

# Zaawansowane ćwiczenia z biblioteki `pandas`

## Wprowadzenie

Ćwiczenia te mają na celu rozwinięcie umiejętności w korzystaniu z biblioteki `pandas` do zaawansowanej analizy danych. Pracujemy na dwóch plikach CSV: `astronomical_data_file1.csv` i `astronomical_data_file2.csv`, które zawierają dane obserwacyjne.

## Ćwiczenia

1. **Wczytaj dane i zweryfikuj strukturę.** Wczytaj oba pliki CSV do obiektów `DataFrame` i zweryfikuj ich struktury, wyświetlając podstawowe informacje o kolumnach oraz brakujących wartościach.
2. **Porównaj rozkład brakujących danych.** Przedstaw graficznie rozkład brakujących danych w obu plikach, wykorzystując np. wykres słupkowy.
3. **Połącz dane w jeden `DataFrame`.** Połącz oba zestawy danych, zachowując informacje o pochodzeniu każdej kolumny poprzez dodanie odpowiednich sufiksów.
4. **Zidentyfikuj wspólne i unikalne obiekty.** Znajdź obiekty (na podstawie `ObjectID`), które występują tylko w jednym z plików oraz te, które są wspólne dla obu plików.
5. **Oblicz mediane i odchylenie standardowe.** Oblicz mediane i odchylenie standardowe dla kolumn `Brightness` i `Distance` w obu plikach, zarówno dla obiektów wspólnych, jak i unikalnych.
6. **Uzupełnij brakujące wartości metoda interpolacji.** Uzupełnij brakujące wartości metoda interpolacji liniowej, a następnie porównaj wyniki ze średnimi oryginalnych kolumn.
7. **Znajdź obiekty o najwyższych i najniższych wartościach.** W obu plikach znajdź obiekty o skrajnych wartościach jasności (`Brightness`) i prędkości (`Velocity`) i sprawdź, czy pojawiają się w obu plikach.
8. **Sprawdź korelacje między kolumnami.** Przeprowadź analizę korelacji między kolumnami (np. `Brightness`, `Distance`, `Velocity`) i przedstaw wyniki w formie macierzy korelacji.

9. **Przeanalizuj zmienność jasności i odległości.** Podziel obiekty na grupy według przedziałów jasności i odległości, a następnie oblicz średnia i odchylenie standardowe dla każdej grupy.
10. **Przefiltruj dane na podstawie warunków.** Znajdź wszystkie obiekty, które mają jasność powyżej 15, odległość poniżej 5 oraz predkość powyżej 100. Przedstaw wyniki na wykresie.
11. **Zidentyfikuj trend w danych.** Utwórz wykres pokazujący, jak jasność zmienia się wraz z odległością dla obiektów wspólnych dla obu plików. Dodaj linie trendu.
12. **Oblicz różnice dla wspólnych obiektów.** Dla obiektów obecnych w obu plikach, oblicz różnice w jasności i predkości między plikami. Zidentyfikuj obiekty o największych różnicach.
13. **Grupowanie danych według kolumny.** Grupuj obiekty według ich jasności i oblicz statystyki dla każdego przedziału jasności, m.in. średnia odległość i średnia predkość.
14. **Porównaj zmienne między plikami.** Zrób porównanie rozkładów jasności, odległości i predkości między plikami za pomocą wykresu pudełkowego (boxplot).
15. **Analiza wielkości brakujących danych.** Zidentyfikuj kolumny o największej liczbie brakujących wartości w każdym pliku. Przedstaw wykres z rozkładem braków.
16. **Wyznacz wartości odstające.** Znajdź wartości odstające w kolumnach `Brightness` i `Distance` na podstawie zdefiniowanych progów. Wyświetl je i wskaż, czy występują w obu plikach.
17. **Wypełnij wartości brakujące średnia lub mediana.** Uzupełnij brakujące wartości w kolumnach `Brightness` i `Distance` na podstawie ich średniej lub mediany. Porównaj wyniki obu metod.
18. **Eksport przetworzonych danych.** Po wykonaniu wszystkich powyższych analiz, zapisz przetworzone dane w trzech plikach: wszystkie dane, dane wspólne i dane unikalne.
19. **Tworzenie funkcji analizy danych.** Stwórz funkcje, która automatycznie analizuje dane wejściowe i wykonuje operacje takie jak wypełnianie braków, wyznaczanie statystyk i znajdowanie wspólnych obiektów.
20. **Porównanie za pomocą wykresów.** Użyj wykresów histogramów, aby porównać rozkład wartości `Brightness` oraz `Velocity` między obiektami wspólnymi i unikalnymi w obu plikach.

## Podpowiedzi

Do wykonania powyższych zadań pomocne mogą być następujące funkcje biblioteki `pandas`:

- `pd.read_csv()` - wczytywanie danych z plików CSV,
- `df.info()` - uzyskanie informacji o strukturze danych,
- `df.isna()` i `df.dropna()` - sprawdzanie i usuwanie brakujących wartości,
- `df.fillna()` - wypełnianie brakujących wartości,
- `pd.merge()` - łączenie zbiorów danych na podstawie klucza,
- `df.groupby()` - grupowanie danych,
- `df.corr()` - obliczanie korelacji między kolumnami,
- `df.describe()` - wyznaczanie podstawowych statystyk,
- `df.median()` i `df.mean()` - obliczanie mediany i średniej,
- `df.plot()` i `df.hist()` - tworzenie wykresów,
- `df.duplicated()` - znajdowanie duplikatów,
- `df.apply()` - stosowanie funkcji do kolumn,
- `df.to_csv()` - eksport przetworzonych danych do plików CSV.