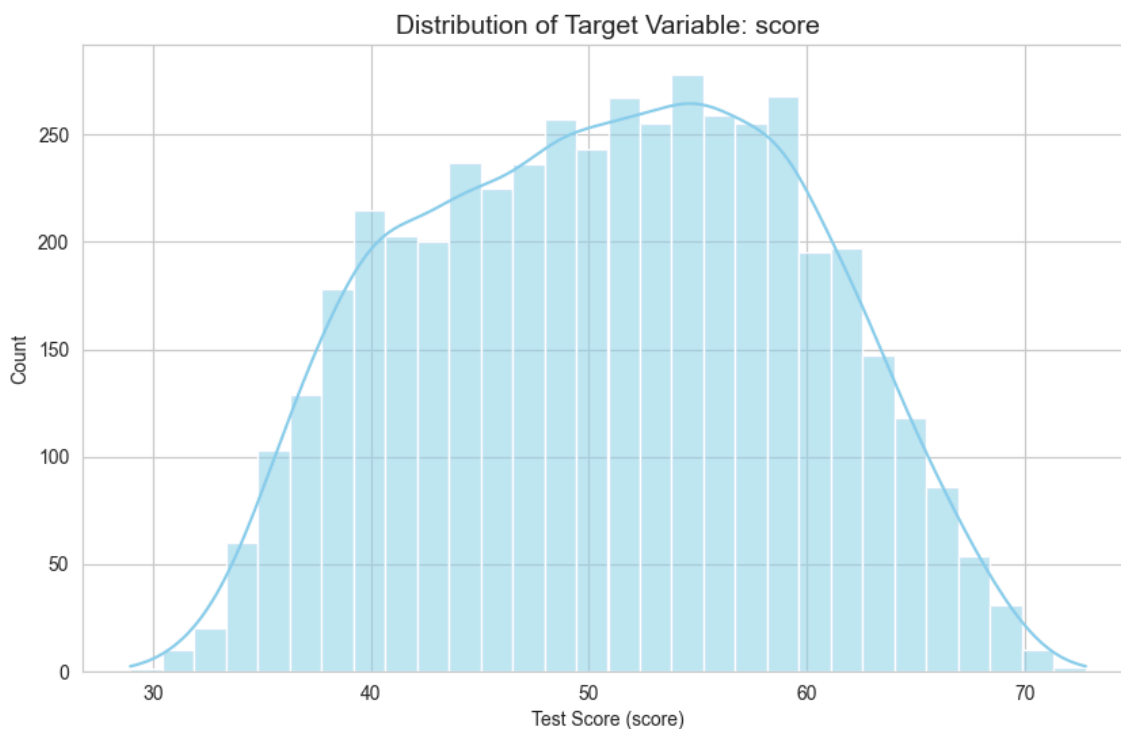


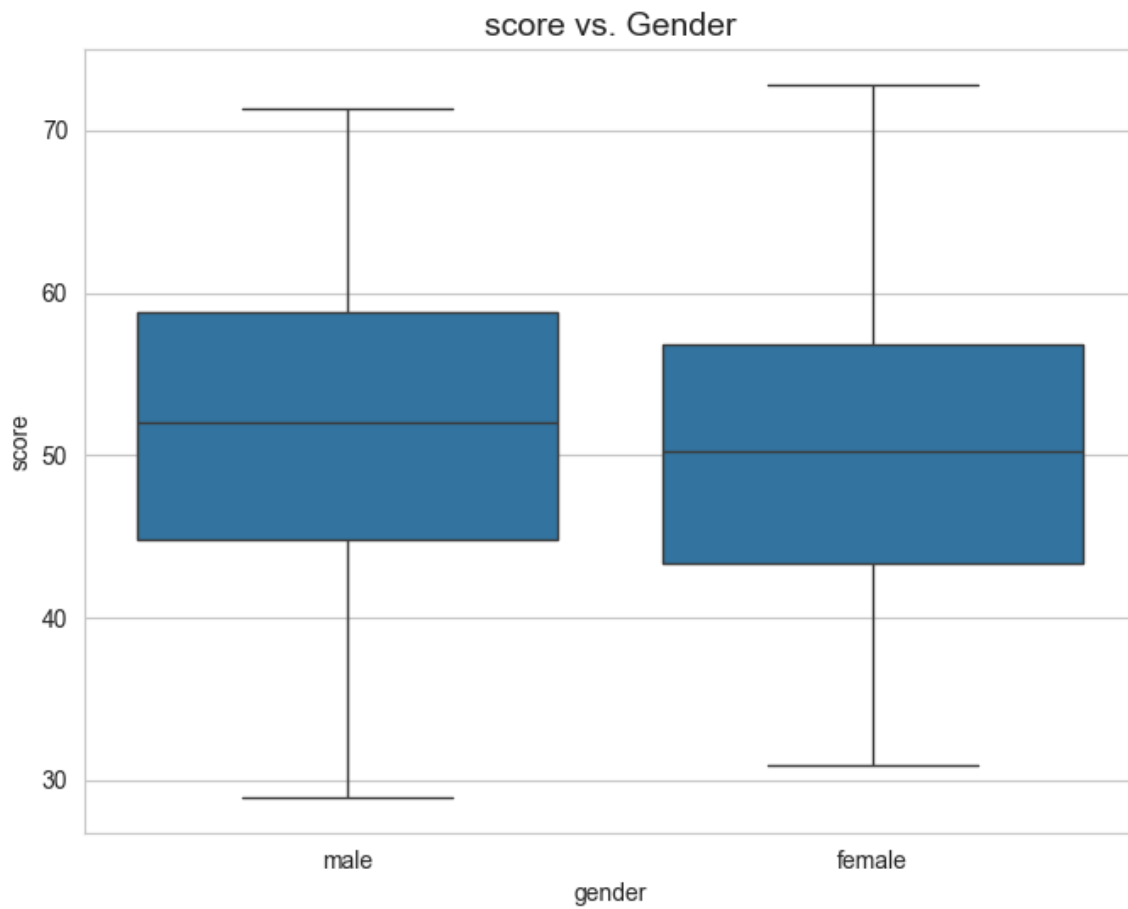
1. Wstępna Analiza Danych (EDA)

Wczytanie i braki: Załadowano zbiór `CollegeDistance`. **Nie zidentyfikowano żadnych brakujących wartości (NaN)**, więc dane były od razu czyste.

Charakterystyka: Najsilniejsza liniowa zależność występuje między `score` a `education`.







2. Przygotowanie Danych i Inżynieria Cech

Standaryzacja: Cechy numeryczne przeskalowano za pomocą `StandardScaler`.

Kodowanie: Cechy kateryczne ('gender', 'ethnicity', 'fcollege', 'mcollege', 'home', 'urban', 'region') zakodowano metodą One-Hot Encoding.

Podział: Dane podzielono na zbiory **Treningowy (80%)** i **Testowy (20%)**.

3. Wybór i Wyniki Modeli

Wytrenowano trzy modele regresji (Regresja Liniowa, Random Forest, Gradient Boosting) w celu wyboru najlepszego algorytmu.

Wyniki R2 (Test Set):

- Regresja Liniowa: 0.3523
- Random Forest: 0.2924
- Gradient Boosting: 0.3659

Wybór: Model **Gradient Boosting** osiągnął najwyższy wynik R^2 i został wybrany jako model końcowy, ponieważ najlepiej dopasowuje się do nieliniowych wzorców w danych.

4. Ocena Końcowa i Wnioski

Najlepszy Model: Gradient Boosting Regressor ($R^2=0.3659$).

Końcowe Metryki (Test Set):

- R^2 : 0.3659
- MAE: 5.7015
- MSE: 48.0831

Wnioski: Uzyskany R^2 jest umiarkowany. Sugeruje to, że brakuje kluczowych czynników predykcyjnych poza dostarczonymi danymi (np. czynników psychologicznych, indywidualnej motywacji).