SPRAWOZDANIE

Zajęcia: Nauka o danych I

Prowadzący: prof. dr hab. Vasyl Martsenyuk

Laboratorium Nr 3 Data 19.10.2024

Temat: "Wykorzystanie pakietu

Pandas do manipulacji i

przetwarzania danych w Pythonie"

Wariant 11

Szymon Nycz Informatyka II stopień, niestacjonarne, 1 semestr, gr.1b

1. Polecenie:

Premise General Population COVID-19 Health Services Disruption Survey 2020 http://ghdx.healthdata.org/record/ihme-data/premise-general-population-covid-19-health-services-disruption-survey-2020

2. Link do repozytorium:

Link: https://github.com/Maciek332/Semestr 1 Nycz/tree/master/NoD/Lab 3

3. Opis programu opracowanego

- Wczytywanie danych i wyświetlanie podstawowych informacji
 - Wczytaj dane z pliku



Wyświetl pierwsze 5 wierszy

```
observation_id
                                     submitted_time gender
  gp_4503617949401088 2020-07-07 14:48:29.83 UTC gp_4503631639609344 2020-07-09 13:22:37.107 UTC
                                                       Male
                                                      Female
  gp_4503700758593536 2020-07-04 18:53:36.471 UTC
  gp_4503737805832192 2020-07-12 17:58:20.798 UTC
                                                       Male
  gp_4503819343101952 2020-07-06 00:20:22.983 UTC
                                              geography \
            Under 16
                                    Suburban/Peri-urban
  26 to 35 years old City center or metropolitan area
  36 to 45 years old City center or metropolitan area
  26 to 35 years old
  26 to 35 years old
                                    Suburban/Peri-urban
                                 financial_situation
                                                                   education
  I can afford food and regular expenses, but no... Secondary/high school
          I cannot afford enough food for my family College or university
                                                             Primary school
  I can comfortably afford food, clothes, and fu...
  I can afford food and regular expenses, and bu...
                                                            Technical school
  I can afford food and regular expenses, and bu...
                                                           Technical school
   employment_status
                                ethnicity
                                               religion ...
  Employed full-time
                                 Ankole Christianity ...
          Unemployed
                                  Mestizo Catholicism ...
  Employed full-time Non-hispanic White
                                           Agnosticism ...
            12.0
                                  Colombia 1.691686 gp_4703669741420544
       400000.0
                                  Colombia 1.691686 gp 4762741153988608
[5 rows x 48 columns]
```

Sprawdź podstawowe informacje o danych

```
D ~
        print(data_frame.info())
    <class 'pandas.core.frame.DataFrame'>
    RangeIndex: 52490 entries, 0 to 52489
    Data columns (total 48 columns):
     # Column
                                          Non-Null Count Dtype
         observation id
                                          52490 non-null object
         submitted time
                                         52490 non-null object
     2 gender
                                         52469 non-null object
                                         52490 non-null object
     3 age
     4 geography
                                         52490 non-null object
     5 financial situation
                                        52490 non-null object
     6 education
                                         52490 non-null object
                                         52490 non-null object
        employment status
     8 ethnicity
                                          52490 non-null object
         religion
                                          52490 non-null object
     10 gp hh
                                         52478 non-null object
                                         52490 non-null object
     11 gp_pre_provider_need
     12 gp_pre_provider_condition
                                         21777 non-null object
     13 gp pre provider condition other 2982 non-null
                                                         object
     14 gp_pre provider visit
                                         21777 non-null object
     15 gp_pre_provider_where
                                         9972 non-null object
     16 gp_pre_provider_where_other
                                         803 non-null
                                                         object
     17 gp_pre_provider_num_visit
                                         19037 non-null object
     18 gp pre provider why
                                          8363 non-null
                                                         object
     19 gp_pre_provider_why_other
                                          398 non-null
                                                         object
     47 user id
                                          52490 non-null object
    dtypes: float64(3), object(45)
    memory usage: 19.2+ MB
    None
```

Wyświetl podstawowe statystyki opisowe

```
print(data frame.describe())
     ✓ 0.0s
[40]
           gp pre income gp post income
                                                 weight
            5.249000e+04
                             5.249000e+04 52490.000000
    count
    mean
            1.905125e+56
                             1.905125e+59
                                               1.795820
    std
            4.364774e+58
                             4.364774e+61
                                               0.384507
    min
            0.000000e+00
                             0.000000e+00
                                               1.000000
    25%
            1.500000e+02
                             7.000000e+01
                                               1.544666
    50%
            3.425000e+03
                             2.000000e+03
                                               1.730488
    75%
            2.500000e+04
                             2.000000e+04
                                               1.908871
    max
            1.0000000e+61
                             1.000000e+64
                                               6.411420
```

- Obliczanie podstawowych statystyk
 - o Oblicz średnią

Oblicz mediane

Oblicz odchylenie standardowe

- Identyfikacja i obsługa brakujących danych
 - Sprawdź brakujące wartości

```
missing_values = data_frame.isnull().sum()
   print("Brakujce wartoci w kadej kolumnie:")
   print(missing_values)
Brakujce wartoci w kadej kolumnie:
observation id
                                        0
submitted time
                                        0
gender
                                       21
                                        0
age
geography
                                        0
financial situation
education
                                        0
employment_status
                                        0
ethnicity
                                        0
religion
                                        0
gp hh
                                       12
gp_pre_provider need
                                        0
gp_pre_provider_condition
                                    30713
gp_pre_provider_condition_other
                                    49508
gp pre provider visit
                                    30713
gp_pre_provider_where
                                    42518
gp_pre_provider_where_other
                                    51687
gp_pre_provider_num_visit
                                    33453
                                    44127
gp_pre_provider_why
gp pre provider why other
                                    52092
```

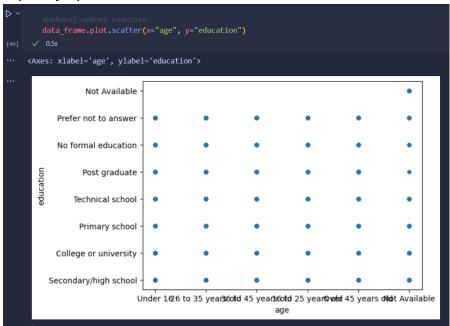
Uzupełnij brakujące wartości średnią

Usuń wiersze w których brakuje danych

- Wykrywanie wartości odstających
 - Oblicz IQR

o Zidentyfikuj wartości odstające

- Analiza zależności między kolumnami
 - Wykonaj wykres rozrzutu



- Przekształcenie danych
 - Dodaj nową kolumnę

```
#Dodaj nową kolumnę
data_frame["nowa"] = data_frame["gp_pre_income"] / 2

[50] 

0.0s
```

Grupuj dane według kolumny

```
#Grupuj dane wedug kolumny 'region' i oblicz średnią
grouped = data_frame.groupby("gender")["gp_pre_income"].mean()
print("Coś:")
print(grouped)

**Coś:
gender
Female
2.366555e+08
Male
2.867301e+56
Prefer not to respond
7.089249e+12
Name: gp_pre_income, dtype: float64
```

Posortuj dane według kolumny

```
df_sorted = data_frame.sort_values(by="age", ascending=False)
   print("Dane posortowane według wieku:")
Dane posortowane według wieku:
           observation_id
                                          submitted_time gender
                                                                        age
      gp_4503617949401088 2020-07-07 14:48:29.83 UTC Male Under 16
39805 gp_6212113459838976 2020-07-01 16:26:44.48 UTC
                                                            Male Under 16
39774 gp_6210954120658944 2020-07-18 20:27:53.831 UTC Female Under 16 7745 gp_4835939500425216 2020-07-03 07:36:54.128 UTC Male Under 16
39769 gp 6210720447594496 2020-07-07 19:59:10.754 UTC Male Under 16
                 geography
                                                            financial_situation \
       Suburban/Peri-urban I can afford food and regular expenses, but no...
0
39805
                     Rural
                                         I can afford food, but nothing else
                     Rural I can afford food and regular expenses, and bu...
                      Rural I can afford food and regular expenses, but no...
7745
39769
                     Rural
                                    I cannot afford enough food for my family
                   education employment status ethnicity
                                                                  religion ...
       Secondary/high school Employed full-time Ankole Christianity ...
0
39805 Secondary/high school
                                         Student Sarakole
                                                                   Muslim ...
            Technical school Employed part-time Igbo Christianity ...
ndary/high school Self-employed Bisaya Other ...
39774
7745
       Secondary/high school
39769
            Primary school
                                   Self-employed Malinke
                                                                    Muslim ...
      gp post labor force gp pre labor force gp unemployment why \
0
                       No
                                           No
                                                               NaN
                        No
                                           No
                                                               NaN
39805
```

4. Wnioski

Rozpoczęliśmy od załadowania danych z pliku CSV i zaprezentowania podstawowych informacji o zbiorze danych. Następnie obliczyliśmy kluczowe statystyki opisowe dla wybranych kolumn, aby lepiej zrozumieć rozkład danych. Brakujące dane mogły wpłynąć na jakość analizy, dlatego musieliśmy je zidentyfikować i odpowiednio się nimi zająć. Wartości odstające mogły zniekształcać wyniki, więc również je wykryliśmy. Analizując zależności między różnymi kolumnami, obliczyliśmy współczynniki korelacji. Na końcu dokonaliśmy transformacji danych, tworząc nowe kolumny, grupując je i sortując.

Podczas realizacji tych kroków opanowaliśmy zarówno podstawowe, jak i bardziej zaawansowane techniki manipulacji danymi w Pandas. Zrozumieliśmy, jak efektywnie ładować, analizować i przekształcać dane, co stanowi kluczową umiejętność w pracy z danymi oraz analizie danych.