

Your latest: 100% • Your highest: 100% • To pass you need at least 80%. We keep your highest score.

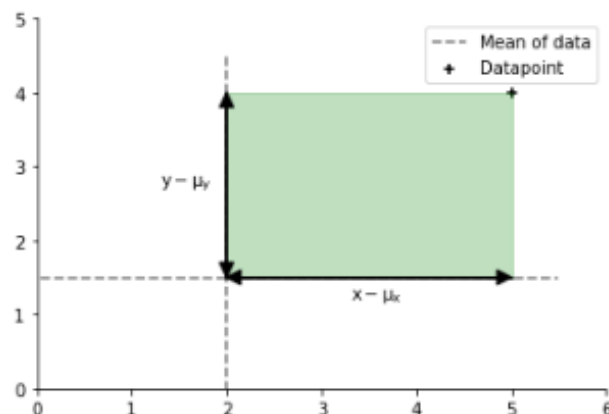
1. In this question we will look at a two-dimensional dataset  $\mathcal{D} = \{\mathbf{x}_i\}_{i=1}^N$  with  $N$  samples. Each sample  $\mathbf{x}_i$  in the dataset is a two-dimensional vector with coordinates  $x, y$ , i.e., the first component of the vector is denoted by  $x$  and the other one by  $y$ .

1 / 1 point

The covariance between two *scalar* random variables is

$$\text{cov}[x, y] = E[(x - \mu_x)(y - \mu_y)] \approx \frac{1}{N} \sum_{i=1}^N (x - \mu_x)(y - \mu_y).$$

In the formula for covariance, we can think of each individual multiplication as the calculation of an area, a rectangle with sides  $x - \mu_x$  and  $y - \mu_y$ .



For this datapoint, an increase in  $x$  from the mean is linked to an increase in  $y$ . Where  $x - \mu_x$  and  $y - \mu_y$  have the same sign, the contribution to the covariance is positive and in green, while if the signs are opposite it will be negative and in red. In other words, green means that  $x$  and  $y$  are positively correlated, while red means they're negatively correlated.

The total sum of areas, divided by the number of points  $n$ , will be the value of the covariance.

Run the code once to see this, then uncomment the line that will show the rectangles and run again.

```
1 # RUN THE CODE ONCE, THEN UNCOMMENT LINE 29 TO VISUALISE COVARIANCE
2
3 fig, ax = plt.subplots()
4
5
```

```

6 #Choose an array by deleting the # in front of the word "data" below.
7 #To switch, put the # back and delete another one
8
9 #Random:
10 data = np.array([[1,2],[5,4],[-2,-3],[4,-2],[2,3],[8,-9]])
11
12 #Straight line:
13 #data = np.array([[1,1],[-3,-3],[2,2],[7,7]])
14
15 #Q1: square
16 #data = np.array([[0,0],[4,4],[0,4],[4,0]])
17
18 #Feel free to input your own array or modify the ones above!
19
20 # First calculate the mean with NumPy function np.mean().
21 # The first argument is the dataset and "axis" specifies the direction
22 # Variance in 1D can be calculated similarly with np.var()
23 mean_data = np.mean(data, axis=0)
24 create_plot(data) #which also adds 1d variances
25
26 area=0
27 mean = mean_data
28
29 for i in range(len(data)):
30     # show_rectangle(mean, data[i])
31     # and a calculation that adds (or subtracts)
32     # the value of the area to our value of the covariance:
33     area += calculate_area(mean, data[i])
34
35 plt.show()

```

Run

Reset

The dashed lines meet at the mean of the dataset. The blue lines represent the magnitude of the variance of the x (horizontal) and y (vertical) components of the dataset.

If red and green balance out, the covariance will be 0. Otherwise the sign of the covariance will give a direction in which the points appear to correlate.

What is  $\text{cov}(x, y)$  for the dataset in the array labelled "Q1: square"? Is it what you would expect from the plot?

0

✓ Correct

Correct! Since the points are evenly distributed around the mean they balance out and there is no way to determine a direction of correlation.

2. The covariance matrix is given by

$$\begin{bmatrix} \text{cov}(x, x) & \text{cov}(x, y) \\ \text{cov}(y, x) & \text{cov}(y, y) \end{bmatrix} = \begin{bmatrix} \text{var}(x) & \text{cov}(x, y) \\ \text{cov}(y, x) & \text{var}(y) \end{bmatrix}$$

Compute the covariance matrix for the following dataset

$$\mathcal{D} = \left\{ \begin{bmatrix} 1 \\ 2 \end{bmatrix}, \begin{bmatrix} 5 \\ 4 \end{bmatrix} \right\}$$

Here, every column vector represents a data point.

Do the exercise using pen and paper. You can check if your answer makes sense with this codeblock.

```

1  data = np.array([[1,2],[5,4]])
2
3  mean_data = np.mean(data, axis=0)
4  create_plot(data)
5
6  area=0
7  mean = mean_data
8
9  for i in range(len(data)):
10     show_rectangle(mean, data[i])
11     area += calculate_area(mean, data[i])
12
13  plt.show()
```

[Run](#)
[Reset](#)

☐  $\begin{bmatrix} 1 & 2 \\ 2 & 4 \end{bmatrix}$

☒  $\begin{bmatrix} 4 & 2 \\ 2 & 1 \end{bmatrix}$

☐  $\begin{bmatrix} 2 & 2 \\ 4 & 1 \end{bmatrix}$

✓ **Correct**  
Good job!

3. Consider a data set  $\mathcal{D}$  with covariance matrix  $\begin{bmatrix} 3 & 2 \\ 2 & 4 \end{bmatrix}$ .

What is the covariance matrix if we multiply every vector in  $\mathcal{D}$  by 2?

Run the codeblock below to observe what happens to the example in Question 2 (this will NOT give you the answer to this question but might aid intuition).

```

1 data = np.array([[1,2],[5,4]])
2 #data *= 2
3 #Uncomment the line above to multiply by 2 and run again
4
5 mean_data = np.mean(data, axis=0)
6 create_plot(data)
7
8 area=0
9 mean = mean_data
10
11 for i in range(len(data)):
12     show_rectangle(mean, data[i])
13     area += calculate_area(mean, data[i])
14
15 plt.show()
```

Run

Reset

- ☒  $\begin{bmatrix} 12 & 8 \\ 8 & 16 \end{bmatrix}$
- ☐  $\begin{bmatrix} 16 & 8 \\ 8 & 12 \end{bmatrix}$
- ☐  $\begin{bmatrix} 4 & 2 \\ 2 & 3 \end{bmatrix}$
- ☐  $\begin{bmatrix} 3 & 2 \\ 2 & 4 \end{bmatrix}$

✓ Correct

Yes, every element in the covariance matrix is multiplied by 4.

4. Consider the data set  $\mathcal{D} = \left\{ \begin{bmatrix} 1 \\ 2 \end{bmatrix}, \begin{bmatrix} 7 \\ 4 \end{bmatrix} \right\}$  with covariance matrix  $\begin{bmatrix} 9 & 3 \\ 3 & 1 \end{bmatrix}$ .

Compute the new covariance matrix when we add  $\begin{bmatrix} 2 \\ 2 \end{bmatrix}$  to each element in  $\mathcal{D}$ .

Run the codeblock below to observe what happens to the example in Question 2 when a vector is added to every point (this will NOT give you the answer but might aid intuition).

```
1 data = np.array([[1,2],[5,4]])
2 #data += [2,2]
3 #Uncomment line above after first run to add [2,2], then run again
4
5 mean_data = np.mean(data, axis=0)
6 create_plot(data)
7
8 area=0
9 mean = mean_data
10
11 for i in range(len(data)):
12     show_rectangle(mean, data[i])
13     area += calculate_area(mean, data[i])
14
15 plt.show()
```

Run

Reset

- ☒  $\begin{bmatrix} 9 & 3 \\ 3 & 1 \end{bmatrix}$
- ☐  $\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$
- ☐  $\begin{bmatrix} 11 & 5 \\ 5 & 3 \end{bmatrix}$

✓ Correct

Well done. The covariance will not change.

5. We are looking at a data set  $\mathcal{D}$  where every element in  $\mathcal{D}$  consists of an  $x$  and  $y$  coordinate. The data covariance matrix is given by

$$\begin{bmatrix} 1 & 0.8 \\ 0.8 & 1 \end{bmatrix}$$

Which of the following statements is correct?

- ☒  $x$  and  $y$  are positively correlated, i.e., when  $x$  increases then  $y$  increases on average, and vice versa.
- ☐  $x$  and  $y$  are negatively correlated, i.e., when  $x$  increases then  $y$  decreases on average, and vice versa.
- ☐  $x$  and  $y$  are uncorrelated, i.e., when  $x$  increases then  $y$  does not change on average (and vice versa).

✓ **Correct**

Well done!