

# SPRAWOZDANIE

Zajęcia: Nauka o danych I

Prowadzący: prof. dr hab. Vasyl Martsenyuk

Laboratorium Nr 5 Data 23.11.2024 Temat: Wykorzystanie narzędzi do eksploracyjnej analizy danych (EDA) Wariant drugi (2)	Bartosz Bieniek Informatyka II stopień, stacjonarne, 1 semestr, gr.A
---	---

## 1. Polecenie

- Zidentyfikować wartości odstające za pomocą algorytmu Isolation Forest,
- Zredukować wymiarowość danych z użyciem PCA,
- Stworzyć zaawansowane interaktywne wizualizacje danych,
- Zwizualizować dane wielowymiarowe za pomocą zaawansowanych algorytmów (t-SNE, UMAP),
- Stworzyć interaktywne wizualizacje danych w 2D i 3D,
- Przeanalizować zależności między zmiennymi za pomocą macierzy korelacji.
- Przeprowadzać testy statystyczne dla analizy różnic w grupach.

## 2. Opis programu opracowanego

<https://github.com/mindgoner/Studia/tree/master/Nauka%20o%20Danych/Laboratorium%205>

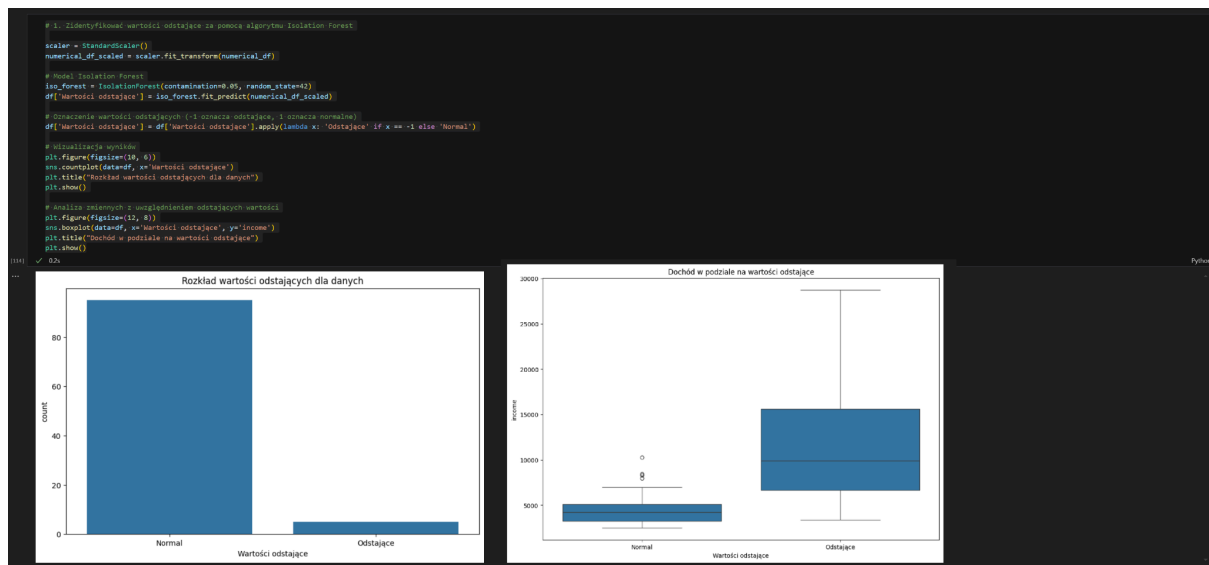
```
# 0. Przygotowanie danych
import pandas as pd
import seaborn as sns
import plotly.express as px
import plotly.graph_objects as go
import matplotlib.pyplot as plt
from sklearn.ensemble import IsolationForest
from sklearn.preprocessing import StandardScaler
from sklearn.decomposition import PCA

df = pd.read_csv('data.csv')
print(df.info())
print(df.describe())
numerical_features = ['age', 'income', 'outcome']
numerical_df = df[numerical_features]

if numerical_df.isnull().sum().any():
    numerical_df.fillna(numerical_df.mean(), inplace=True)
```

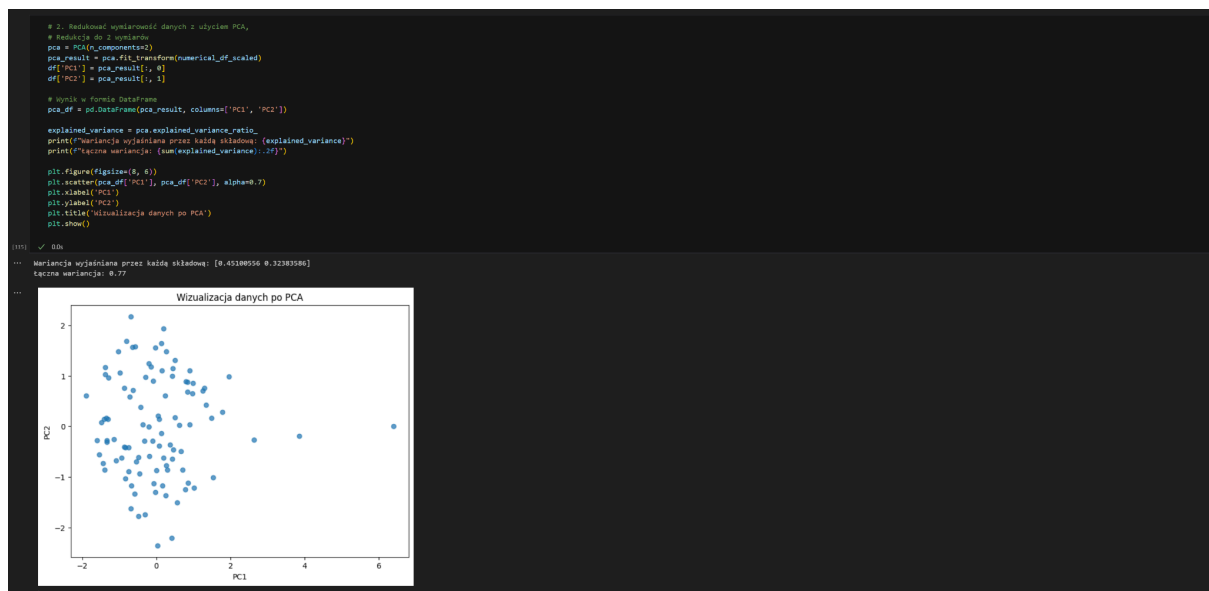
Rys. 1. Kod Źródłowy.

Zaimportowano potrzebne biblioteki do analizy danych, wizualizacji i przetwarzania wstępnego. Wczytano dane z pliku CSV do DataFrame i wybrano kolumny numeryczne: age, income, outcome. Sprawdzono brakujące wartości i w razie potrzeby wypełniono je średnimi dla każdej kolumny.



Rys. 2. Kod Źródłowy i efekt jego wykonania.

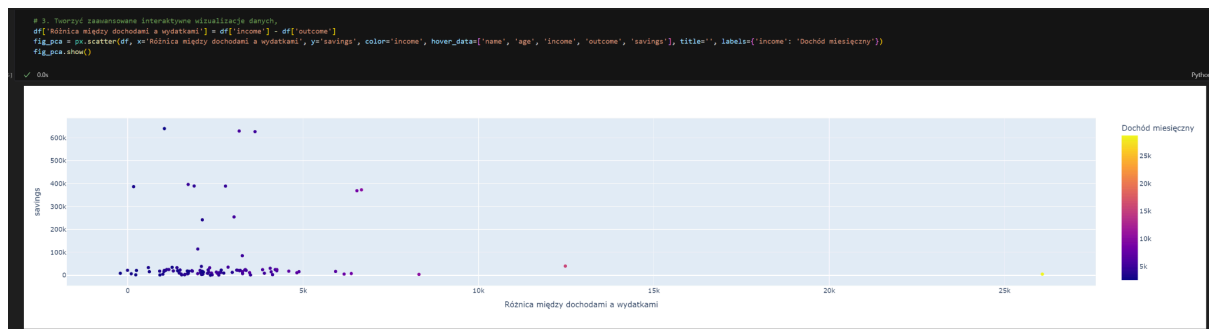
Przeskalowano dane numeryczne, aby algorytm Isolation Forest mógł poprawnie identyfikować wartości odstające. Wartości te oznaczono jako "Odstajace" lub "Normal" i zapisano w nowej kolumnie DataFrame. Wyniki przedstawiono na wykresach: rozkład odstających oraz dochody w podziale na grupy.



Rys. 3. Kod Źródłowy i efekt jego wykonania.

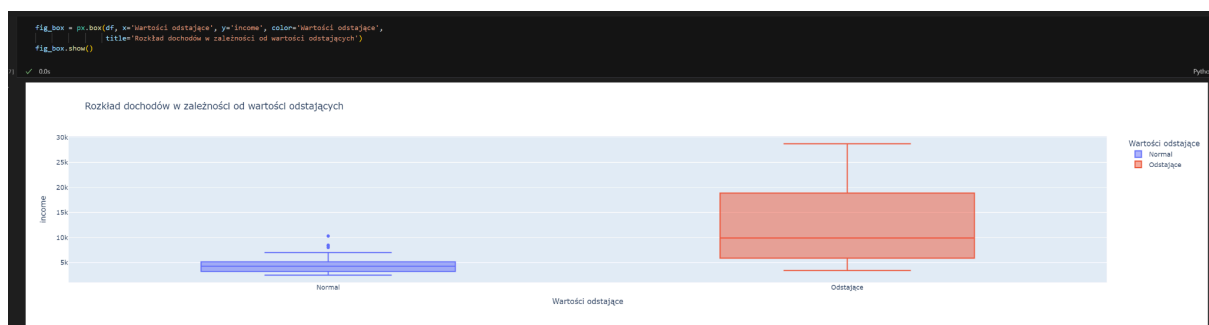
Zredukowano wymiarowość danych numerycznych do dwóch głównych składowych za pomocą PCA, zapisując wyniki jako PC1 i PC2 w DataFrame. Wyświetlono udział wariancji wyjaśnianej przez każdą składową oraz jej łączną

wartość. Wizualizację danych po redukcji wymiarowości wykonano za pomocą wykresu punktowego.



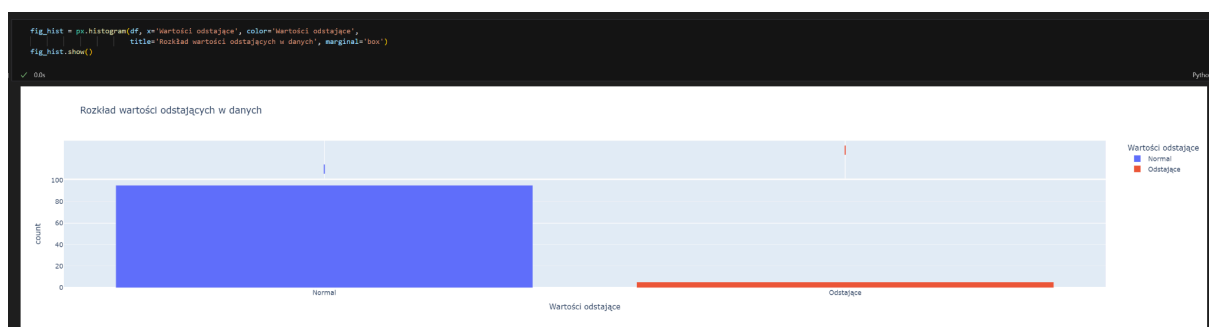
Rys. 4. Kod Źródłowy i efekt jego wykonania.

Obliczono różnicę między dochodami a wydatkami i zapisano ją w nowej kolumnie. Utworzono interaktywną wizualizację punktową, przedstawiającą związek między różnicą dochodów i wydatków a oszczędnościami, z użyciem koloru do reprezentacji poziomu dochodów. Włączono dodatkowe dane do wyświetlania w chmurkach, takie jak imię, wiek czy szczegóły finansowe.



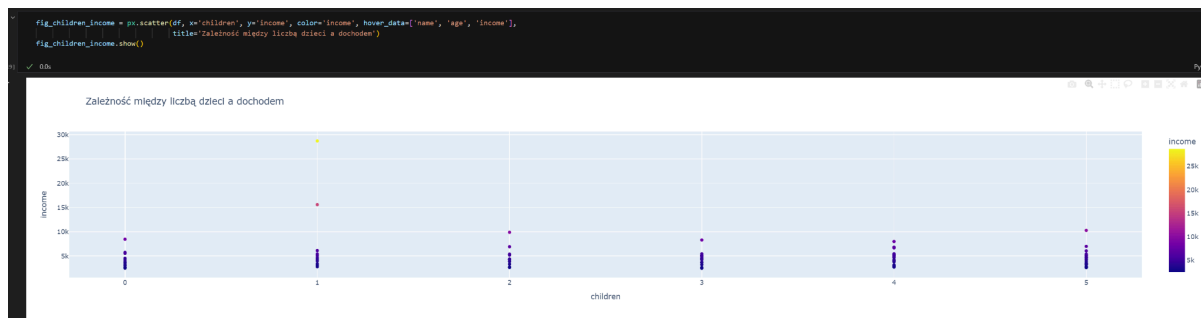
Rys. 5. Kod Źródłowy i efekt jego wykonania.

Utworzono interaktywny wykres pudełkowy, który przedstawia rozkład dochodów w podziale na wartości oznaczone jako odstające i normalne. Kolorami oznaczono różne grupy, aby ułatwić porównanie. Wizualizacja pozwala zaobserwować wpływ wartości odstających na dochody.



*Rys. 6. Kod Źródłowy i efekt jego wykonania.*

Stworzono interaktywny histogram przedstawiający rozkład wartości odstających w danych. Dodano margines w formie wykresu pudełkowego, aby uzupełnić analizę o dodatkowe informacje o rozkładzie. Kolory pozwalają łatwo odróżnić grupy wartości odstających i normalnych.



*Rys. 7. Kod Źródłowy i efekt jego wykonania.*

Utworzono interaktywną wizualizację punktową, przedstawiającą zależność między liczbą dzieci a dochodem. Kolorami oznaczono różne poziomy dochodów, a dodatkowe dane wyświetlane w chmurkach pozwalają na bardziej szczegółową analizę.

```
# 4. Dzielimy dane na dwie grupy: dane o dzieciach i dane o dochodach (t-SNE, UMAP)
from sklearn.manifold import TSNE
from umap import UMAP

columns_of_interest = ['age', 'income', 'outcome', 'savings', 'credit_score', 'spending_score']
X = df[columns_of_interest]

# Skalowanie danych (standardizacja)
scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)

# Użyjemy UMAP do redukcji wymiarowości
umap_model = UMAP(n_components=2, random_state=42)
umap_results = umap_model.fit_transform(X_scaled)

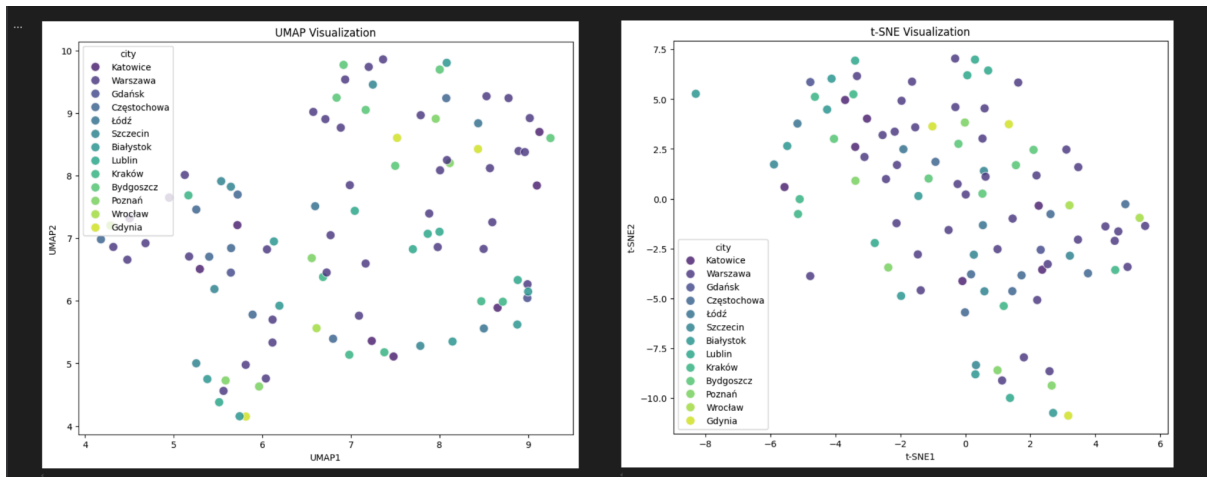
tsne_model = TSNE(n_components=2, random_state=42, perplexity=30, n_iter=1000)
tsne_results = tsne_model.fit_transform(X_scaled)

# Połączymy wyniki z danymi oryginalnymi
umap_df = pd.DataFrame(umap_results, columns=['UMAP1', 'UMAP2'])
tsne_df = pd.DataFrame(tsne_results, columns=['t-SNE1', 't-SNE2'])
df_umap = pd.concat([df, umap_df], axis=1)
df_tsne = pd.concat([df, tsne_df], axis=1)

# UMAP
plt.figure(figsize=(10, 8))
sns.scatterplot(data=df_umap, x='UMAP1', y='UMAP2', hue='city', palette='viridis', s=100, alpha=0.8)
plt.title('UMAP Visualization')
plt.show()

# t-SNE
plt.figure(figsize=(10, 8))
sns.scatterplot(data=df_tsne, x='t-SNE1', y='t-SNE2', hue='city', palette='viridis', s=100, alpha=0.8)
plt.title('t-SNE Visualization')
plt.show()
```

*Rys. 8. Kod Źródłowy.*



Rys. 9. Efekt wykonania kodu z poprzedniego rysunku.

Przygotowano dane, wybierając kluczowe kolumny i skalując je za pomocą StandardScaler. Następnie zastosowano UMAP i t-SNE do redukcji wymiarowości do 2D. Wyniki obydwu algorytmów zwizualizowano, kolorując punkty według miasta, co pozwala lepiej zobrazować złożone relacje w danych. Wykorzystano wykresy punktowe z seaborn do przedstawienia wyników dla UMAP i t-SNE, co umożliwia analizę zróżnicowania między grupami.

```
# 3. Wczyść i przetwóżyć dane
df = pd.read_csv('data.csv')

# 4. Wybrać kolumny do analizy
df = df[['income', 'outcome', 'city']]

# 5. Skalowanie danych
scaler = StandardScaler()
df[['income', 'outcome']] = scaler.fit_transform(df[['income', 'outcome']])

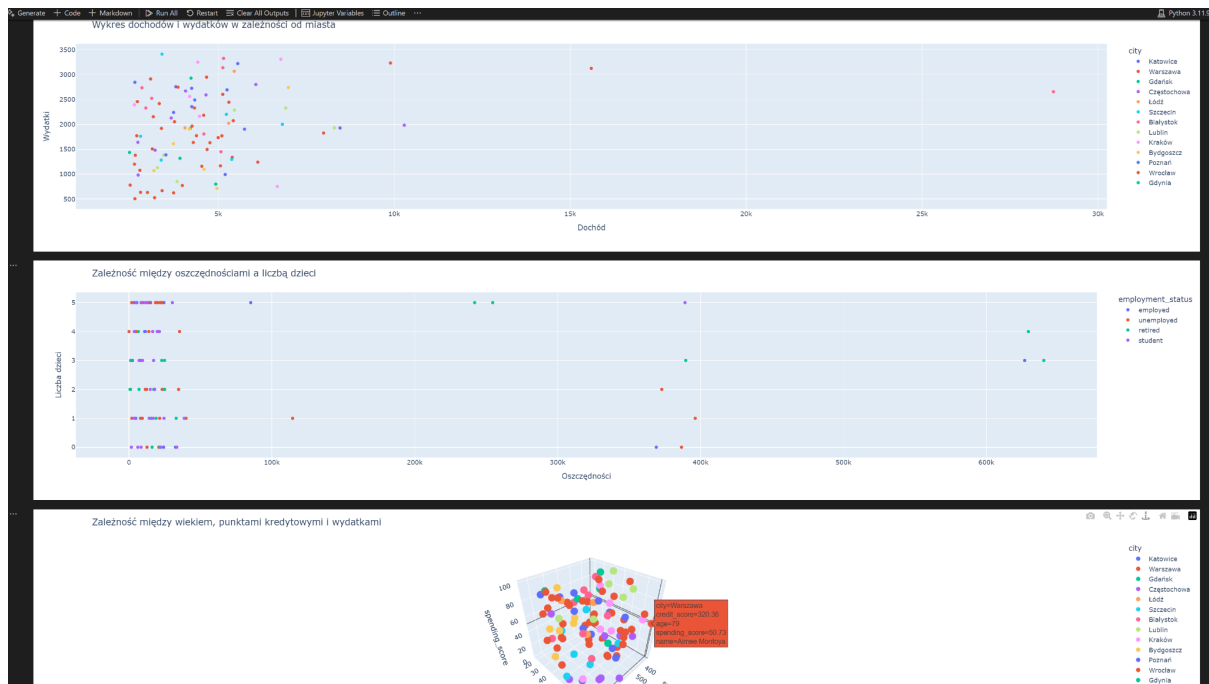
# 6. Wykresy punktowe z UMAP
fig = plt.figure(figsize=(10, 10))
fig = px.scatter(df, x='income', y='outcome', color='city', hover_data=[name, 'city'])
fig.update_layout(title='Wykres dochodów i wydatków w zależności od miasta',
                  xaxis_title='Dochód',
                  yaxis_title='Wydatki')
fig.show()

# 7. Wykresy punktowe z t-SNE
fig = px.scatter(df, x='income', y='outcome', color='city', hover_data=[name, 'city'])
fig.update_layout(title='Wykres dochodów i wydatków w zależności od miasta',
                  xaxis_title='Dochód',
                  yaxis_title='Wydatki')
fig.show()

# 8. Wykresy punktowe z UMAP
fig = px.scatter(df, x='income', y='outcome', color='city', hover_data=[name, 'city'])
fig.update_layout(title='Wykres dochodów i wydatków w zależności od miasta',
                  xaxis_title='Dochód',
                  yaxis_title='Wydatki')
fig.show()

# 9. Wykresy punktowe z t-SNE
fig = px.scatter(df, x='income', y='outcome', color='city', hover_data=[name, 'city'])
fig.update_layout(title='Wykres dochodów i wydatków w zależności od miasta',
                  xaxis_title='Dochód',
                  yaxis_title='Wydatki')
fig.show()
```

Rys. 10. Kod Źródłowy.



Rys. 11. Efekt wykonania kodu z poprzedniego rysunku.

W pierwszym kroku utworzono interaktywną wizualizację 2D, przedstawiającą zależność między dochodami a wydatkami, z uwzględnieniem miasta. Następnie stworzono drugi wykres 2D pokazujący związek między oszczędnościami a liczbą dzieci, z różnicowaniem na status zatrudnienia. Na końcu wykonano interaktywny wykres 3D, który ukazuje zależność między wiekiem, punktami kredytowymi i wydatkami, również z uwzględnieniem miasta.



Rys. 12. Kod Źródłowy i efekt jego wykonania.

Zaimportowano niezbędne biblioteki i przygotowano dane numeryczne, uwzględniając kolumny income, outcome, savings. Następnie obliczono

macierz korelacji dla tych zmiennych, a wynik wizualizowano za pomocą interaktywnego heatmap. Wykres przedstawia związek między zmiennymi poprzez kolorystyczny rozkład wartości korelacji.

```
# 7. Przeprowadź testy statystyczne dla analizy różnic w grupach.

from scipy import stats
from statsmodels.formula.api import ols
from statsmodels.stats.anova import anova_lm
from scipy.stats import chi2_contingency

# Test T-studenta
group1 = df[df['employment_status'] == 'employed']['income']
group2 = df[df['employment_status'] == 'retired']['income']
t_stat, p_value = stats.ttest_ind(group1, group2)

print(f"Test T-studenta: (t_stat)")
print(f"P-wartość: (p_value)")
if p_value < 0.05:
    print("Istnieje istotna różnica w średnich income między grupą zatrudnionych a emerytami.")
else:
    print("Brak istotnej różnicy w średnich income między grupą zatrudnionych a emerytami.")

# Test ANOVA
model = ols('outcome ~ employment_status', data=df).fit()
anova_table = anova_lm(model)

print(anova_table)

if anova_table['PR(>F)'][0] < 0.05:
    print("Istnieje istotna różnica w średnich outcome między przynajmniej dwoma grupami employment_status. (Pracujący zarabiają więcej niż emeryci)")
else:
    print("Brak istotnej różnicy w średnich outcome między przynajmniej dwoma grupami employment_status.")

# Test Chi-kwadrat
contingency_table = pd.crosstab(df['employment_status'], df['city'])
chi2_stat, p_val, dof, expected = chi2_contingency(contingency_table)

print(f"Chi-kwadrat statystyka: (chi2_stat)")
print(f"P-wartość: (p_val)")

if p_val < 0.05:
    print("Istnieje istotna zależność między employment_status a city. (np. w Warszawie jest więcej bezrobotnych)")
else:
    print("Brak istotnej zależności między employment_status a city.")

✓ 0/0 Python
```

Teststatystyka: 2.475167674044107  
P-wartość: 0.01789215021119085  
Istnieje istotna różnica w średnich income między grupą zatrudnionych a emerytami.  
df non\_em emeryci PR(>F)  
employment\_status 3.0 5.985835e+00 1.995812e+00 3.056484 0.015169  
Residuals 90.0 1.237616e+07 5.405795e+05 NaN NaN  
Istnieje istotna różnica w średnich outcome między przynajmniej dwoma grupami employment\_status. (Pracujący zarabiają więcej niż emeryci)  
Chi-kwadrat statystyka: 188.93872127781896  
P-wartość: 2.810950612191500e-09  
Istnieje istotna zależność między employment\_status a city. (np. w Warszawie jest więcej bezrobotnych)

Rys. 13. Kod Źródłowy i efekt jego wykonania.

Ten kod przeprowadza różne testy statystyczne, aby zidentyfikować istotne różnice w danych:

Test T-studenta sprawdza, czy istnieje istotna różnica w średnich dochodach między grupą zatrudnionych a emerytami. Wynik pokazuje, że różnice są istotne, ponieważ p-wartość jest poniżej 0.05.

Test ANOVA analizuje, czy różnice w średnich wydatkach są istotne pomiędzy przynajmniej dwoma grupami zatrudnienia. W tym przypadku różnice są istotne, a Pracujący mają wyższe średnie dochody niż emeryci.

Test Chi-kwadrat sprawdza zależność między statusem zatrudnienia a miastem, ujawniając, że w Warszawie jest więcej bezrobotnych.

Każdy test dostarcza kluczowych informacji na temat struktury danych, co pomaga w głębszym zrozumieniu ich dystrybucji oraz związków.

### 3. Wnioski

Istotne różnice w dochodach i wydatkach między grupami zatrudnienia: Test T-studenta oraz ANOVA pokazują, że istnieje znacząca różnica w dochodach między pracującymi a emerytami, a także w średnich wydatkach pomiędzy tymi grupami. Pracujący mają wyższe dochody i większe wydatki niż emeryci.

Zależność między miastem a statusem zatrudnienia: Test Chi-kwadrat wskazuje na istotną zależność między employment\_status a miastem, co może sugerować, że w większych miastach, takich jak Warszawa, więcej osób jest bezrobotnych.

Wpływ liczby dzieci i oszczędności na dochody: Analiza wykazała, że liczba dzieci i oszczędności wpływają na dochody, co może wskazywać na różnice w zarządzaniu finansami w zależności od demografii i sytuacji rodzinnej.