

Dokumentacja Lab3

Wstępna analiza danych

Informacje o danych:

- Liczba wierszy: 4739
- Liczba kolumn: 15

Kolumny zawierają dane numeryczne oraz katégoryczne:

- Numeryczne: rownames (identyfikator), score (wynik do przewidzenia), unemp, wage, distance, tuition, education
- Katégoryczne: gender, ethnicity, fcollege, mcollege, home, urban, income, region

Brakujące dane:

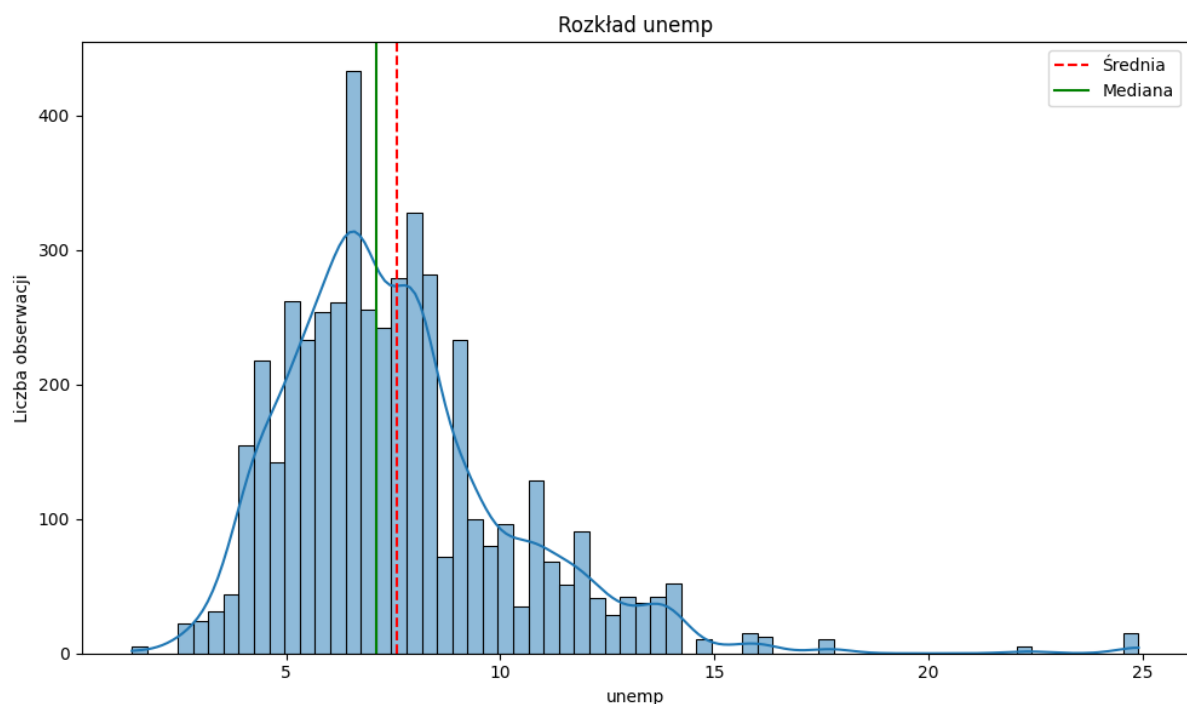
- Wszystkie kolumny mają pełny zbiór danych, nie ma brakujących wartości.
- Brak danych może być reprezentowany poprzez wartości 0 w kolumnach numerycznych.

Analiza statystyczna:

Zmienne numeryczne:

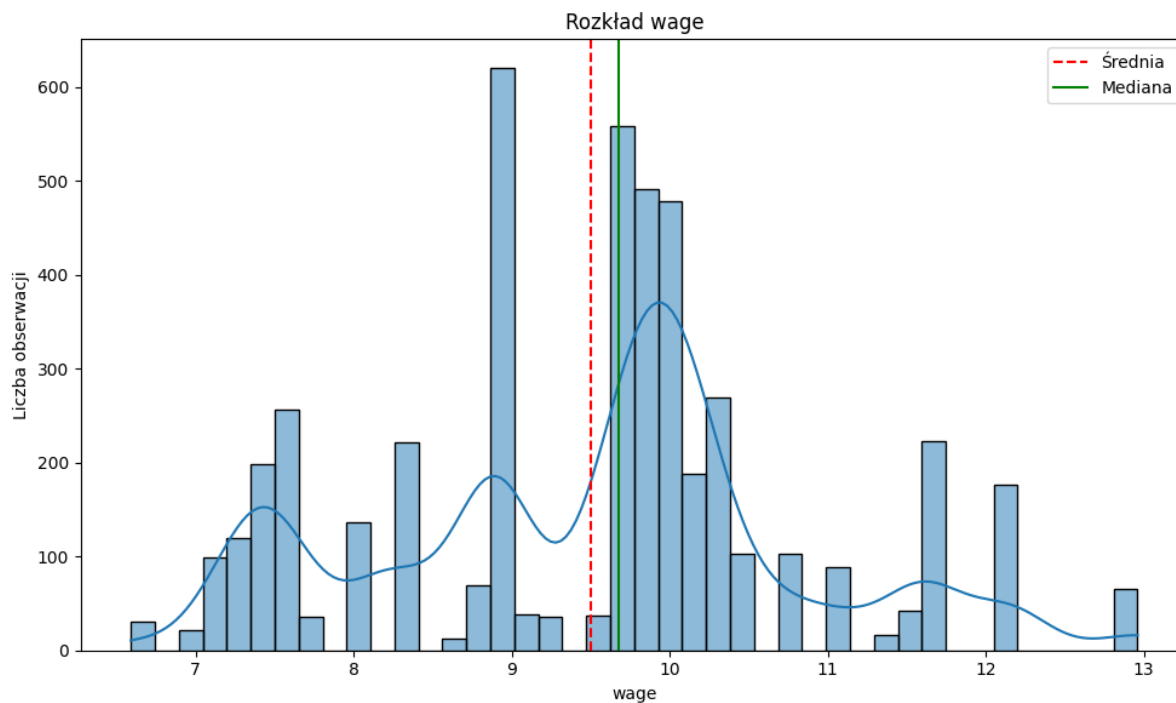
1. Unemployment (unemp):

- Średnia: 7.6
- Mediana: 7.1
- Zakres wartości: 1.4 - 24.9



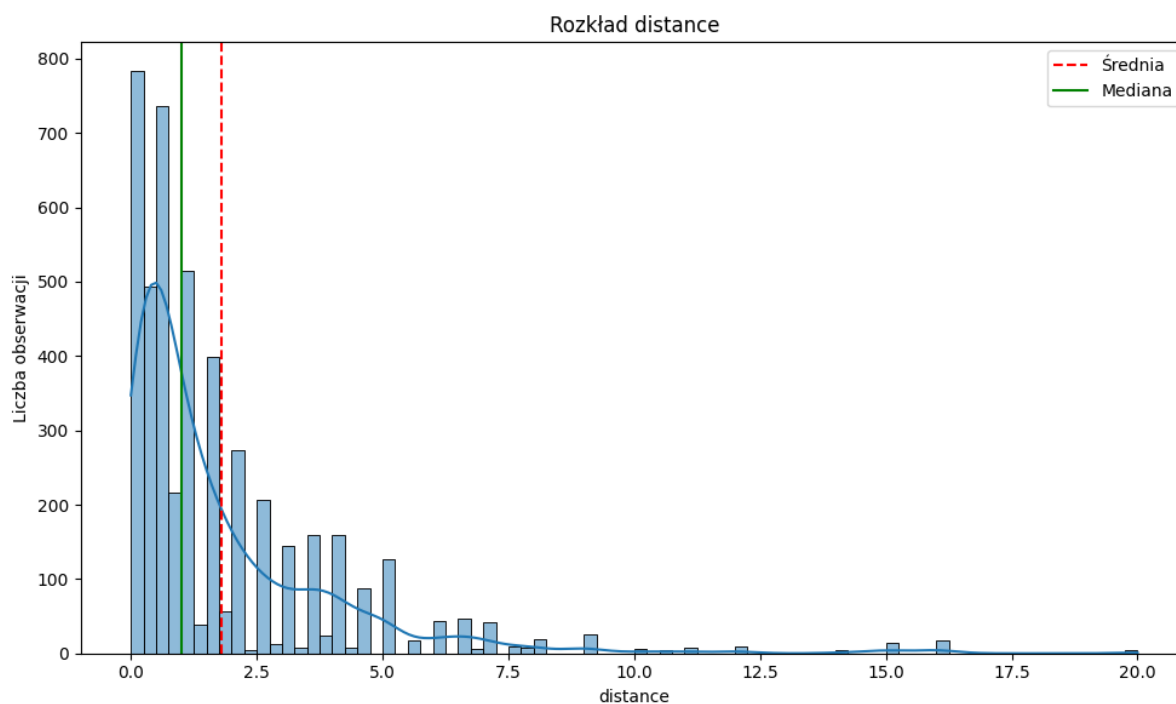
2. Wynagrodzenie (wage):

- Średnia: 9.5
- Mediana: 9.7
- Zakres: 6.59 - 12.96



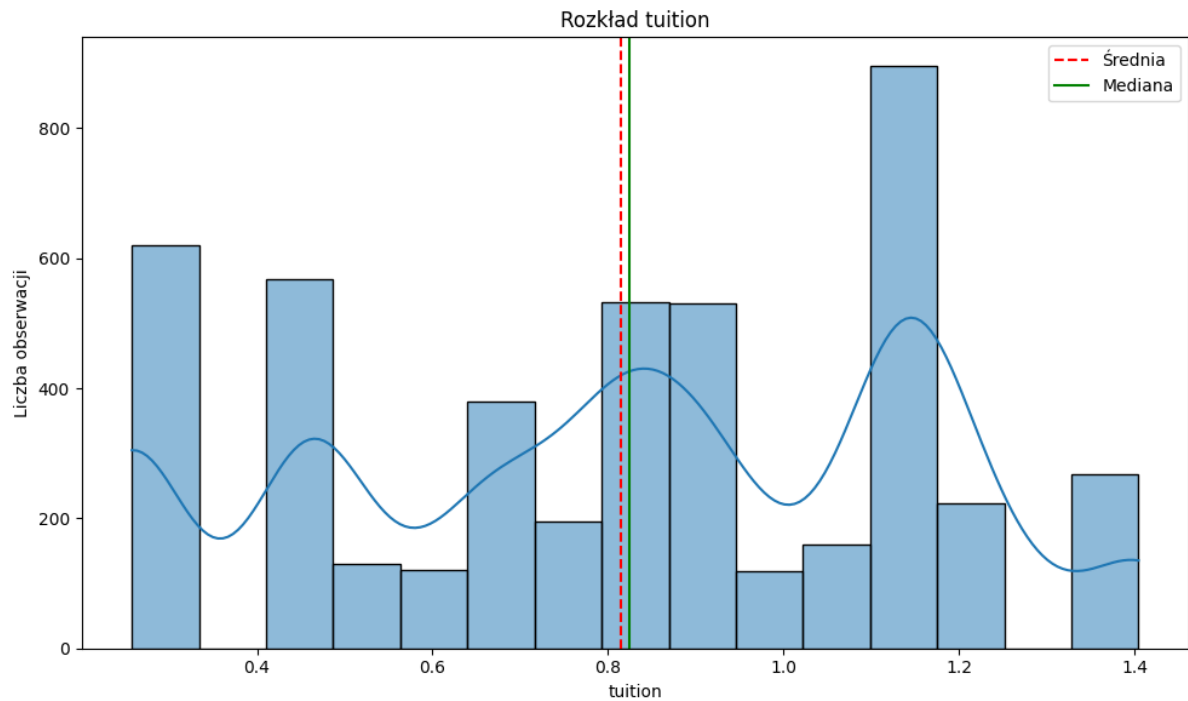
3. Dystans do uczelni (distance):

- Średnia: 1.8
- Mediana: 1.0
- Zakres: 0.0 - 20.0



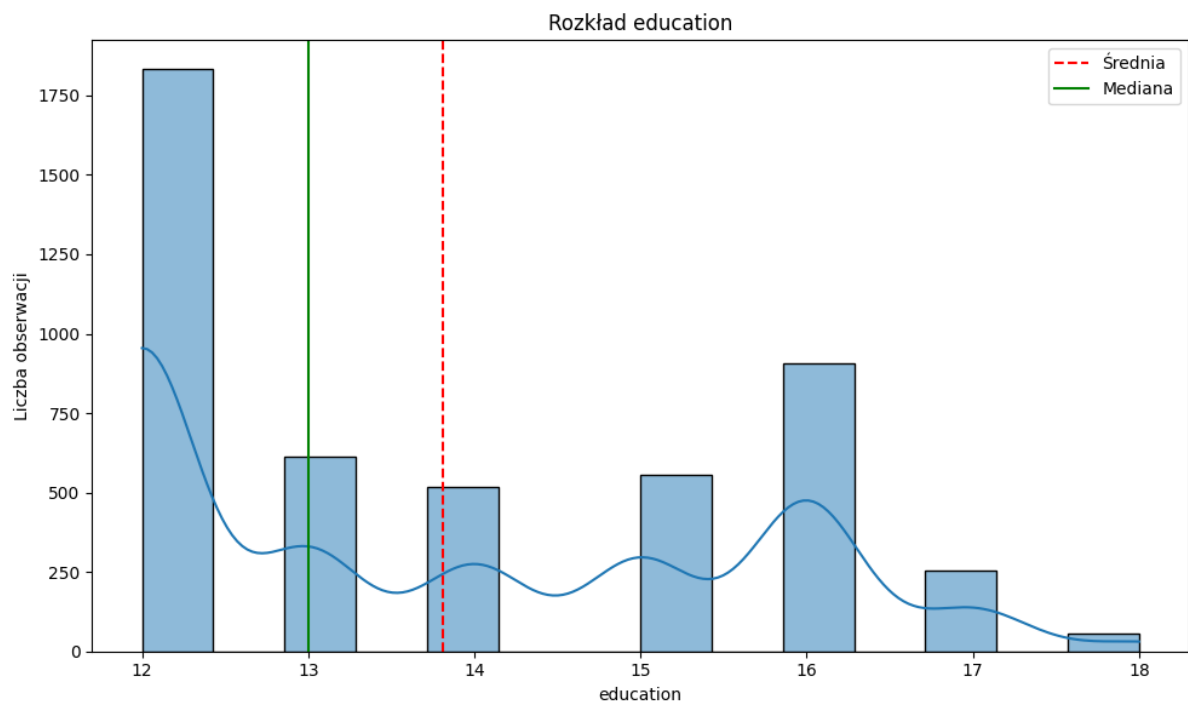
4. Czesne (tuition):

- Średnia: 0.81
- Mediana: 0.82
- Zakres: 0.26 - 1.40



5. Edukacja (education):

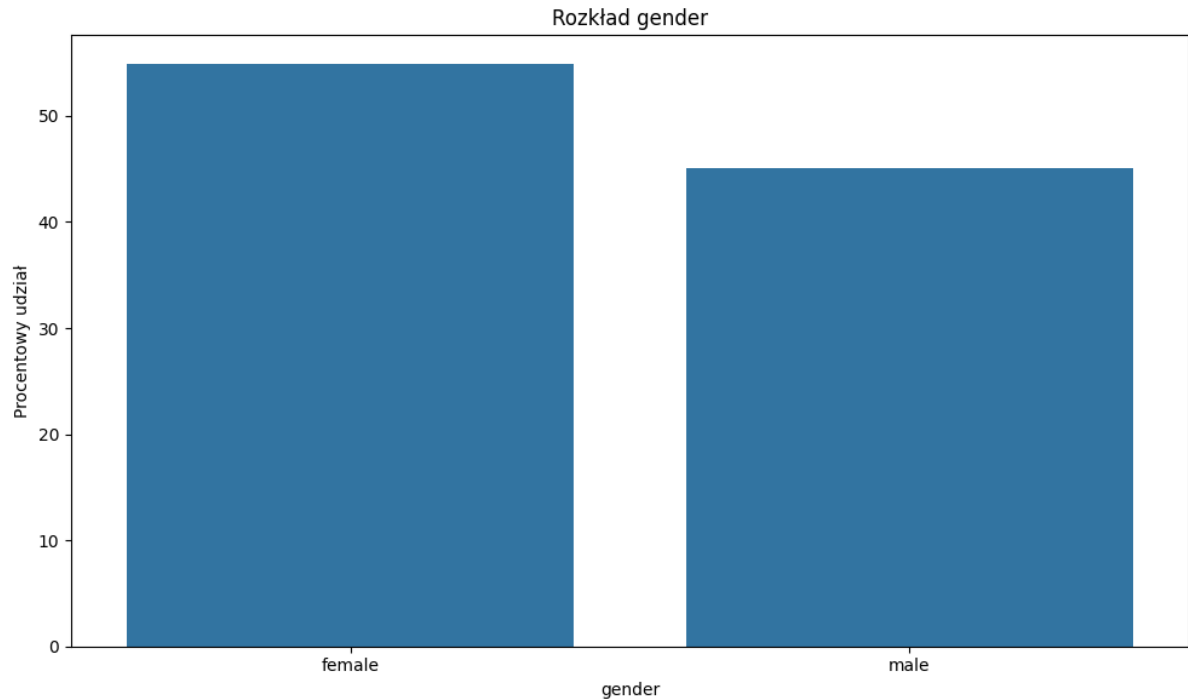
- Średnia: 13.8
- Mediana: 13
- Zakres: 12 - 18



Zmienne kategoryczne:

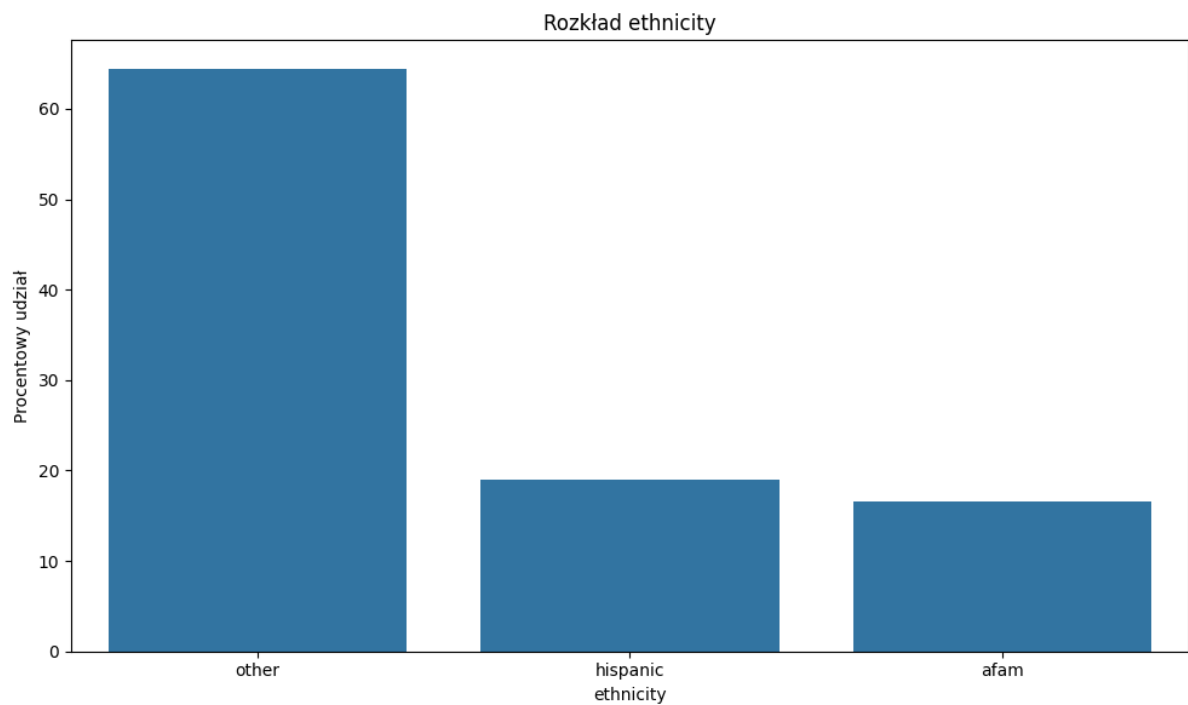
1. Płeć (gender):

- Kategoria: male, female



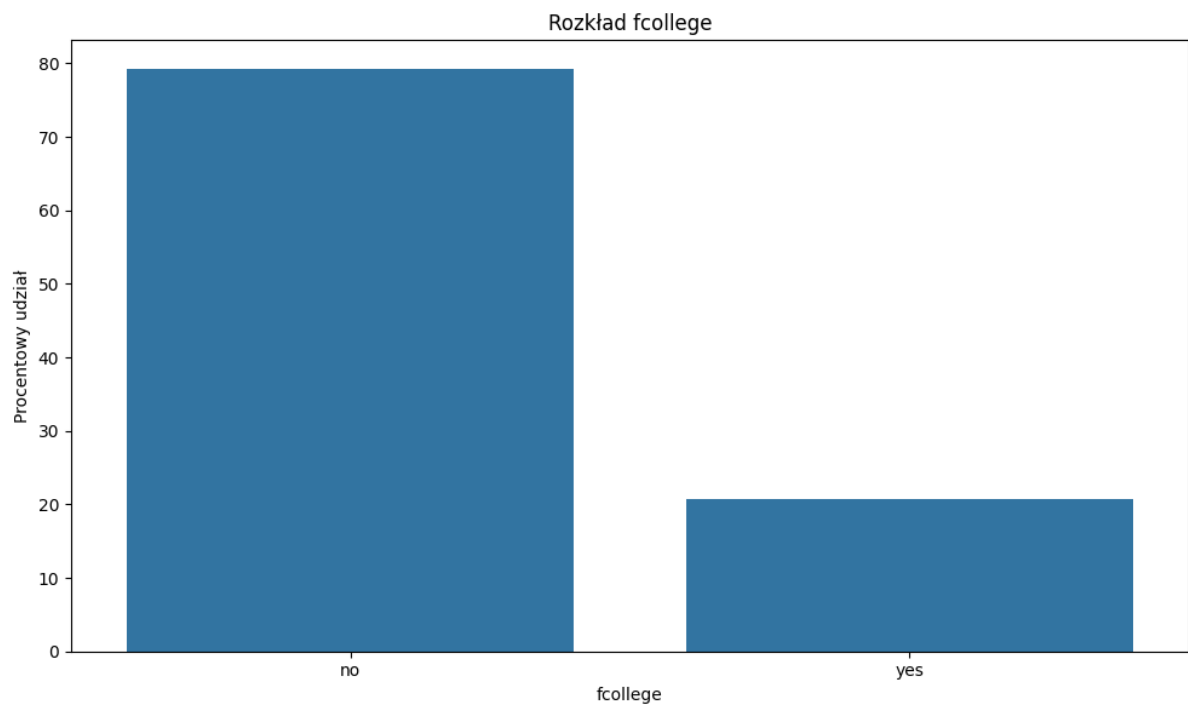
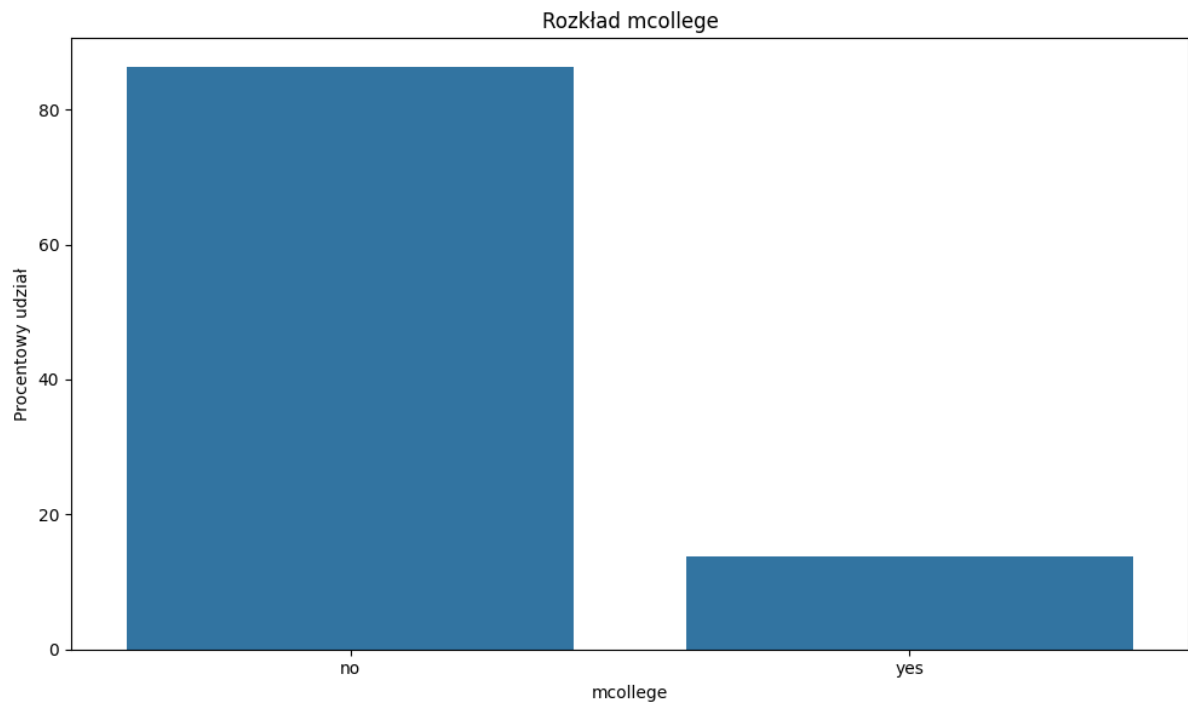
2. Etniczność (ethnicity):

- Kategoria: afam (Afroamerykanin), hispanic, other



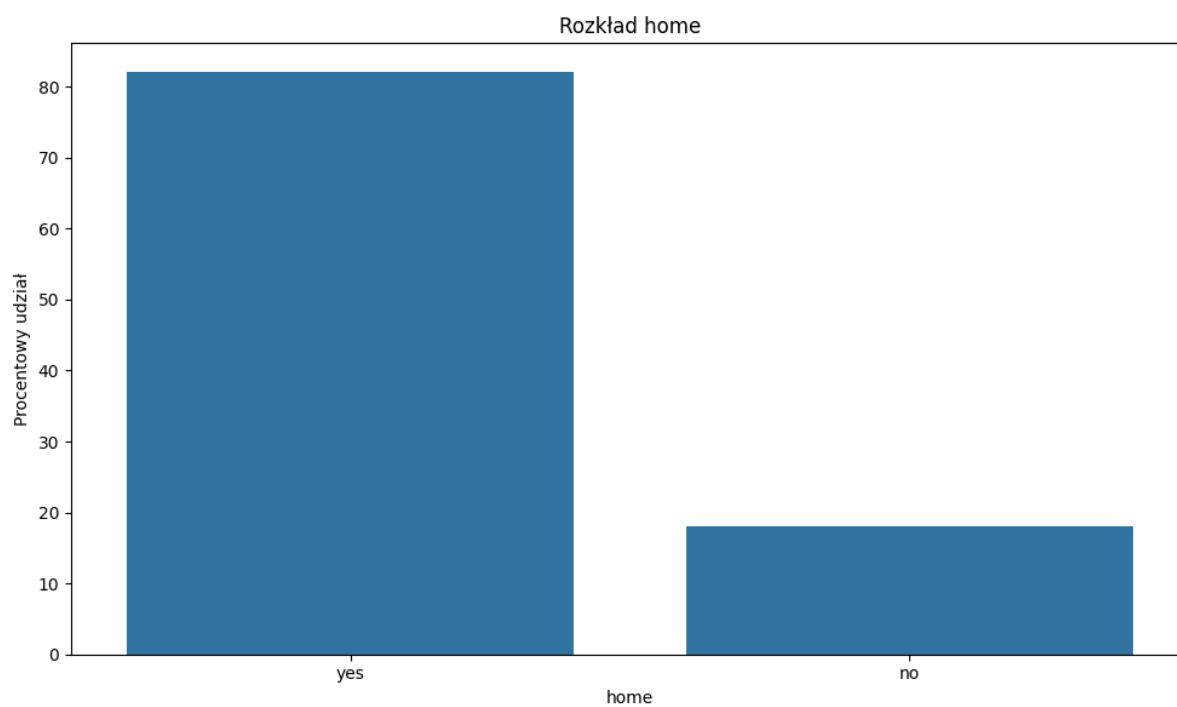
3. Rodzic ukończył studia (fcollege, mcollege):

- Kategoria: yes, no



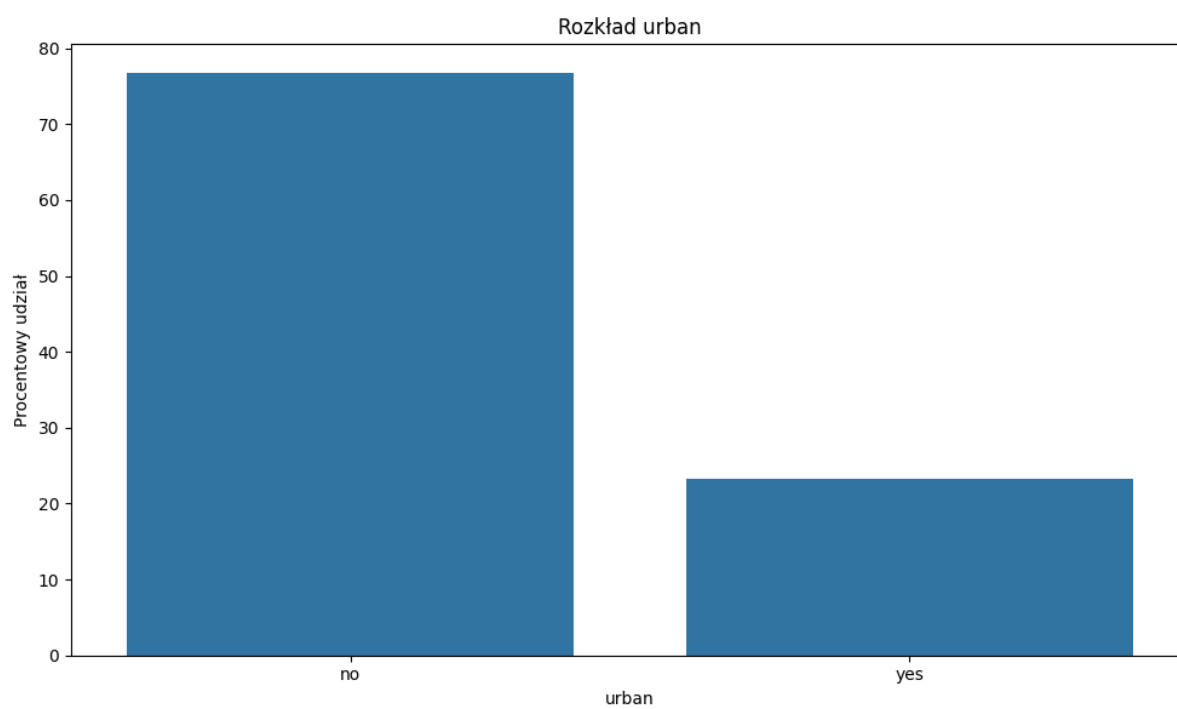
4. Dom (home):

- Kategoria: yes, no



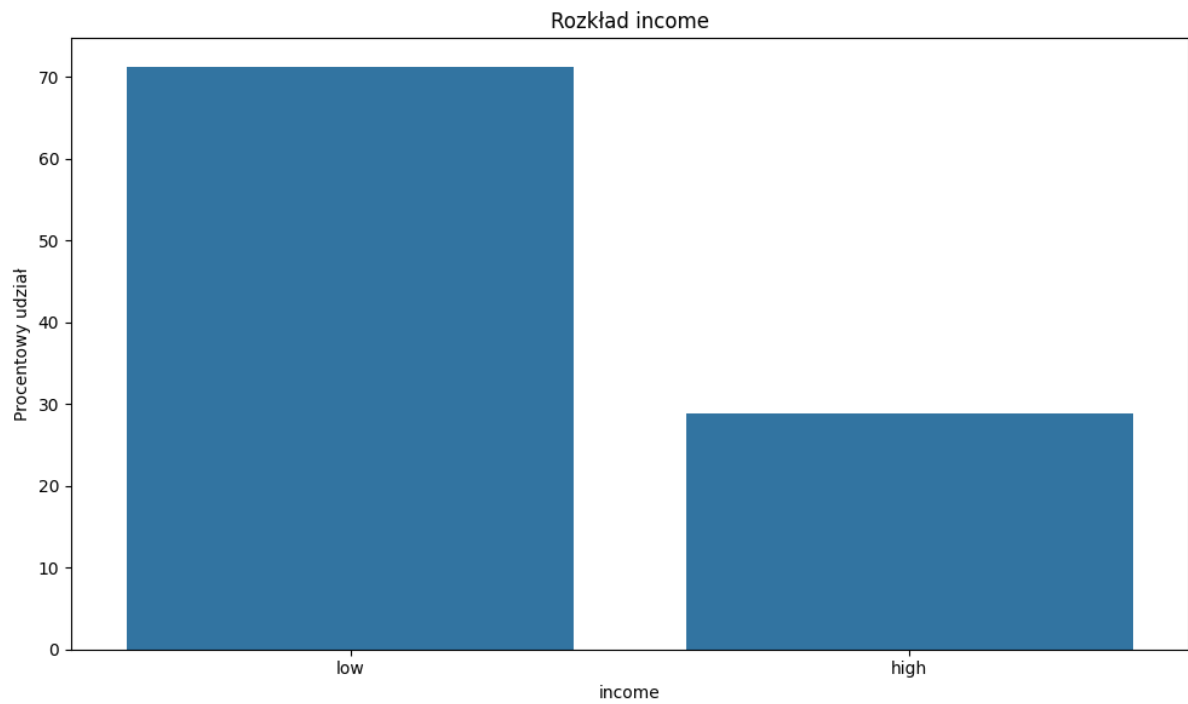
5. Obszar miejski (urban):

- Kategoria: yes, no



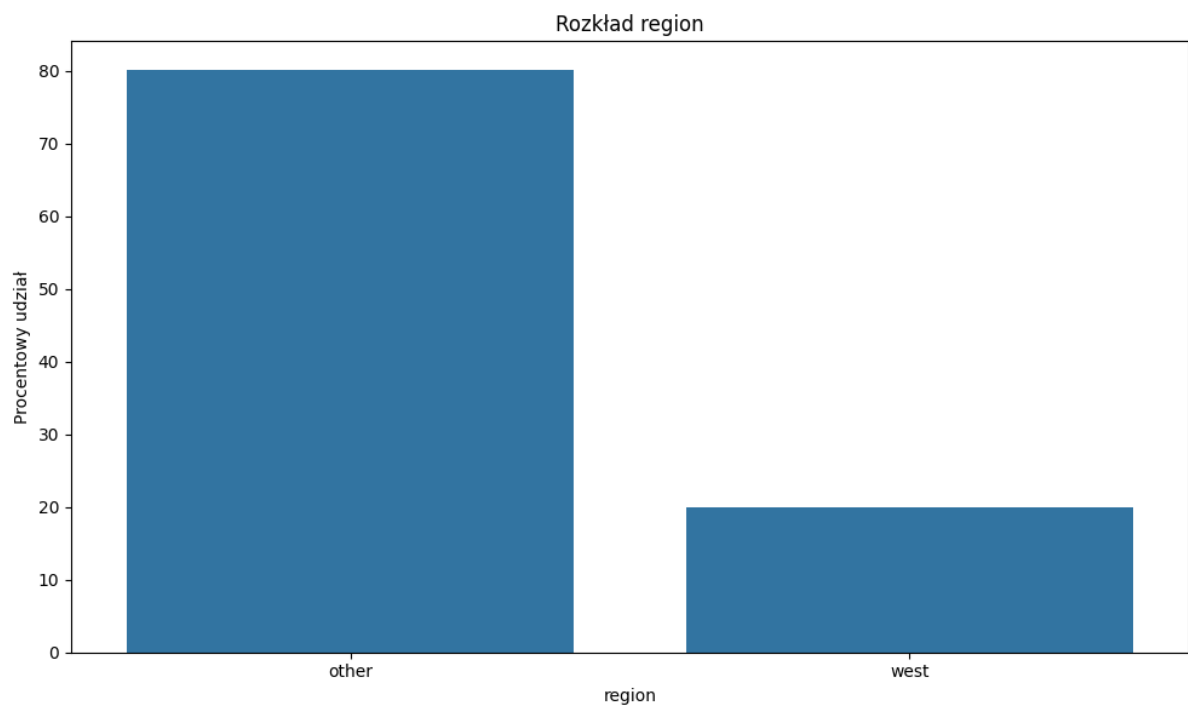
6. Dochód rodziny (income):

- Kategoria: high, low



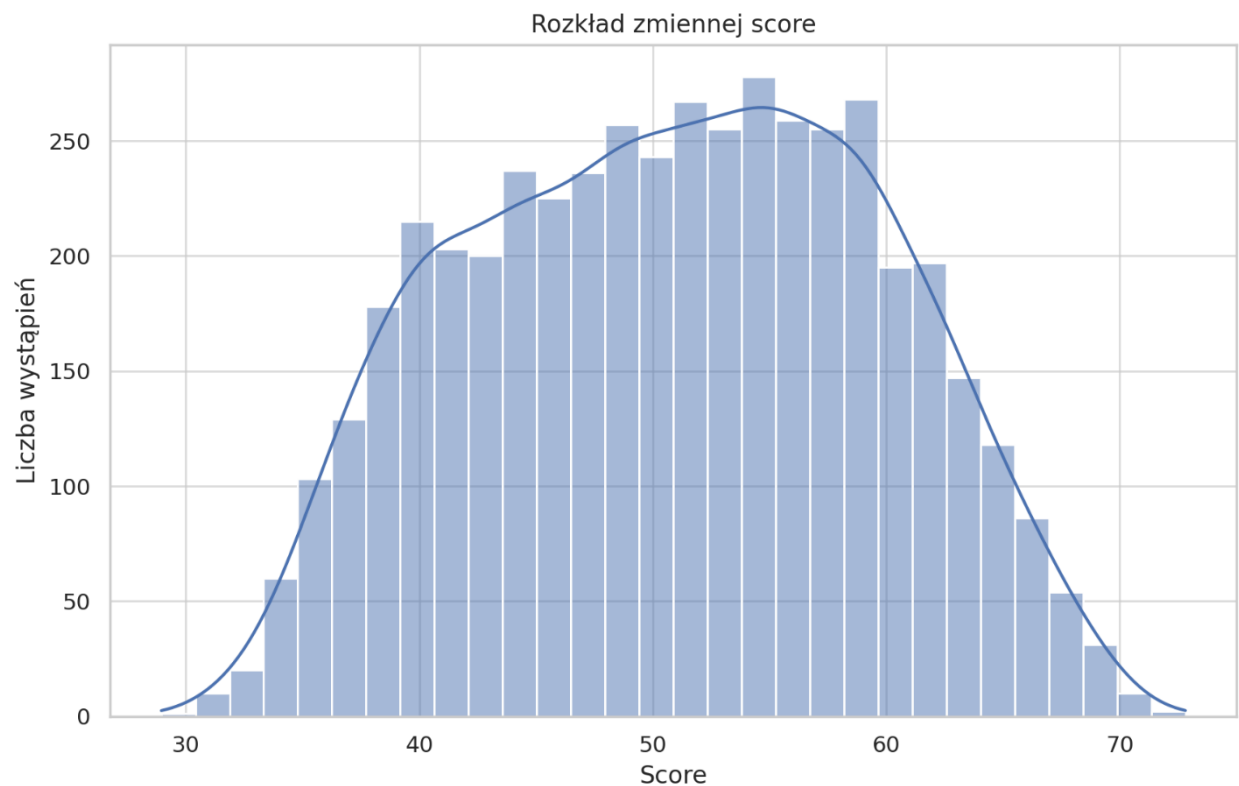
7. Region (region):

- Kategoria: west, other

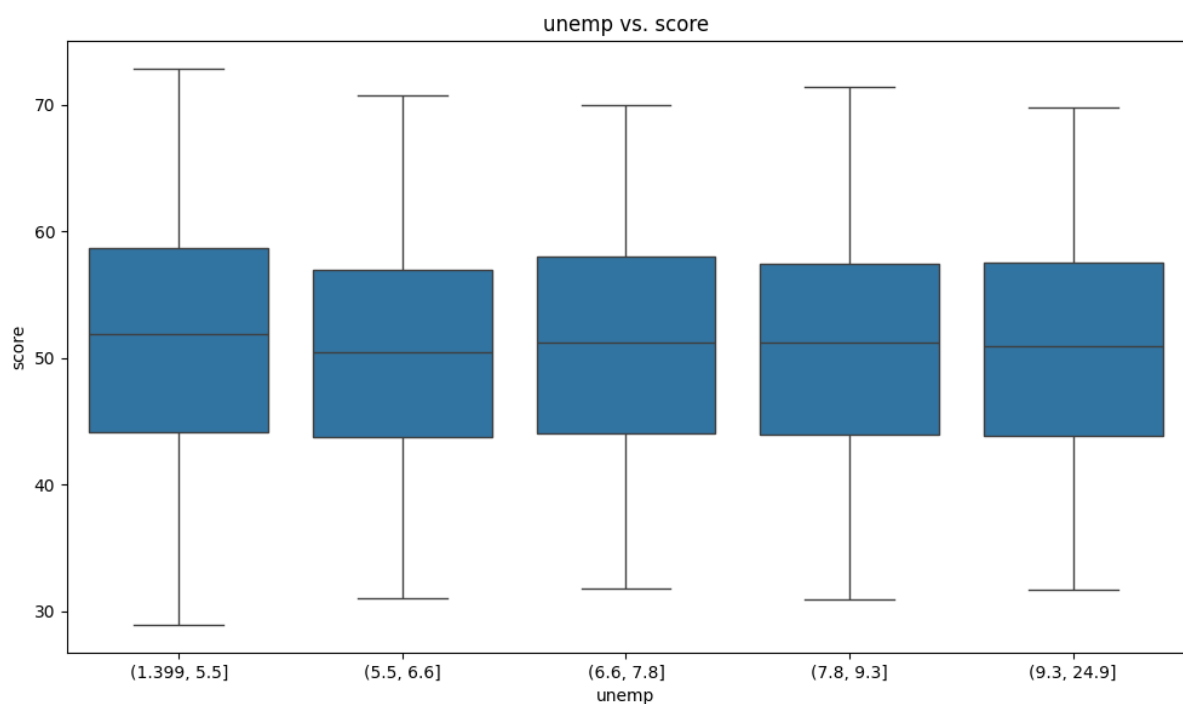


Wizualizacja zależności danych

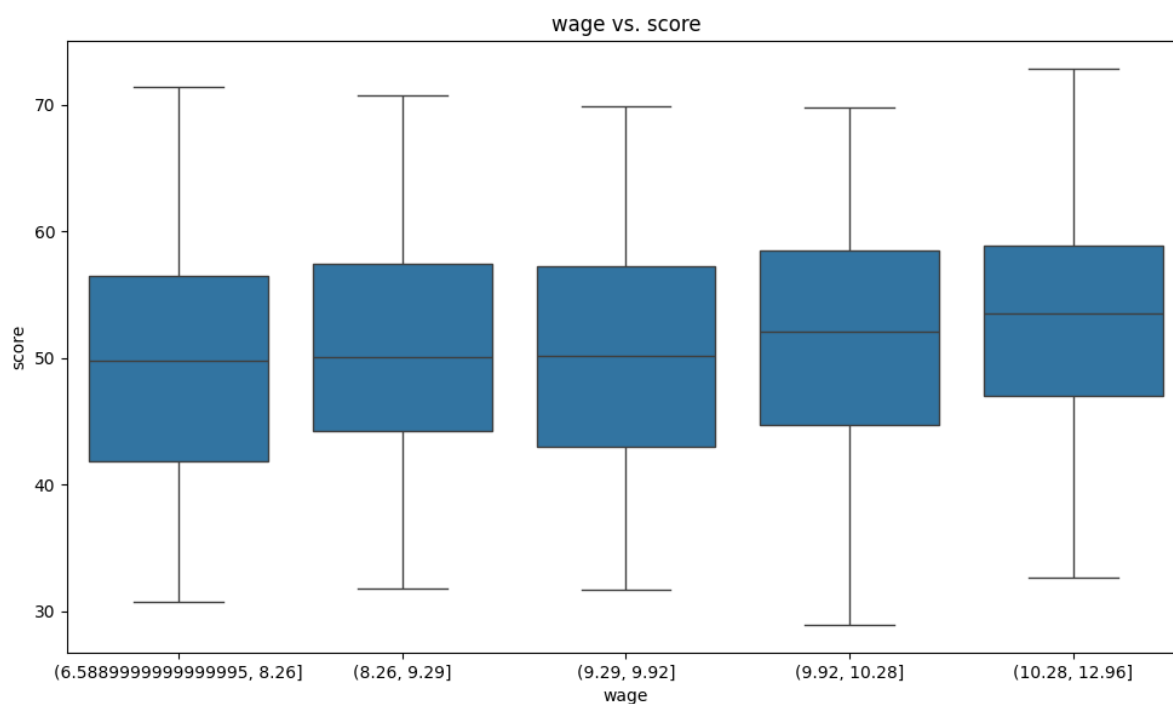
Rozkład zmiennej score jest zbliżony do normalnego, jednak istnieją pewne asymetrie i wartości odstające. Większość wyników znajduje się w przedziale od około 40 do 60.



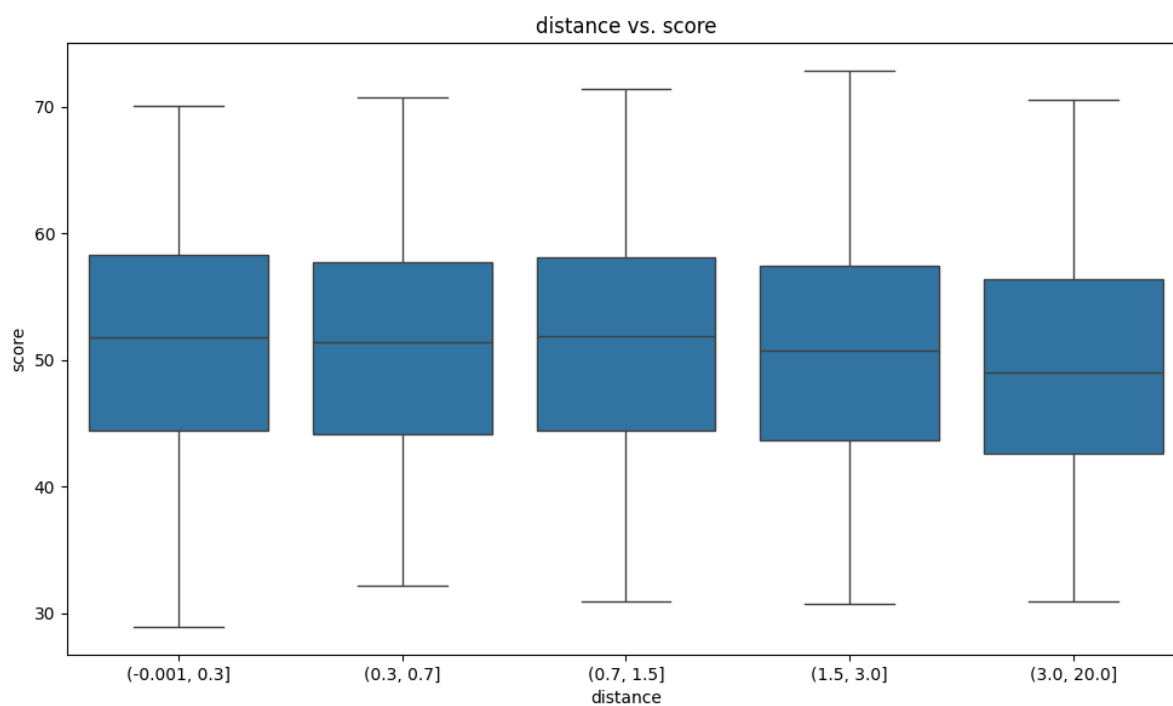
Grupy dla poziomów bezrobocia pokazują, że wyższy poziom bezrobocia nie ma jednoznacznej korelacji ze wzrostem lub spadkiem wartości zmiennej score.



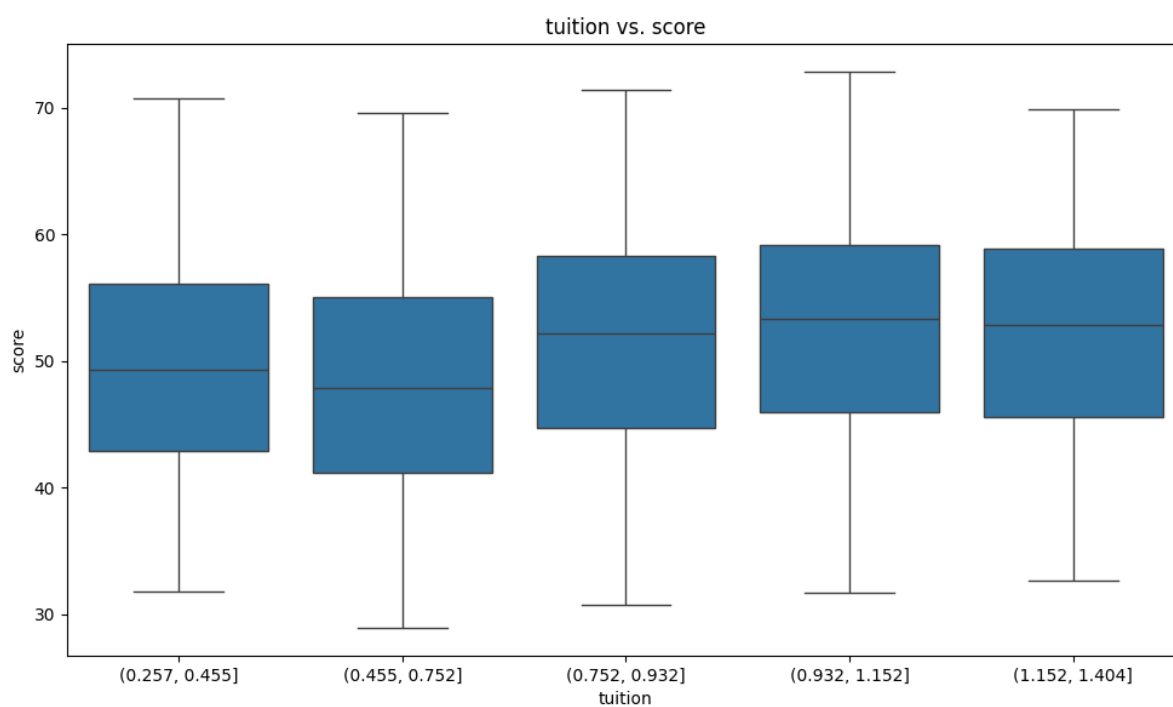
Wyższe wynagrodzenie ma umiarkowaną zależność z wyższymi wynikami score, jednak zależność nie jest linearna.



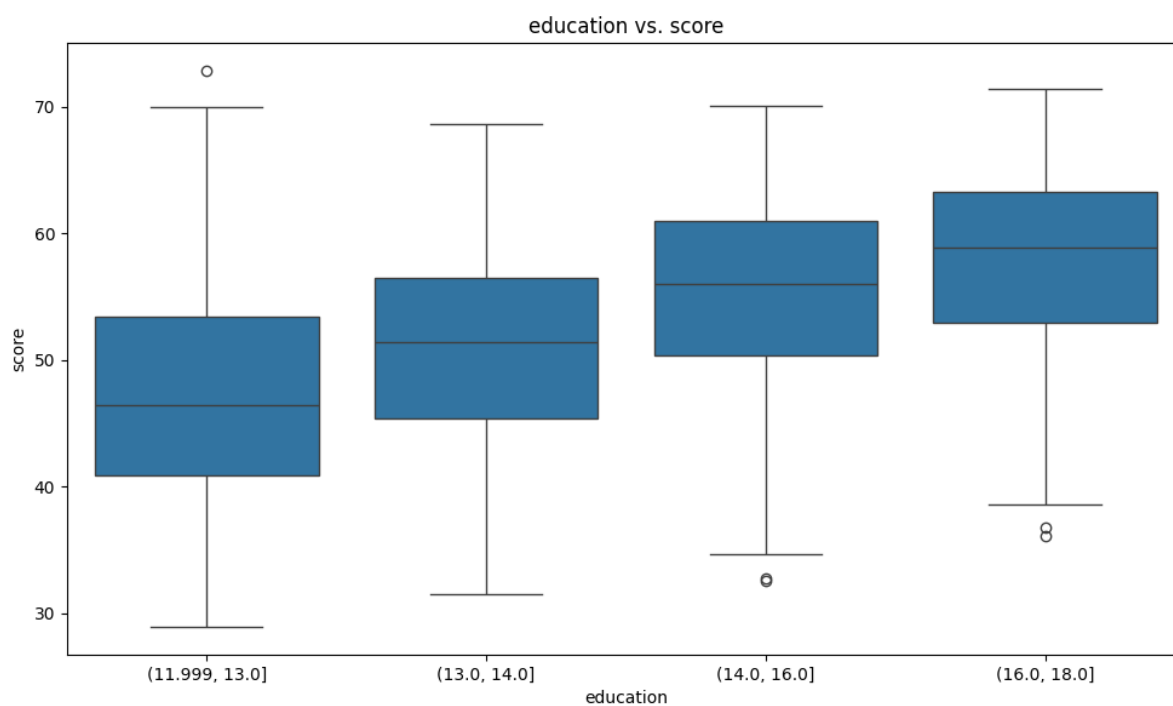
Dłuższy dystans ma tendencję do obniżenia wyniku, jednak w grupach powyżej 2.5 km wynik staje się bardziej zróżnicowany.



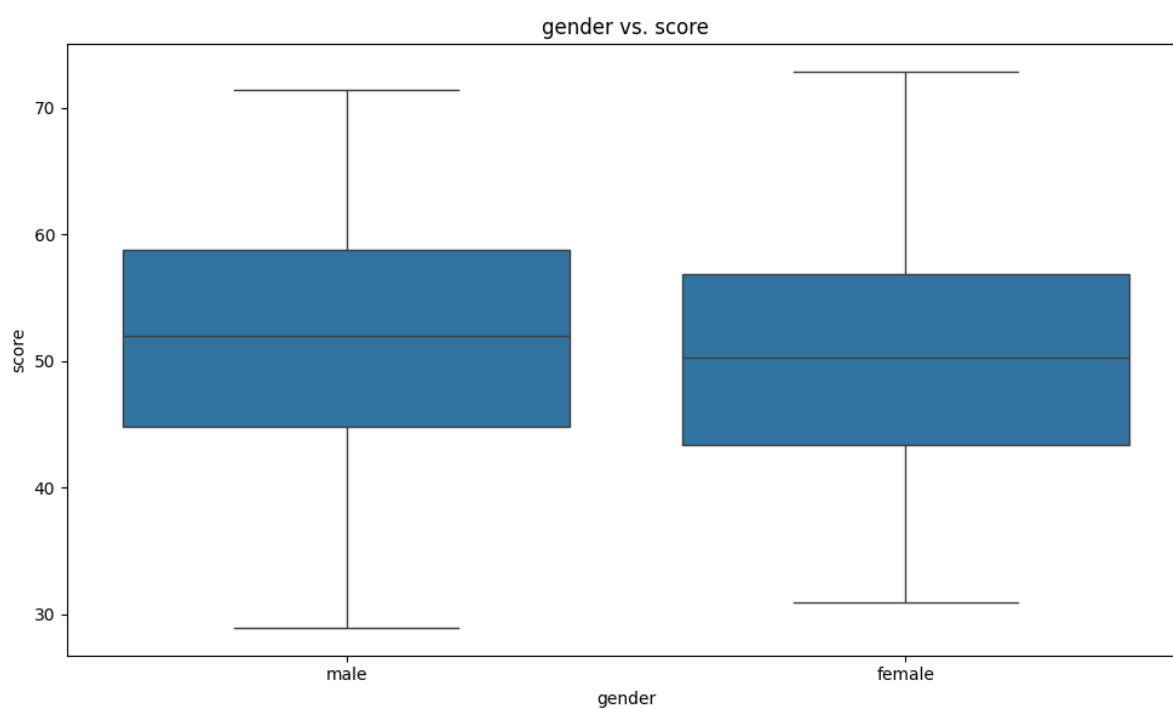
Wyższe czesne koreluje z nieznacznie wyższym wynikiem score.



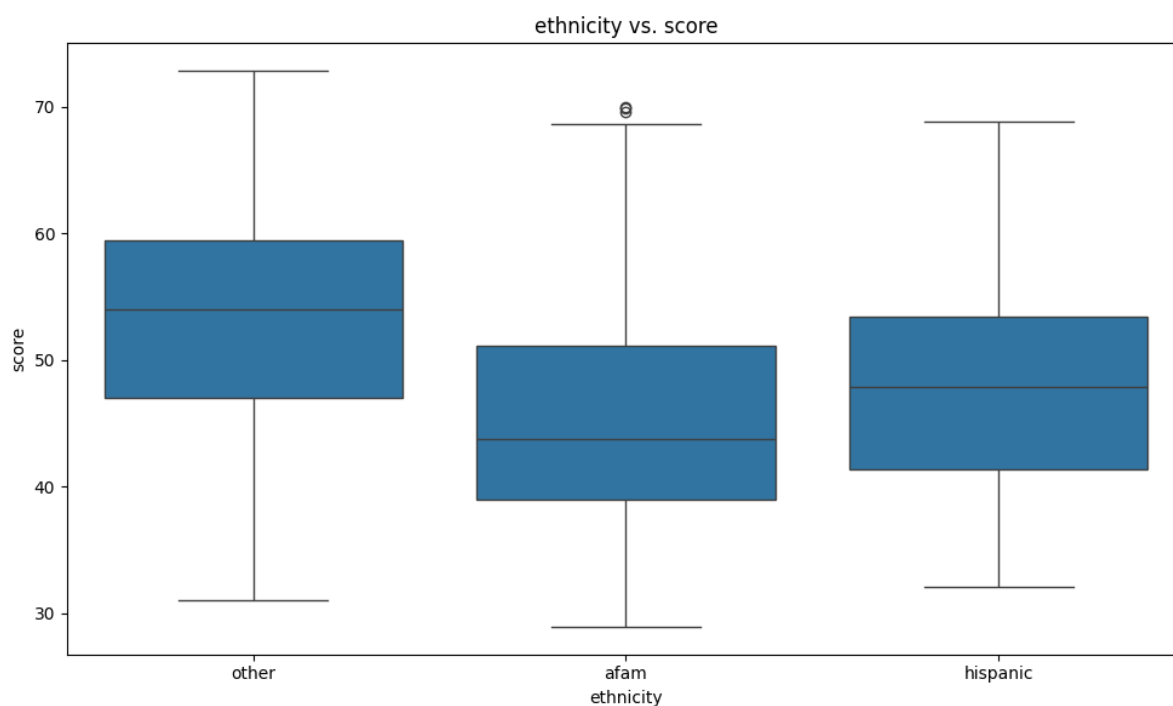
Większa wartość edukacji pokazuje trend w kierunku wyższych wyników score.



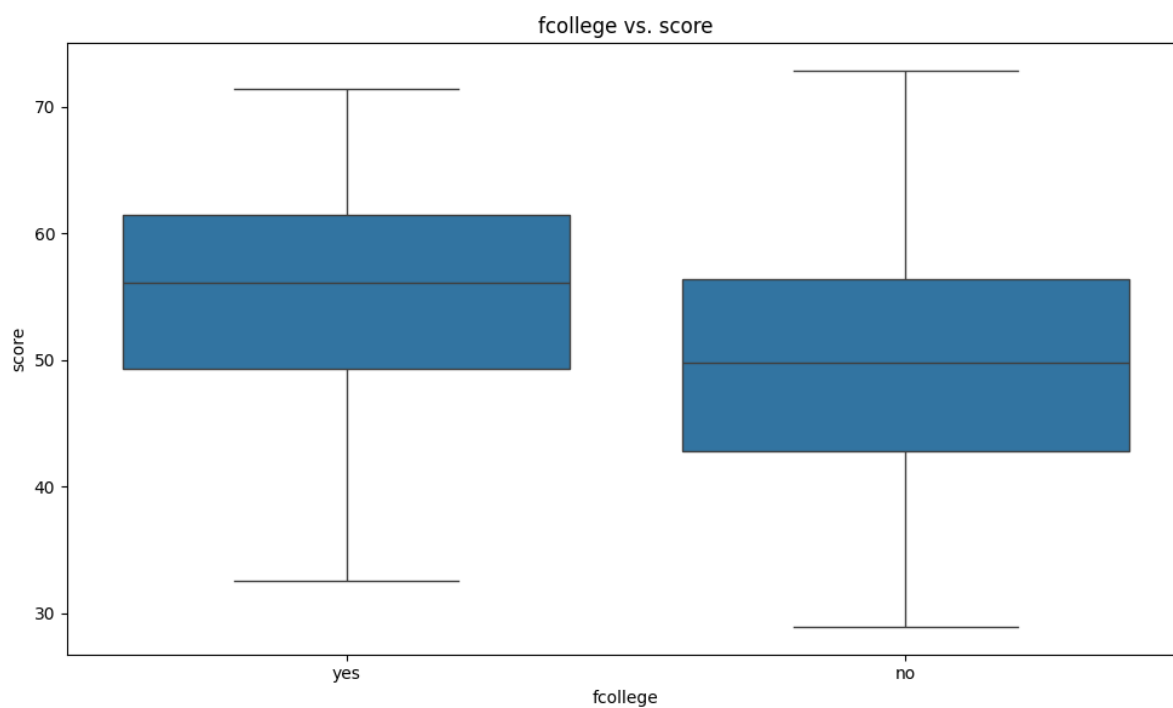
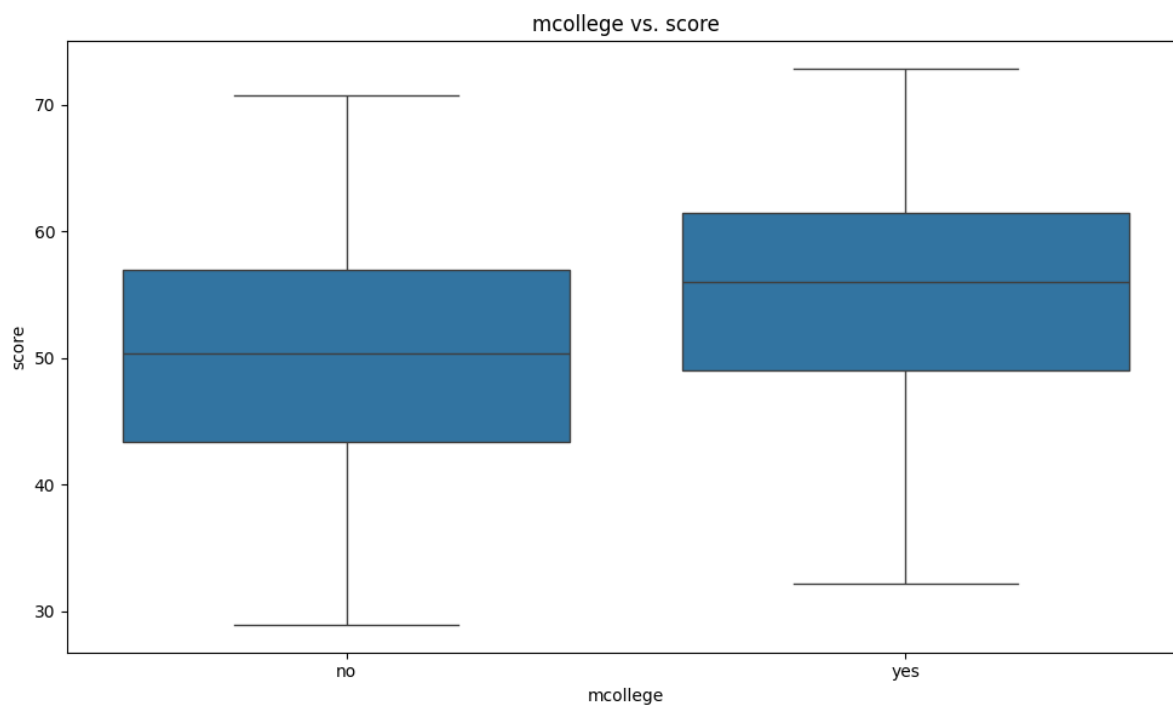
Kobiety mają nieco wyższe wyniki score w porównaniu do mężczyzn.



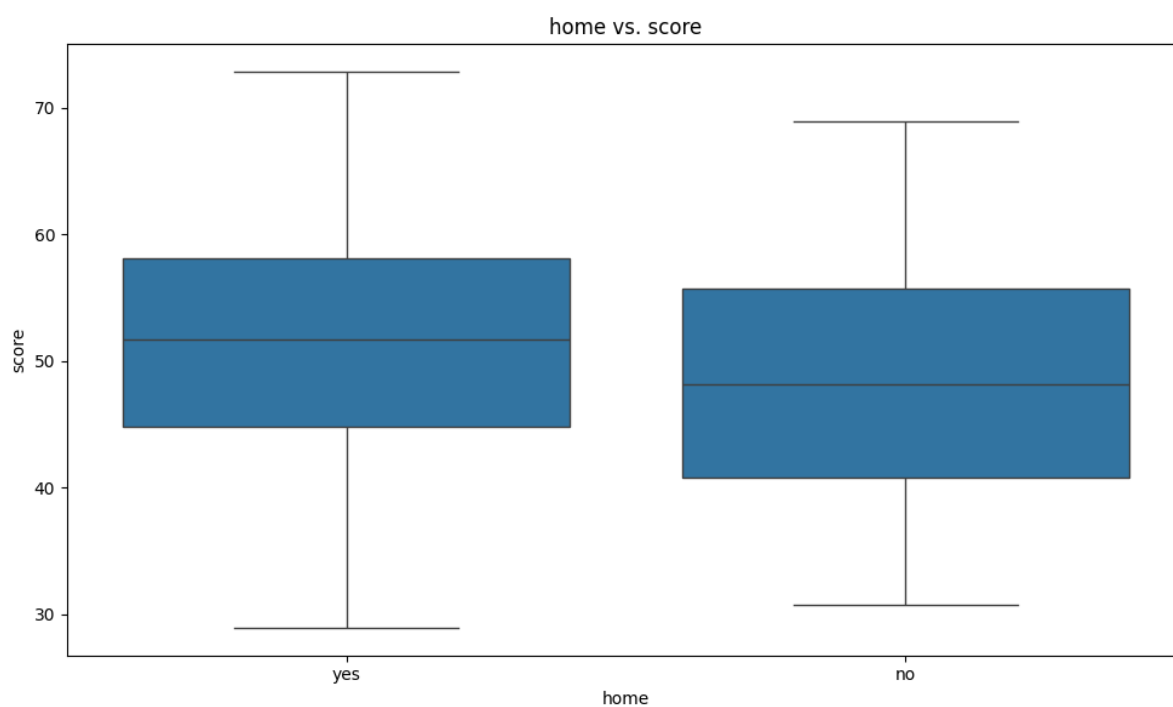
Różnice etniczne mają wpływ na wynik score, Afroamerykanie mają tendencję do osiągania niższych wyników.



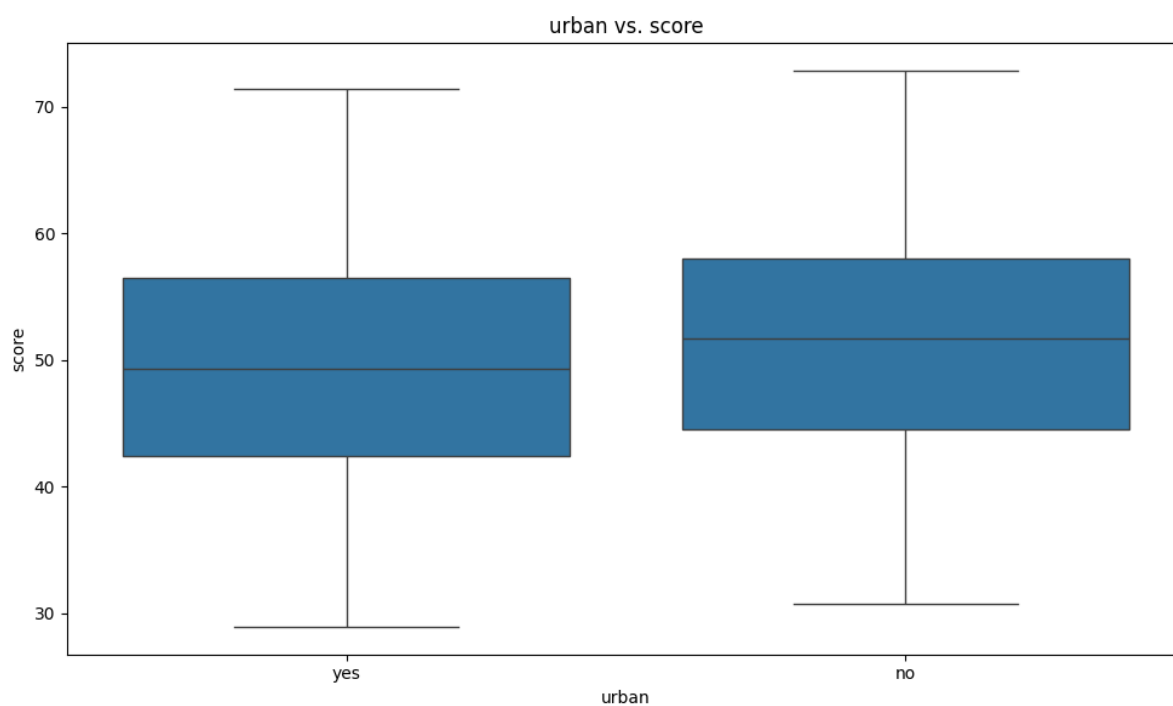
Rodzice z wykształceniem wyższym (zarówno matka, jak i ojciec) mają pozytywną korelację z wyższym wynikiem score.



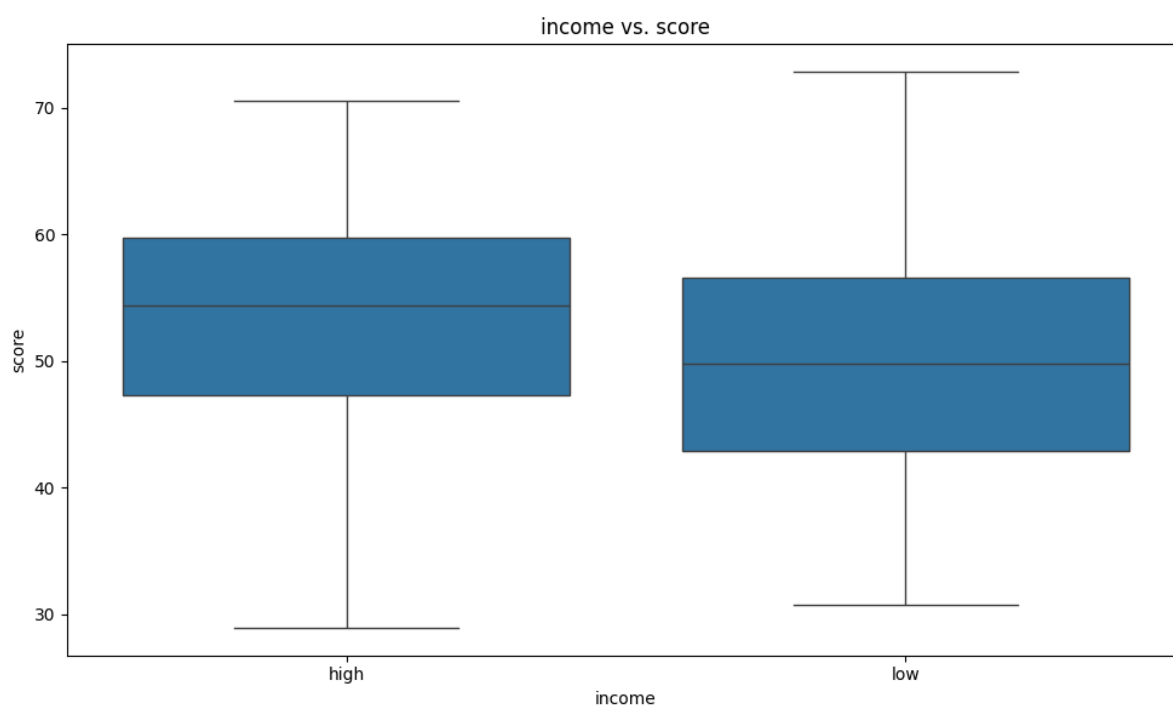
Studenci mieszkający w domu mają tendencję do niższych wyników score.



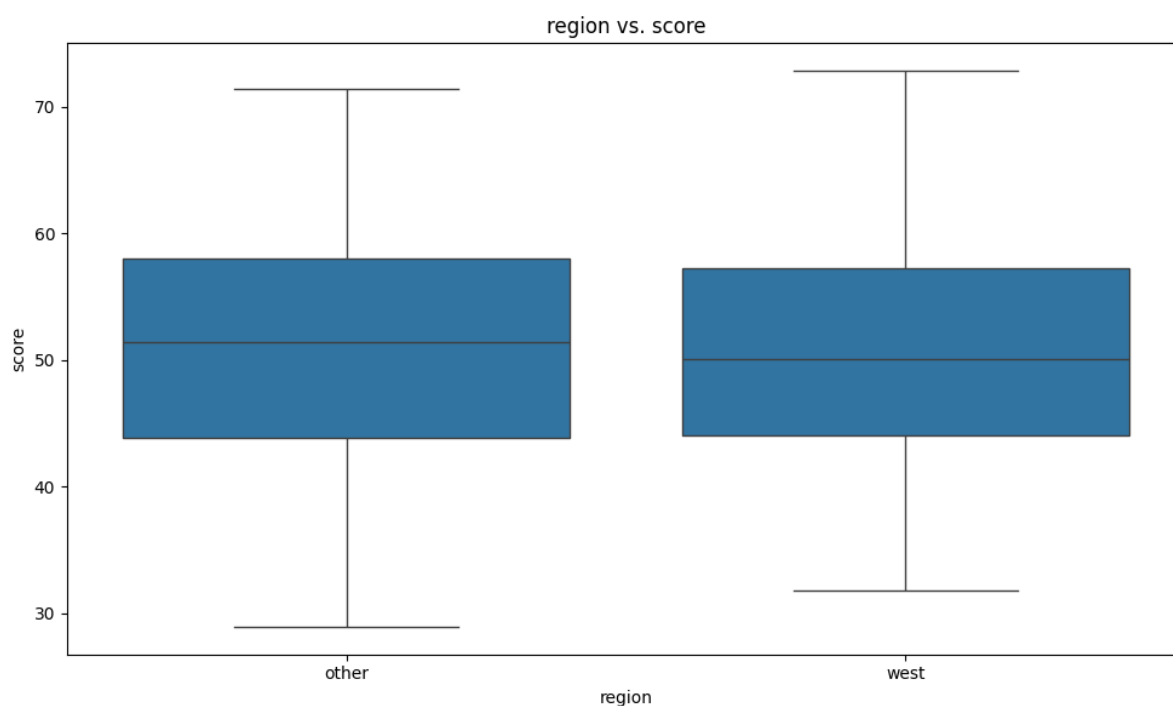
Mieszkanie w obszarze miejskim nieznacznie wpływa na wyższy wynik.



Wyższy dochód koreluje z wyższym wynikiem score.



Region ma wpływ na wyniki score, region "west" pokazuje wyższe wyniki w porównaniu do "other".



Inżynieria cech i przygotowanie danych

Podjęto następujące kroki, aby przygotować dane to pracy z modelem:

Imputacja

- W zmiennych numerycznych brakujące wartości są uzupełniane średnią (SimpleImputer(strategy='mean')).
- W zmiennych kategoriycznych brakujące wartości są uzupełniane najczęściej występującą kategorią (SimpleImputer(strategy='most_frequent')).

Standaryzacja:

- Dla zmiennych numerycznych używana jest standaryzacja za pomocą StandardScaler. Każda zmienna jest przekształcana tak, aby miała średnią 0 i odchylenie standardowe 1.

One-Hot Encoding:

- Dla zmiennych kategoriycznych używana jest technika one-hot encoding, która przekształca kategorie w zmienne binarne.

Podział danych:

- Dane są podzielone na zbiór treningowy i testowy w proporcji 80% na trening i 20% na test.

Kod odpowiedzialny za inżynierie cech i przygotowanie danych – data_prediction.py

```
21 # Wczytanie datasetu
22 df = pd.read_csv('./CollegeDistance.csv')
23
24 # Podział zmiennych
25 num_cols = ['unemp', 'wage', 'distance', 'tuition', 'education']
26 cat_cols = ['gender', 'ethnicity', 'fcollege', 'mcollege', 'home', 'urban', 'income', 'region']
27
28 # Zmienna przewidywana (target)
29 target = 'score'
30
31 # Tworzymy pipeline dla zmiennych numerycznych i kategoriycznych
32 numeric_transformer = Pipeline(steps=[
33     ('imputer', SimpleImputer(strategy='mean')), # Imputacja wartości brakujących
34     ('scaler', StandardScaler()) # Standaryzacja
35 ])
36
37 categorical_transformer = Pipeline(steps=[
38     ('imputer', SimpleImputer(strategy='most_frequent')), # Imputacja najczęściej występujących wartości
39     ('onehot', OneHotEncoder(handle_unknown='ignore')) # One-hot encoding
40 ])
41
42 # Połączenie transformacji dla zmiennych numerycznych i kategoriycznych
43 preprocessor = ColumnTransformer(
44     transformers=[
45         ('num', numeric_transformer, num_cols),
46         ('cat', categorical_transformer, cat_cols)
47     ])
48
49 # Przygotowanie danych do modelu
50 X = df.drop(columns=[target]) # Zbiór cech (bez zmiennej score)
51 y = df[target] # Zmienna przewidywana
52
53 # Podział danych na zbiór treningowy i testowy
54 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

Wybór modelu predykcyjnego

Zadanie polega na przewidywaniu zmiennej ciągłej score, co wskazuje na regresję, dlatego należy wybrać odpowiedni model regresyjny.

Potencjalne modele do rozważenia:

Regresja liniowa:

- Dobrze działa na stosunkowo prostych danych, gdzie zmienność przewidywanej zmiennej dobrze opisuje jedna lub więcej zmiennych wejściowych w liniowy sposób.
- Szybki i wydajny model do treningu, jednak zakłada liniową zależność, co może ograniczać jego skuteczność przy bardziej skomplikowanych zależnościach.

Lasy losowe (Random Forest):

- Model oparty na drzewach decyzyjnych, działający dobrze z bardziej złożonymi relacjami między zmiennymi oraz w przypadku danych, które nie są idealnie liniowe.
- Odporny na nadmierne dopasowanie (overfitting) dzięki technice baggingu.

Gradient Boosting (XGBoost, LightGBM):

- Zaawansowana technika regresji, skuteczna w modelowaniu skomplikowanych zależności. Wymaga większej mocy obliczeniowej, ale może dawać bardzo dobre wyniki, zwłaszcza na danych z nieliniowymi zależnościami.

Regresja LASSO / Ridge:

- Modele liniowe z regularyzacją. Regularyzacja pozwala na kontrolę nadmiernego dopasowania poprzez karanie modelu za zbyt duże współczynniki.

Random Forest

Random Forest będzie najlepszym modelem do rozwiązania tego problemu, ponieważ:

- Jest odporny na zmienne kategoryczne i dobrze radzi sobie z różnorodnością zmiennych.
- Może modelować złożone relacje między cechami, a regulacja liczby drzew i głębokości drzewa może zapobiec nadmiernemu dopasowaniu.
- W przeciwieństwie do regresji liniowej, nie wymaga założenia liniowości danych, co może być korzystne w przypadku tego zestawu danych, który zawiera zarówno zmienne numeryczne, jak i kategoryczne.

Trenowanie modelu

Po wytrenowaniu modelu stosując podstawowe ustawienia, uzyskano następujące wyniki:

- Wyniki dla zbioru treningowego: MSE = 10.4788, MAE = 2.5032, RMSE = 3.2371, $R^2 = 0.8615$, MAPE = 0.0513
- Wyniki dla zbioru testowego: MSE = 53.6624, MAE = 5.8536, RMSE = 7.3255, $R^2 = 0.2924$, MAPE = 0.1204

Widoczna jest duża dysproporcja pomiędzy zachowaniem modelu na danych treningowych, a testowych. Otrzymane wyniki nie są satysfakcjonujące i potrzebna jest optymalizacja modelu.

Wyjaśnienie metryk:

- Mean Squared Error (MSE): ocenia średni błąd kwadratowy pomiędzy przewidywaną a rzeczywistą wartością.

MAE (Mean Absolute Error):

- Średni błąd absolutny – mówi, jak bardzo w średniej model myli się przy przewidywaniu, bez kwadratowania różnic.

RMSE (Root Mean Squared Error):

- Pierwiastek z błędu średniokwadratowego – miara ta przedstawia błąd w jednostkach oryginalnej zmiennej, co może być łatwiej interpretowalne.

R-squared (R^2):

- Ocenia, jak dobrze nasz model dopasowuje się do danych.

MAPE (Mean Absolute Percentage Error):

- Średni błąd procentowy – miara przydatna, gdy chcemy zobaczyć, jak bardzo procentowo przewidywane wartości różnią się od rzeczywistych.

Kod odpowiedzialny za wytrenowanie modelu – data_prediction.py

```
80 #Trenowanie modelu Random Forest
81 rf_model = RandomForestRegressor(n_estimators=100, random_state=42)
82
83 # Trenowanie modelu na zbiorze treningowym
84 rf_model.fit(X_train_prepared, y_train)
85
86 # Przewidywanie na zbiorze treningowym
87 y_train_pred = rf_model.predict(X_train_prepared)
88
89 # Przewidywanie na zbiorze testowym
90 y_test_pred = rf_model.predict(X_test_prepared)
91
92 # Ocena modelu - MSE, MAE, RMSE, R^2, MAPE
93 train_mse = mean_squared_error(y_train, y_train_pred)
94 train_mae = mean_absolute_error(y_train, y_train_pred)
95 train_rmse = np.sqrt(train_mse)
96 train_r2 = r2_score(y_train, y_train_pred)
97 train_mape = mean_absolute_percentage_error(y_train, y_train_pred)
98
99 test_mse = mean_squared_error(y_test, y_test_pred)
100 test_mae = mean_absolute_error(y_test, y_test_pred)
101 test_rmse = np.sqrt(test_mse)
102 test_r2 = r2_score(y_test, y_test_pred)
103 test_mape = mean_absolute_percentage_error(y_test, y_test_pred)
```

Optymalizacja modelu

Zdecydowanie warto przeprowadzić optymalizację modelu, ponieważ dotychczasowe wyniki na zbiorze testowym wskazują na nadmierne dopasowanie (overfitting) i słabe ogólne dopasowanie na nowych danych. Zastosowano następujące metody:

Walidacja krzyżowa (cross-validation)

- Walidacja krzyżowa pomaga uniknąć nadmiernego dopasowania, dzieląc dane na kilka różnych podzbiorów. Na każdym podzbiorze model jest trenowany i oceniany, co daje bardziej wiarygodną ocenę modelu.

Tunowanie hiperparametrów (Grid Search)

- Do tuningu hiperparametrów użyta została metoda GridSearchCV (metoda przeszukiwania wszystkich możliwych kombinacji), aby znaleźć najlepsze parametry dla modelu Random Forest.

Po zastosowaniu optymalizacji modelu otrzymano następujące wyniki:

- Wyniki na zbiorze testowym: MSE = 49.3515, MAE = 5.7161, RMSE = 7.0251, R^2 = 0.3492, MAPE = 0.1180

Powyższe wyniki dalej znacząco różnią się od tych uzyskanych pierwotnie na zbiorze treningowym.

Kod odpowiedzialny za optymalizację modelu – data_prediction.py

```
110 # Rozpoczęcie optymalizacji na nowym modelu podstawowym
111 rf_base = RandomForestRegressor(random_state=42)
112
113 # Parametry do tuningu
114 param_grid = {
115     'n_estimators': [100, 200, 300],
116     'max_depth': [None, 10, 20, 30],
117     'min_samples_split': [2, 5, 10],
118     'min_samples_leaf': [1, 2, 4],
119     'bootstrap': [True, False]
120 }
121
122 # Użycie GridSearchCV do tuningu hiperparametrów
123 grid_search = GridSearchCV(estimator=rf_base, param_grid=param_grid,
124                             cv=5, n_jobs=-1, verbose=2, scoring='neg_mean_squared_error')
125
126 # Dopasowanie modelu na danych treningowych
127 grid_search.fit(X_train_prepared, y_train)
128
129 # Najlepsze parametry
130 best_params = grid_search.best_params_
131
132 # Logowanie najlepszych parametrów
133 logging.info(f"Najlepsze parametry dla modelu Random Forest: {best_params}")
134
135 # Wytrenowanie modelu z najlepszymi parametrami
136 rf_best = grid_search.best_estimator_
137
138 # Ocena modelu na zbiorze testowym
139 y_test_pred_best = rf_best.predict(X_test_prepared)
140
141 # Obliczanie metryk dla zoptymalizowanego modelu
142 test_mse_best = mean_squared_error(y_test, y_test_pred_best)
143 test_mae_best = mean_absolute_error(y_test, y_test_pred_best)
144 test_rmse_best = np.sqrt(test_mse_best)
145 test_r2_best = r2_score(y_test, y_test_pred_best)
146 test_mape_best = mean_absolute_percentage_error(y_test, y_test_pred_best)
```

Możliwe przyczyny słabych wyników modelu

Skuteczność modelu Random Forest:

- Choć Random Forest jest potężnym algorytmem, czasem jego zdolność do generalizacji może być ograniczona, szczególnie w przypadku trudnych problemów regresyjnych. Może to wynikać z nieadekwatności cech do przewidywania zmiennej docelowej (score).

Ograniczona informacyjność cech:

- Cechy w zbiorze danych mogą nie być wystarczająco informatywne lub mogą zawierać zbyt wiele szumu, co utrudnia dokładne przewidywania.

Złożoność zależności:

- Relacje między zmiennymi a celem (score) mogą być nieliniowe lub zbyt skomplikowane, co powoduje, że Random Forest nie jest w stanie w pełni uchwycić tych zależności.

Alternatywne modele predykcyjne

W ramach testów przeprowadzono również testy modelu typu Gradient Boosting, a dokładniej XGBoost. Poniższe wyniki nadal nie były satysfakcjonujące i nie zostały umieszczone w projekcie z tego powodu.

- Zbiór testowy: MSE = 48.5271, MAE = 5.7489, RMSE = 6.9661, $R^2 = 0.3601$, MAPE = 0.1191