

Міністерство освіти і науки України
Національний технічний університет України «Київський політехнічний інститут
імені Ігоря Сікорського»
Факультет інформатики та обчислювальної техніки

Кафедра інформатики та програмної інженерії

Звіт

з комп'ютерного практикуму № 3 з дисципліни
«Аналіз даних в інформаційних системах»
на тему: «Описова статистика»

Виконав студент ІП-13, Дойчев Костянтин Миколайович
(шифр, прізвище, ім'я, по батькові)

Перевірів Олійник Юрій Олександрович
(прізвище, ім'я, по батькові)

Комп'ютерний практикум 3

Тема – Описова статистика.

Мета – ознайомитись з методикою первинної обробки статистичних даних; проаналізувати вплив способу представлення даних на їх інформативність.

Завдання

Основне:

1. Скачати дані із файлу [Data2.csv](#)
2. Записати дані у data frame
3. Дослідити структуру даних
4. виправити помилки в даних
5. Побудувати діаграми розмаху та гістограми
6. Додати стовпчик із щільністю населення

кількість рядків та колонок, назви всіх колонок, кількість записів в кожній з них, тип даних колонки та використання пам'яті.

Висновки:

1. Помилки в правописі (spelling) колонки Population. Populatiion -> Population
2. Пусті дані (NaN)
3. Від'ємні значення де вони не мають бути. Площа не може бути від'ємною

Виправлення помилок

Змінимо назву колонки та перевіримо наявність пустих значень

Based on what we see the dataset is not clean, which can lead to problems in the future.

So, we need to clean it up.

1. Fix the naming (spelling mistakes, etc.)
2. Check for undefined values
3. Replace those values with mean
4. Verify the data types
5. Convert invalid values (area and GDP per capita cannot be negative)

```
36 1 # fix the naming
2 df = df.rename(columns={'Country Name': 'Country', 'Populatiion': 'Population'})
Executed at 2023.06.25 19:26:42 in 35ms
```

```
37 1 # Check for undefined values
2 df.isna().any()
Executed at 2023.06.25 19:26:42 in 34ms
```

```
37 1 |< < 6 rows > >| Length: 6, dtype: bool pd.Series
2
3 Country      False
4 Region       False
5 GDP per capita True
6 Population    True
7 CO2 emission  True
8 Area         False
```

Замінімо європейський стиль написання десяткових чисел (10,5) на американський (10.5)

```
In 39 1 # replace commas with dots (to convert to decimal)
2 # convert to string
3 df['Population'] = df['Population'].astype(str)
4 # replace commas with dots
5 df['Population'] = df['Population'].str.replace(',', '.')
Executed at 2023.06.25 19:26:42 in 24ms
```

```
In 40 1 # check the types
2 df.dtypes
Executed at 2023.06.25 19:26:42 in 20ms
```

```
Out 40 1 |< < 6 rows > >| Length: 6, dtype: object pd.Series
2
3 Country      object
4 Region       object
5 GDP per capita float64
6 Population    object
7 CO2 emission  float64
8 Area         float64
```

```
In 48 1 # convert to float
2 df['Population'] = df['Population'].astype(float)
3 df.dtypes
Executed at 2023.06.25 19:27:10 in 16ms
```

Замінімо пусті значення середнім арифметичним та приведемо числа до додатніх

```
1 # replace empty values with mean
2 df = df.fillna(df.mean(numeric_only=True))
   Executed at 2023.06.25 19:26:42 in 15ms

1 # check whether we still have undefined values
2 df.isna().any()
   Executed at 2023.06.25 19:26:42 in 10ms
```

▼ |< < 6 rows > >| Length: 6, dtype: bool **pd.Series** * CSV ▾ ⬇ ⌘ 🔍 ⋮

	<unnamed>
Country	False
Region	False
GDP per capita	False
Population	False
CO2 emission	False
Area	False

```
1 # replace floats with absolute values as it's not possible to have negative values in these columns
2 df['GDP per capita'] = df['GDP per capita'].abs()
3 df['Area'] = df['Area'].abs()
4
   Executed at 2023.06.25 19:26:42 in 6ms
```

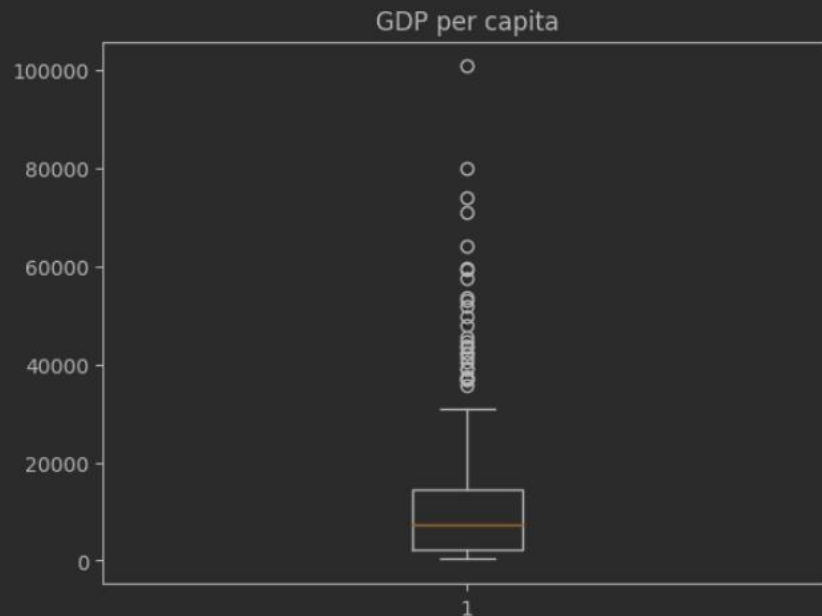
Візуалізація

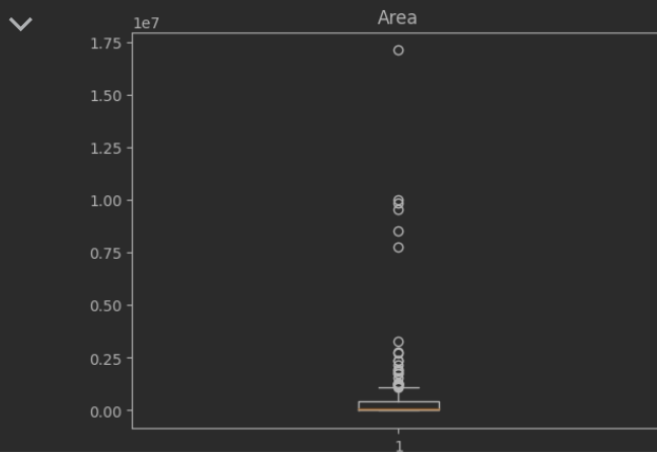
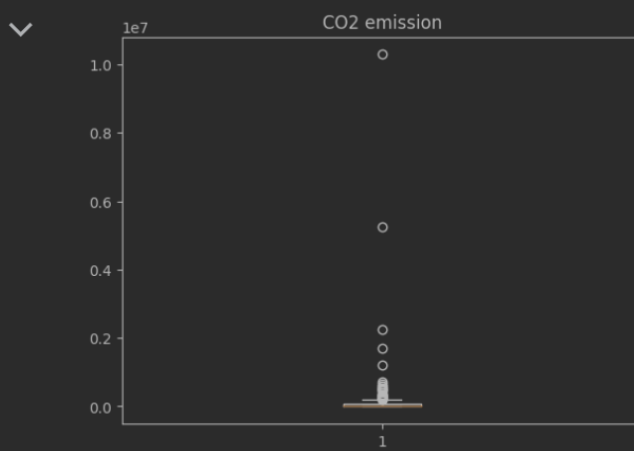
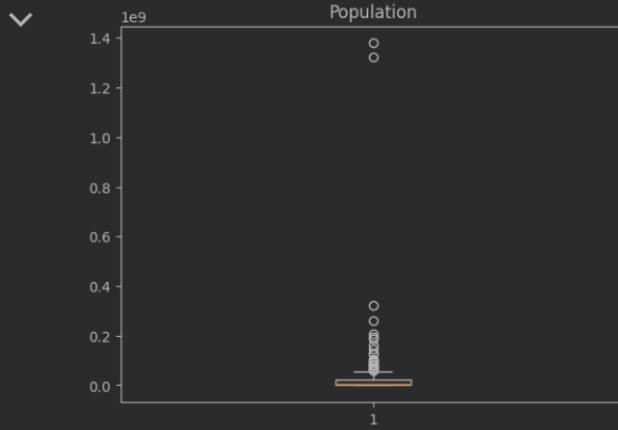
Виведемо діаграми розмаху та гістограми для кожного стовпця з чисельними даними.

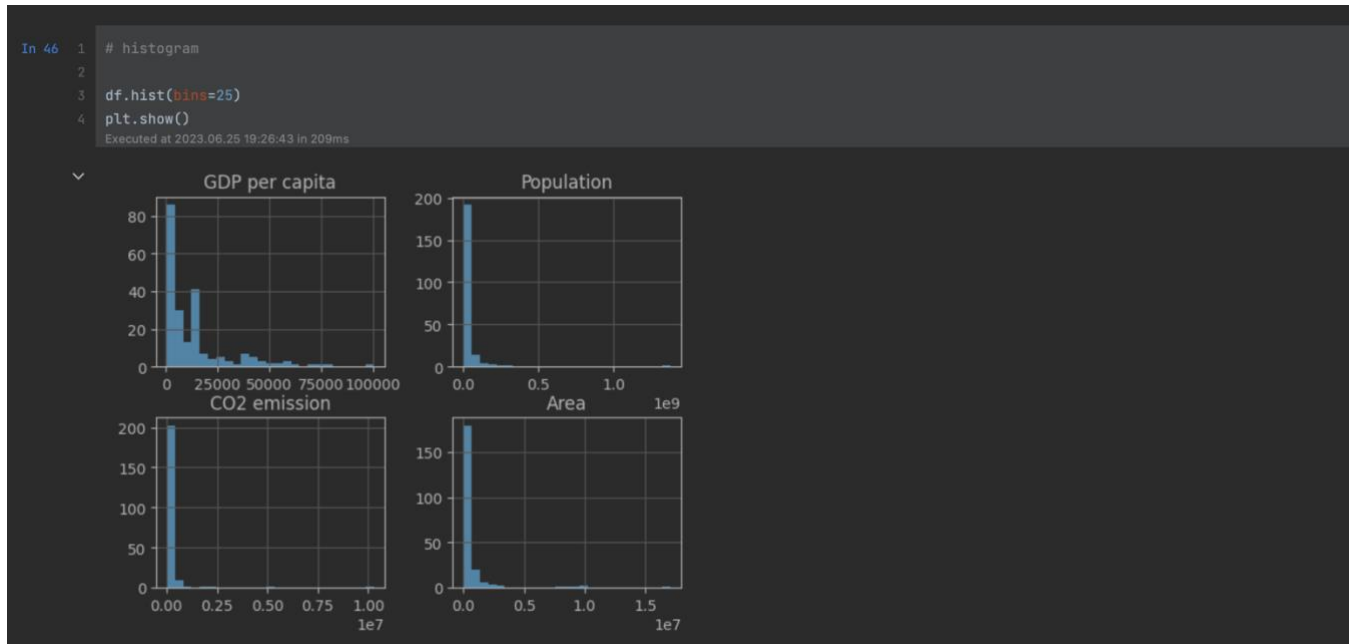
Data visualization

```
In 52 1 # boxplot
      2 for column in df.columns:
      3     if df[column].dtype == float:
      4         plt.figure()
      5         plt.title(column)
      6         plt.boxplot(df[column])
```

Executed at 2023.06.25 19:38:01 in 317ms







Додавання стовпчику із щільністю населення

Додаємо стовпчик із щільністю населення кожної країни, який є просто представленням кількості населення поділеного на площу країни.

```
In 50 1 # Create a population density column
2
3 df['Population Density'] = df['Population'] / df['Area']
Executed at 2023.06.25 19:27:46 in 3ms
```

```
In 51 1 df
Executed at 2023.06.25 19:27:50 in 18ms
```

Out 51 217 rows x 7 columns [pd.DataFrame](#)

	Country	Region	GDP per capita	Population	CO2 emission	Area	Population Density
0	Afghanistan	South Asia	561.778746	34656032.0	9809.225000	652860.0	53.083405
1	Albania	Europe & Central Asia	4124.982390	2876101.0	5716.853000	28750.0	100.038296
2	Algeria	Middle East & North Africa	3916.881571	40606052.0	145400.217000	2381740.0	17.048902
3	American Samoa	East Asia & Pacific	11834.745230	55599.0	165114.116337	200.0	277.995000
4	Andorra	Europe & Central Asia	36988.622030	77281.0	462.042000	470.0	164.427660
5	Angola	Sub-Saharan Africa	3308.700233	28813463.0	34763.160000	1246700.0	23.111786
6	Antigua and Barbuda	Latin America & Caribbean	14462.176280	100963.0	531.715000	440.0	229.461364
7	Argentina	Latin America & Caribbean	12440.320980	43847430.0	204024.546000	2780400.0	15.770188
8	Armenia	Europe & Central Asia	3614.688357	2924016.0	5529.836000	29740.0	98.346200
9	Aruba	Latin America & Caribbean	13374.833168	104822.0	872.746000	180.0	582.344444

Висновок

У цьому комп'ютерному практикуму було вивчено можливості Python, а саме Pandas у роботі з даними. Вхідні дані було записано в data frame, дані були не зовсім «чистими», тому було необхідно змінити назви, заповнити пусті значення, виправити помилки. На діаграмах розмаху було помічено великий розмах між даними. Наприклад, на діаграмі населення є дві країни з кількістю населення значно більшою за всі інші, так само і з викидами CO₂, дані з ВВП на душу населення є найбільш кучними.