

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

Виконав:	Ковалець Олег Сергійович	Перевірів:	Марцафей А.С.
Група	ІПЗ-22	Дата перевірки	
Форма навчання	денна	Оцінка	
Спеціальність	121		
2022			

ЛАБОРАТОРНА РОБОТА № 3

ДВОВИМІРНА СТАТИСТИКА

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

Завдання

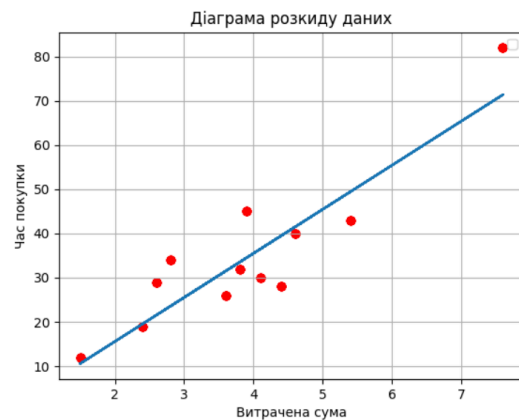
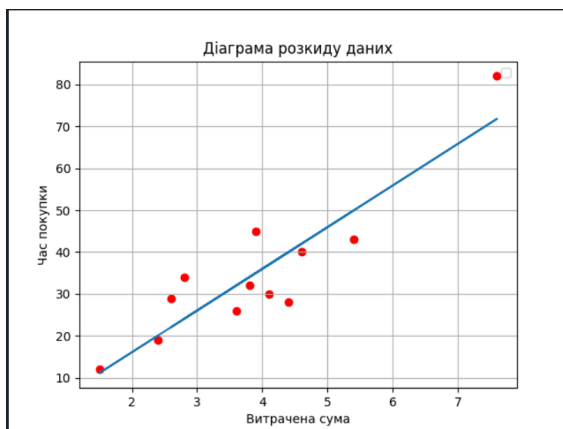
1. Намалювати діаграму розсіювання для даних.

Зчитаємо поданий файл та розіб'ємо його на два масиви, за допомогою яких побудуємо діаграму розсіювання даних. Для побудови використаємо метод scatter бібліотеки matplotlib (python).

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 import math
4
5
6 def filelist(): # Creating two arrays
7     file = input("Enter a name of the file: ")
8     f = open(file, "r")
9     length = int(f.readline())
10    sum = np.zeros(length)
11    time = np.zeros(length)
12
13    for i in range(length):
14        [a, b] = f.readline().split()
15        a = a.replace(',', '.')
16        sum[i] = float(a)
17        time[i] = int(b)
18
19    return sum, time
```

```
47 # Create plots
48 plt.grid(True)
49 plt.scatter(x, y, c='r')
50 plt.plot(x, b1*x+b0)
51 # Labels
52 plt.title("Діаграма розкиду даних")
53 plt.xlabel("Витрачена сума")
54 plt.ylabel("Час покупки")
55 plt.legend()
56 plt.show()
```

Маємо наступні діаграми для файлів **input_10.txt** та **input_100.txt** відповідно:



2. Знайдіть центр ваги і коваріацію.

Центром ваги вважається набір значень $(x_, y_)$, де $x_$ та $y_$ - середні значення для осі x та y відповідно.

Коваріація розраховується за наступною формулою:

$$\text{cov}(X, Y) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) .$$

Отже, знайдемо центр ваги та коваріацію.

```
23 # AverageX, AverageY
24 av_x, av_y = 0, 0
25 for i in range(len(x)):
26     av_x += x[i]
27     av_y += y[i]
28
29 av_x, av_y = round(av_x / len(x), 1), av_y / len(y)
30
31 covariance, var_x, var_y = 0, 0, 0
32 for i in range(len(x)):
33     covariance += (x[i] - av_x)*(y[i] - av_y)
34     var_x += pow(x[i] - av_x, 2)
35     var_y += pow(y[i] - av_y, 2)
36 # Covariance, VariationX, VariationY
37 covariance, var_x, var_y = covariance / len(x), var_x / len(x), var_y / len(y)
```

Маємо наступні значення для **input_10.txt** та **input_100.txt** відповідно:

```
Center of gravity = (3.9, 35.0)
Covariance = 22.999999999999996
```

```
Center of gravity = (3.9, 34.5)
Covariance = 22.592
```

3. Знайти рівняння лінії регресії y від x .

Рівняння лінії регресії $(y = k + mx)$ знаходиться за наступною формулою:

$$m \equiv b_1 = \frac{\text{cov}(X, Y)}{\text{Var}(X)} \text{ and } k \equiv b_0 = \bar{y} - b_1 \bar{x}.$$

Коваріацію вже знайдено, а з нею і дисперсію для x та y в одному циклі (скріншот вище). Отже розрахуємо рівняння лінії регресії y від x .

```

39      b1 = covariance / var_x
40      b0 = av_y - b1*av_x

50      plt.plot(x, b1*x+b0)

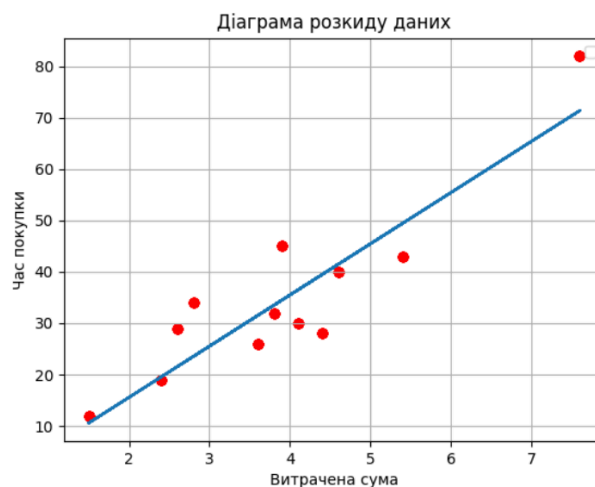
```

Маємо наступні значення для **input_10.txt** та **input_100.txt** відповідно:

```
Regression line: y = 9.953119365308332 + (-3.817165524702496)*x
```

```
Regression line: y = 9.973512272647008 + (-4.3966978633233325)*x
```

Вставивши розраховані коефіцієнти у рівняння для `matplotlib.pyplot`, ми отримали саме ту пряму на діаграмі розсіювання даних у виводі з першого завдання.



4. Розрахуйте коефіцієнт кореляції між даними.

Знайдемо коефіцієнт кореляції за наступною формулою:

$$r = \frac{1}{n-1} \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{s_x} \right) \left(\frac{y_i - \bar{y}}{s_y} \right) ,$$

де $s(x)$ та $s(y)$ - стандартне відхилення для x та y відповідно.

```
42 deviationX, deviationY, correlation = np.sqrt(var_x), np.sqrt(var_y), 0
43 for i in range(len(x)):
44     correlation += ((x[i] - av_x) / deviationX) * ((y[i] - av_y) / deviationY)
45
46 correlation /= (len(x) - 1)
```

Маємо наступні значення для **input_10.txt** та **input_100.txt** відповідно:

```
Correlation coefficient = 0.9106715347800048
```

```
Correlation coefficient = 0.9828959169823492
```

Бачимо, що коефіцієнт кореляції знаходиться в межах від -1 до 1, як і має бути.

5. Зробити висновок про залежності.

Можемо розрахувати висновки за допомогою коду нижче:

```
63 writer.write("\tConclusions\n")
64
65 if correlation == 0:
66     writer.write("x and y are not linearly related\n")
67 elif math.fabs(correlation) > np.sqrt(3) / 2:
68     writer.write("x and y have strong linear relation\n")
69 else:
70     writer.write("x and y have weak linear relation\n")
71
72 if correlation > 0:
73     writer.write("The relation between x and y is positive")
74 elif correlation < 0:
75     writer.write("The relation between x and y is negative")
```

Маємо однакові висновки для двох файлів:

```
Conclusions
x and y have strong linear relation
The relation between x and y is positive
```

Підсумувавши, маємо наступний код для запису результатів у файл output.txt:

```
57 # Write results in file
58 writer = open("output.txt", "w")
59 writer.write(f"Center of gravity = ({av_x}, {av_y})\n")
60 writer.write(f"Covariance = {covariance}\n")
61 writer.write(f"Regression line: y = {b1} + ({b0})*x\n")
62 writer.write(f"Correlation coefficient = {correlation}\n")
63 writer.write("\tConclusions\n")
```

Та наступний вивід:

```
main.py × output.txt × input_10.txt × input_100.txt ×
1 Center of gravity = (3.9, 35.0)
2 Covariance = 22.999999999999996
3 Regression line: y = 9.953119365308332 + (-3.817165524702496)*x
4 Correlation coefficient = 0.9828959169823492
5 Conclusions
6 x and y have strong linear relation
7 The relation between x and y is positive
```

```
main.py × output.txt × input_10.txt × input_100.txt ×
1 Center of gravity = (3.9, 34.5)
2 Covariance = 22.592
3 Regression line: y = 9.973512272647008 + (-4.3966978633233325)*x
4 Correlation coefficient = 0.9106715347800048
5 Conclusions
6 x and y have strong linear relation
7 The relation between x and y is positive
```

Висновок:

Протягом даної лабораторної роботи було ознайомлено з двовимірною статистикою. Було побудовано діаграму розсіювання даних, лінію регресії та розраховано центр ваги для неї за вказаними даними. Також на основі цих даних було знайдено коваріантність та коефіцієнт кореляції. Після виконання основних було підведено підсумки. Результати записані у вихідний файл і всі вимоги виконання були дотримані.