two principal components are:

$$\mathbf{u}_1 = \begin{pmatrix} 1 \\ 1 \\ 2 \end{pmatrix} \quad (\text{corresponding to } \lambda_1 = 9.00) \quad (42)$$

$$\mathbf{u}_2 = \begin{pmatrix} -1 \\ 1 \\ 0 \end{pmatrix} \quad (\text{corresponding to } \lambda_2 = 1.00) \quad (43)$$

3. Now we want to represent the data $\mathbf{x}_1, \dots, \mathbf{x}_3$ using a 2-dimensional subspace instead of a 3-dimensional one. PCA gives us the 2-D plane which minimizes the difference between the original data and the data projected to the 2-dimensional plane. In other words, \mathbf{x}_i can be approximated as:

$$\tilde{\mathbf{x}}_i = a_{i1}\mathbf{u}_1 + a_{i2}\mathbf{u}_2 + \bar{\mathbf{x}}, \quad (44)$$

where \mathbf{u}_1 and \mathbf{u}_2 are the principal components we found in 3.b. Figure 1 gives an example of what this might look like.

Figure 1: Example of 2-D plane spanned by the first two principal components.

Find a_{i1}, a_{i2} for $i = 1, 2, 3$. Then, find the $\tilde{\mathbf{x}}_i$'s and the difference between $\tilde{\mathbf{x}}_i$ and \mathbf{x}_i , i.e., $\|\tilde{\mathbf{x}}_i - \mathbf{x}_i\|_2$ for $i = 1, 2, 3$. (Again, feel free to use any numerical analysis program to get the final answer. But, show your calculation process.)

To find a_{i1} and a_{i2} , we can project the original data points onto the principal components:

Let us project \mathbf{x}_1 onto the first two principal components:

$$\mathbf{x}_1 = a_{11}\mathbf{u}_1 + a_{12}\mathbf{u}_2 + \bar{\mathbf{x}} \quad (45)$$

$$\mathbf{x}_1 - \bar{\mathbf{x}} = a_{11}\mathbf{u}_1 + a_{12}\mathbf{u}_2 \quad (46)$$

$$[-1, -1, -2] = a_{11} \begin{pmatrix} 1 \\ 1 \\ 2 \end{pmatrix} + a_{12} \begin{pmatrix} -1 \\ 1 \\ 0 \end{pmatrix} \quad (47)$$

$$[-1, -1, -2] = \begin{pmatrix} a_{11} - a_{12} \\ a_{11} + a_{12} \\ 2a_{11} \end{pmatrix} \quad (48)$$

Now, we can set up a system of equations:

$$-1 = a_{11} - a_{12} \quad (49)$$

$$-1 = a_{11} + a_{12} \quad (50)$$

$$-2 = 2a_{11} \quad (51)$$

From the third equation, we can find a_{11} :

$$2a_{11} = -2 \quad (52)$$

$$a_{11} = -1 \quad (53)$$

Now, we can substitute a_{11} into the first two equations to find a_{12} :

$$-1 = -1 - a_{12} \quad (54)$$

$$-1 = -1 + a_{12} \quad (55)$$

$$a_{12} = 0 \quad (56)$$

So, for \mathbf{x}_1 , we have:

$$a_{11} = -1 \quad (57)$$

$$a_{12} = 0 \quad (58)$$

Now, we can find $\tilde{\mathbf{x}}_1$:

$$\tilde{\mathbf{x}}_1 = a_{11}\mathbf{u}_1 + a_{12}\mathbf{u}_2 + \bar{\mathbf{x}} \quad (59)$$

$$= -1 \begin{pmatrix} 1 \\ 1 \\ 2 \end{pmatrix} + 0 \begin{pmatrix} -1 \\ 1 \\ 0 \end{pmatrix} + [1, 0, 0] \quad (60)$$

$$= [-1, -1, -2] + [1, 0, 0] \quad (61)$$

$$= [0, -1, -2] \quad (62)$$

Now, we can find the difference between $\tilde{\mathbf{x}}_1$ and \mathbf{x}_1 :

$$\|\tilde{\mathbf{x}}_1 - \mathbf{x}_1\|_2 = \|[0, -1, -2] - [-1, -1, -2]\|_2 \quad (63)$$

$$= \|[0, 0, 0]\|_2 \quad (64)$$

$$= 0 \quad (65)$$

Now, we can repeat the process for \mathbf{x}_2 :

$$\mathbf{x}_2 = a_{21}\mathbf{u}_1 + a_{22}\mathbf{u}_2 + \bar{\mathbf{x}} \quad (66)$$

$$\mathbf{x}_2 - \bar{\mathbf{x}} = a_{21}\mathbf{u}_1 + a_{22}\mathbf{u}_2 \quad (67)$$

$$[0, 1, 1] = a_{21} \begin{pmatrix} 1 \\ 1 \\ 2 \end{pmatrix} + a_{22} \begin{pmatrix} -1 \\ 1 \\ 0 \end{pmatrix} \quad (68)$$

$$[0, 1, 1] = \begin{pmatrix} a_{21} - a_{22} \\ a_{21} + a_{22} \\ 2a_{21} \end{pmatrix} \quad (69)$$

Now, we can set up a system of equations:

$$0 = a_{21} - a_{22} \quad (70)$$

$$1 = a_{21} + a_{22} \quad (71)$$

$$1 = 2a_{21} \quad (72)$$

From the third equation, we can find a_{21} :

$$2a_{21} = 1 \quad (73)$$

$$a_{21} = \frac{1}{2} \quad (74)$$

Now, we can substitute a_{21} into the first two equations to find a_{22} :

$$0 = \frac{1}{2} - a_{22} \quad (75)$$

$$0 = \frac{1}{2} + a_{22} \quad (76)$$

$$a_{22} = \frac{1}{2} \quad (77)$$

So, for \mathbf{x}_2 , we have:

$$a_{21} = \frac{1}{2} \quad (78)$$

$$a_{22} = \frac{1}{2} \quad (79)$$

Now, we can find $\tilde{\mathbf{x}}_2$:

$$\tilde{\mathbf{x}}_2 = a_{21}\mathbf{u}_1 + a_{22}\mathbf{u}_2 + \bar{\mathbf{x}} \quad (80)$$

$$= \frac{1}{2} \begin{pmatrix} 1 \\ 1 \\ 2 \end{pmatrix} + \frac{1}{2} \begin{pmatrix} -1 \\ 1 \\ 0 \end{pmatrix} + [1, 0, 0] \quad (81)$$

$$= \begin{pmatrix} \frac{1}{2} - \frac{1}{2} \\ \frac{1}{2} + \frac{1}{2} \\ \frac{1}{2} \end{pmatrix} + [1, 0, 0] \quad (82)$$

$$= [0, 1, 1] + [1, 0, 0] \quad (83)$$

$$= [1, 1, 1] \quad (84)$$

Now, we can find the difference between $\tilde{\mathbf{x}}_2$ and \mathbf{x}_2 :

$$\|\tilde{\mathbf{x}}_2 - \mathbf{x}_2\|_2 = \|[1, 1, 1] - [1, 1, 1]\|_2 \quad (85)$$

$$= \|[0, 0, 0]\|_2 \quad (86)$$

$$= 0 \quad (87)$$

Now, we can repeat the process for \mathbf{x}_3 :

$$\mathbf{x}_3 = a_{31}\mathbf{u}_1 + a_{32}\mathbf{u}_2 + \bar{\mathbf{x}} \quad (88)$$

$$\mathbf{x}_3 - \bar{\mathbf{x}} = a_{31}\mathbf{u}_1 + a_{32}\mathbf{u}_2 \quad (89)$$

$$[1, 0, 1] = a_{31} \begin{pmatrix} 1 \\ 1 \\ 2 \end{pmatrix} + a_{32} \begin{pmatrix} -1 \\ 1 \\ 0 \end{pmatrix} \quad (90)$$

$$[1, 0, 1] = \begin{pmatrix} a_{31} - a_{32} \\ a_{31} + a_{32} \\ 2a_{31} \end{pmatrix} \quad (91)$$

Now, we can set up a system of equations:

$$1 = a_{31} - a_{32} \quad (92)$$

$$0 = a_{31} + a_{32} \quad (93)$$

$$1 = 2a_{31} \quad (94)$$

From the third equation, we can find a_{31} :

$$2a_{31} = 1 \quad (95)$$

$$a_{31} = \frac{1}{2} \quad (96)$$

Now, we can substitute a_{31} into the first two equations to find a_{32} :

$$1 = \frac{1}{2} - a_{32} \quad (97)$$

$$0 = \frac{1}{2} + a_{32} \quad (98)$$

$$a_{32} = -\frac{1}{2} \quad (99)$$

So, for \mathbf{x}_3 , we have:

$$a_{31} = \frac{1}{2} \quad (100)$$

$$a_{32} = -\frac{1}{2} \quad (101)$$

Now, we can find $\tilde{\mathbf{x}}_3$:

$$\tilde{\mathbf{x}}_3 = a_{31}\mathbf{u}_1 + a_{32}\mathbf{u}_2 + \bar{\mathbf{x}} \quad (102)$$

$$= \frac{1}{2} \begin{pmatrix} 1 \\ 1 \\ 2 \end{pmatrix} - \frac{1}{2} \begin{pmatrix} -1 \\ 1 \\ 0 \end{pmatrix} + [1, 0, 0] \quad (103)$$

$$= \begin{pmatrix} \frac{1}{2} + \frac{1}{2} \\ \frac{1}{2} - \frac{1}{2} \\ \frac{1}{2} \end{pmatrix} + [1, 0, 0] \quad (104)$$

$$= [1, 0, 1] + [1, 0, 0] \quad (105)$$

$$= [2, 0, 1] \quad (106)$$

Now, we can find the difference between $\tilde{\mathbf{x}}_3$ and \mathbf{x}_3 :

$$\|\tilde{\mathbf{x}}_3 - \mathbf{x}_3\|_2 = \|[2, 0, 1] - [2, 0, 1]\|_2 \quad (107)$$

$$= \|[0, 0, 0]\|_2 \quad (108)$$

$$= 0 \quad (109)$$

So, we have:

$$\tilde{\mathbf{x}}_1 = [0, -1, -2] \quad \|\tilde{\mathbf{x}}_1 - \mathbf{x}_1\|_2 = 0 \quad (110)$$

$$\tilde{\mathbf{x}}_2 = [1, 1, 1] \quad \|\tilde{\mathbf{x}}_2 - \mathbf{x}_2\|_2 = 0 \quad (111)$$

$$\tilde{\mathbf{x}}_3 = [2, 0, 1] \quad \|\tilde{\mathbf{x}}_3 - \mathbf{x}_3\|_2 = 0 \quad (112)$$

Thus, the differences between the projected data points and the original data points are all zero. This means that the PCA projection perfectly represents the original data in the 2-dimensional subspace spanned by the first two principal components.