

МІНІСТЕРСТВО ОСВІТИ І НАУКИ УКРАЇНИ
КИЇВСЬКИЙ НАЦІОНАЛЬНИЙ УНІВЕРСИТЕТ імені Тараса Шевченка
ФАКУЛЬТЕТ ІНФОРМАЦІЙНИХ ТЕХНОЛОГІЙ
Кафедра програмних систем і технологій

Дисципліна
«Ймовірнісні основи програмної інженерії»

Лабораторна робота № 3
«Двовимірна статистика»

Виконав:	АНТОНОВ Олександр Лаврентійович	Перевірила:	Марцафей Анна Сергіївна
Група	ІПЗ-24(2)	Дата перевірки	
Форма навчання	денна	Оцінка	
Спеціальність	121		

Мета: Навчитись використовувати на практиці набуті знання про міри в двовимірній статистиці.

Завдання 1

1. Намалюйте діаграму розсіювання для даних. Укажіть, чи існує тренд у даних. Якщо так, то вкажіть, чи є це негативним трендом, чи позитивним.
2. Знайдіть коваріацію.
3. Знайти рівняння лінії регресії y від x .
4. Розрахуйте коефіцієнт кореляції між даними.

Завдання 1

1. Намалуйте діаграму розсіювання для даних. Укажіть, чи існує тренд у даних. Якщо так, то вкажіть, чи є це негативним трендом, чи позитивним.

Діаграма розсіювання (точкова діаграма) — один з типів математичних діаграм, що використовує декартову систему координат для відображення значень двох змінних для набору даних. Дані показані у вигляді набору точок, кожен з яких має значення однієї змінної, тобто визначає її положення на горизонтальній осі та значення іншої змінної — її положення на вертикальній осі.

Тренд - тенденція зміни показників часового ряду. Тренди можуть бути описані різними функціями

Виведення даних з файлу, діаграма розсіювання

```
7     averagex, averagey = 0.0, 0.0
8     def connect_txt(nameoffile):
9         inputdata = []
10        input = open("input_103.txt")
11        input.seek(1)
12        for line in input:
13            inputdata.append(input.read(3))
14            input.read(1)
15            inputdata.append(input.read(2))
16        for i in range(int(len(inputdata))):
17            inputdata[i] = inputdata[i].replace(',', ' ')
18        data = [[0 for i in range(2)] for j in range(int(len(inputdata) / 2))]
19        index0 = 0
20        index1 = 0
21        for i in range(int(len(inputdata))):
22            if i % 2 == 0:
23                data[index0][0] = float(inputdata[i])
24                index0 += 1
25            elif i % 2 != 0:
26                data[index1][1] = int(inputdata[i])
27                index1 += 1
28        return data
29
30    def dataX(data):
31        inputdatadata = []
32        for i in range(len(data)):
33            inputdatadata.append(data[i][0])
34        return inputdatadata
35
```

визначення тренду

```
def trend(data):  
    if max(data) == data[len(data)-1]:  
        print("Trend of data is positive")  
        f.write("\nTrend of data is positive")  
    elif min(data) == data[len(data)-1]:  
        print("Trend of data is negative")  
        f.write("\nTrend of data is negative")  
    else:  
        print("The data does not have any trend")  
        f.write("\nThe data does not have any trend")
```

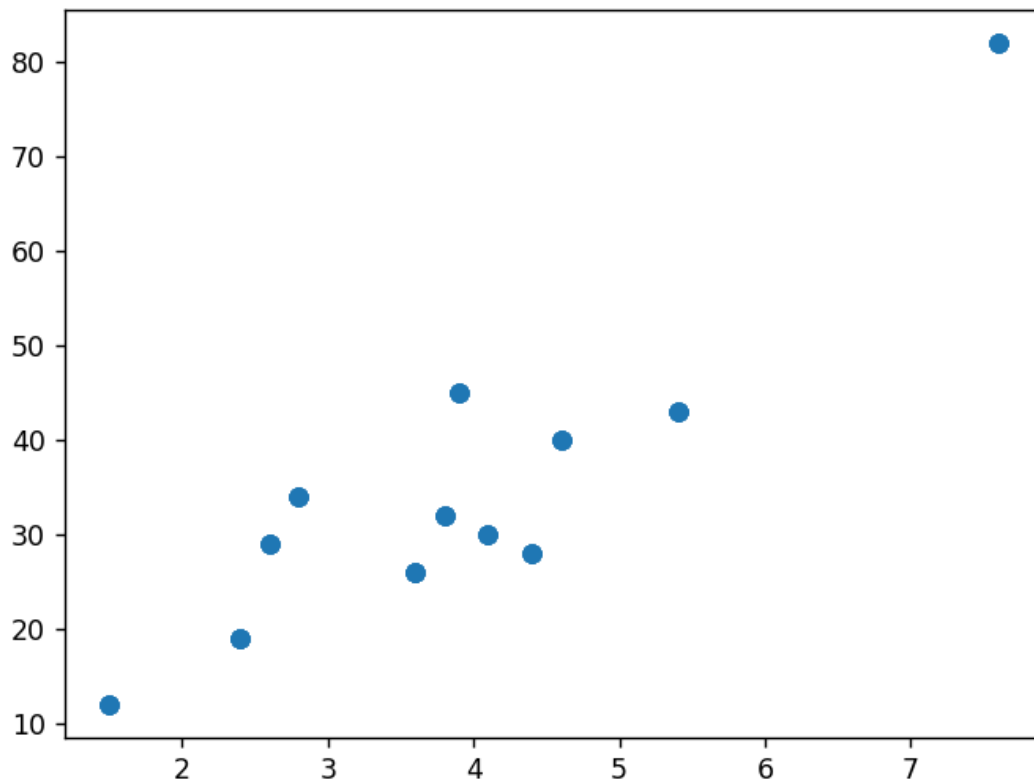
Код створення діаграми розсіювання

```
plt.scatter(infoX, infoY)  
plt.show()
```

Результат виконання

```
Input the name of input file: input_10.txt  
Sorted data: [[1.5, 12], [2.4, 19], [2.6, 29], [2.8, 34], [3.6, 26], [3.8, 32], [3.9, 45], [4.1, 30], [4.4, 28], [4.6, 40], [5.4, 43], [7.6, 82]]  
  
Trend of data is positive
```

Figure 1



2. Знайдіть коваріацію.

Коваріація — це міра спільної мінливості двох випадкових змінних. Якщо більші значення однієї змінної здебільшого відповідають більшим значенням іншої, й те саме виконується для менших значень, тобто змінні схильні демонструвати подібну поведінку, то коваріація є додатною. В протилежному випадку, коли більші значення однієї змінної здебільшого відповідають меншим значенням іншої, тобто змінні схильні демонструвати протилежну поведінку, коваріація є від'ємною. Отже, знак коваріації показує тенденцію в лінійному взаємозв'язку між цими змінними.

Коваріацію (рис. 8) було знайдено за допомогою наступної формули (рис. 6) та реалізовано відповідним кодом (рис. 7).

$$\text{Cov}(X, Y) = \frac{\sum (X_i - \bar{X})(Y_j - \bar{Y})}{n}$$

Формула знаходження коваріації

Знаходження коваріації

```
def covariance(x, y):
    global averagex, averagey
    covariance = 0.0
    for i in range(len(x)):
        averagex += x[i]
        averagey += y[i]
    averagex = averagex / len(x)
    averagey = averagey / len(y)
    for i in range(len(x)):
        covariance += (x[i] - averagex) * (y[i] - averagey)
    covariance = covariance / (len(x)-1)
    print("Коваріацію: ", covariance)
    sus.write("Коваріацію: " + str(covariance))
```

Результат виконання

```
Коваріацію: 22.820202020202036
```

3. Знайти рівняння лінії регресії у від х.

Рівняння, що відображує зміну середньої величини однієї ознаки (у) в залежності від другої (х), називається рівнянням регресії або рівнянням кореляційного зв'язку.

Рівняння лінії регресії (рис. 12) було знайдено за допомогою наступних формул (рис. 9) (рис. 10) та реалізовано відповідним кодом (рис. 11).

$$b_{yx} = \frac{n \sum xy - \sum x \sum y}{n \sum x^2 - (\sum x)^2}$$

Формула знаходження b для рівняння лінії регресії

$$y - \bar{y} = b(x - \bar{x})$$

Формула знаходження лінії регресії

Знаходження рівняння лінійної регресії

```

83 def lineofregression(X, Y):
84     global averagex, averagey
85     byx, sumx, sumy, sumxy, sumx2, sumy2 = 0.0, 0.0, 0.0, 0.0, 0.0, 0.0
86     for i in range(len(X)):
87         sumx += X[i]
88         sumy += Y[i]
89         sumxy += X[i] * Y[i]
90         sumx2 += X[i] * X[i]
91         sumy2 += Y[i] * Y[i]
92     byx = (len(X) * sumxy - (sumx * sumy)) / (len(X) * sumx2 - sumx2)
93     x, y = sp.symbols("x,y")
94     line = sp.Eq(y-averagey, byx*(x-averagex))
95     linex = sp.solve(line, y)
96     liney = sp.solve(line, x)
97     strlinex = str(linex)
98     strliney = str(liney)
99     strlinex = strlinex.replace("[", "")
100    strlinex = strlinex.replace("]", "")
101    strliney = strliney.replace("[", "")
102    strliney = strliney.replace("]", "")
103    print("Лінія регресії y від x. ")
104    print("x = " + strliney)
105    print("y = " + strlinex, "\t(y від x)")
106    sus.write("Лінія регресії y від x.\n")
107    sus.write("x = " + strliney)
108    sus.write("\ny = " + strlinex + "\t(y від x )")
109

```

Лінія регресії у від х.

$x = 0.750738314447591 * y - 22.0444718484419$

$y = 1.33202206515305 * x + 29.3637229167698$ (у від х)

4. Розрахуйте коефіцієнт кореляції між даними.

Коефіцієнт кореляції Пірсона — в статистиці, показник кореляції (лінійної залежності) між двома змінними X та Y, який набуває значень від -1 до +1 включно. Він широко використовується в науці для вимірювання ступеня лінійної залежності між двома змінними.

наскільки залежні дві величини, що більше значення, то залежні значення. Визначає чи залежить траплення велечини

$$r = \frac{\sum (x_i - \bar{x}) (y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2}}$$

Рисунок 13 | Формула розрахунку коефіцієнту кореляції між даними

```
70 def correlation(x, y): |
71     global averagex, averagey
72     corcoef, sum1, sum2, sum3 = 0.0, 0.0, 0.0, 0.0
73     for i in range(len(x)):
74         sum1 += (x[i] - averagex) * (y[i] - averagey)
75         sum2 += (x[i] - averagex) * (x[i] - averagex)
76         sum3 += (y[i] - averagey) * (y[i] - averagey)
77     sum2 = sum2 * sum3
78     corcoef = sum1/math.sqrt(sum2)
79     print("Коефіцієнт кореляції:", corcoef)
80     sus.write("\nКоефіцієнт кореляції:" + str(corcoef))
81
```

Функція розрахунку коефіцієнту кореляції між даними

Коефіцієнт кореляції: 0.901950337071579

Приклад виконання, input_10.txt

Висновок: Під час виконання третьої лабораторної роботи було повторено операції з вхідними даними записаних у txt файл, реалізовано потрібні формули для знаходження інформації за завданнями, також побудовано діаграму розсіювання.

Також було проведено дослідження, де було порівняно коваріацію з кореляцією.

Коваріація та кореляція в основному оцінюють зв'язок між змінними. Найближчою аналогією зв'язку між ними є зв'язок між дисперсією та стандартним відхиленням.

Коваріація вимірює загальну варіацію двох випадкових змінних від їхніх очікуваних значень. Використовуючи коваріацію, ми можемо лише оцінити напрямок зв'язку (чи змінні мають тенденцію рухатися в тандемі чи показують зворотний зв'язок). Однак це не вказує ні на силу зв'язку, ні на залежність між змінними.

З іншого боку, кореляція вимірює силу зв'язку між змінними. Кореляція – це масштабна міра коваріації. Він безрозмірний. Іншими словами, коефіцієнт кореляції завжди є чистим значенням і не вимірюється в жодних одиницях.

Зв'язок між двома поняттями можна виразити за допомогою формули нижче:

$$\rho(X, Y) = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y}$$

Рисунок 16 | Формула зображення зв'язку між коваріацією та кореляцією

Де:

- $\rho(X, Y)$ – кореляція між змінними X і Y
- $\text{Cov}(X, Y)$ – коваріація між змінними X і Y
- σ_X – стандартне відхилення X-змінної
- σ_Y – стандартне відхилення Y-змінної