

Projekt zaliczeniowy



Zaawansowane programowanie w języku Python Rok akademicki 2024/2025

Autorzy:

Mateusz Strojek

[Analiza najlepszych klubów piłkarskich – projekt funkcyjny]

Projekt skupia się na wstępnej analizie danych dotyczących klubów piłkarskich oraz na analizie skupień. Dane są scrapowane ze strony <https://www.whoscored.com/Statistics>. Pobierane są poszczególne zmienne: drużyna, liga, liczba goli, strzały na mecz, żółte kartki, średnie posiadanie piłki na mecz, średni odsetek dokładnych podań, liczba wygranych powietrznych pojedynków na mecz oraz średnia ocena za mecz. Strona zawiera 4 zakładki; na pierwszej pokazane są podstawowe informacje o ramce danych, na drugiej pokazany jest wykres pudełkowy, wykres punktowy oraz wykres typu 'countplot' wraz z korelogramem. Na trzeciej zakładce pokazana jest analiza k-średnich zawierająca wartości indeksu sylwetkowego w zależności od liczby skupień, wykres punktowy w zależności od poszczególnych zmiennych oraz statystyki opisowe każdej grupy. Na ostatniej zakładce zostało przeprowadzone porządkowanie liniowe metodą Hellwiga.

Podział ról

Mateusz Strojek: Całość.

Opis funkcjonalności

Projekt składa się z pięciu plików:

Strona.py – plik zawierający funkcje czyszczące ramkę danych, układ strony, takie jak:

1. `transform_variable_to_categorical(column, data)` - Zmienia typ danej kolumny na kategorie.
2. `change_premier_league_to_pl(column, data)` - Zmienia nazwę "Premier League" na "PL" w kolumnie "Tournament".
3. `plot_scatter_plot(df, x_axis, y_axis, hue=False)` - Tworzy wykres punktowy między dwiema wybranymi zmiennymi.
4. `show_dtypes(data)` - Wyświetla typy danych w ramce danych.
5. `show_stats(data)` - Wyświetla podstawowe statystyki opisowe dla danych.

6. `silhouette_print(data)` - Oblicza i wyświetla wartość indeksu sylwetkowego dla różnych liczb klastrow.
7. `plot_correlation_matrix(df)` - Tworzy macierz korelacji dla danych numerycznych i wyświetla ją w formie korelogramu.
8. `plot_boxplot(variable_box, df, hue=True)` - Tworzy wykres pudełkowy dla wskazanej zmiennej, z możliwością dodania podziału według ligi.
9. `plot_countplot(df)` - Tworzy wykres słupkowy przedstawiający liczbę drużyn w różnych ligach.
10. `plot_kmeans_clusters(data, num_clusters, x_axis, y_axis)` - Przeprowadza analizę K-means na danych i rysuje wykres dla dwóch zmiennych.
11. `generate_h1(text)` - Generuje nagłówek H1 w aplikacji Streamlit.
12. `generate_h2(text)` - Generuje nagłówek H2 w aplikacji Streamlit.
13. `generate_h3_centered(text)` - Generuje nagłówek H3 w aplikacji Streamlit.
14. `generate_head_table(df:pd.DataFrame, slides=0)` - Wyświetla pierwsze wiersze tabeli z danymi.
15. `clean_data(df_raw)` - Czyści dane, zmieniając wartości w kolumnie "Tournament" i usuwając duplikaty.
16. `get_data()` - Ładuje dane, z możliwością skorzystania z istniejącego pliku CSV lub skrapowania nowych danych.
17. `run_tab1(df)` - wygenerowanie pierwszej zakładki „Overview”
18. `run_tab2(df)` - wygenerowanie drugiej zakładki „Plots”
19. `run_tab3(df)` - wygenerowanie trzeciej zakładki „K-means analysis”
20. `run_tab4(df)` - wygenerowanie czwartej zakładki „Linear ordering”
21. `change_variables(variable_types, df)` - zamiana wartości zmiennych w zależności czy zmienna jest stymulantą czy destymulantą
22. `stimulant(column)` - funkcja do stymulanty
23. `destimulant(column)` - funkcja do destymulanty
24. `distances(dataframe, maxima, weights=None)` - Funkcja licząca dystans pomiędzy wzorcem a aktualną wartością
25. `di_count(distance_frame)` - Funkcja licząca mi wartość d1 do wyniku indeksu
26. `linear_ordering_Hellwig(df_original, df, weights, variable_types)` - Funkcja do porządkowania liniowego
27. `run_app()` - generowanie architektury strony z wykresami, itd.
28. `main()` - Główna funkcja aplikacji Streamlit, która generuje interaktywne wykresy i statystyki.

ScrapingData.py – plik zawierający funkcje odpowiedzialne za scrapowanie danych, takie jak:

1. `setup_driver(url)` - Inicjalizuje przeglądarkę Chrome i otwiera stronę internetową podaną w argumencie url.
2. `extract_all_data(driver)` - Wyciąga dane ze wszystkich kolumn tabeli oraz z kolumny 'Discipline'.
3. `extract_cards(discipline)` - Wyciąga informacje o żółtych i czerwonych kartkach z HTML w kolumnie 'Discipline'.
4. `clean_data(dataframe)` - Czyści dane, usuwając zbędne kolumny oraz numerację drużyn.
5. `scrape_statistics(file_name)` - Główna funkcja do scrapowania danych statystycznych z zewnętrznej strony i zapisania ich do pliku CSV.

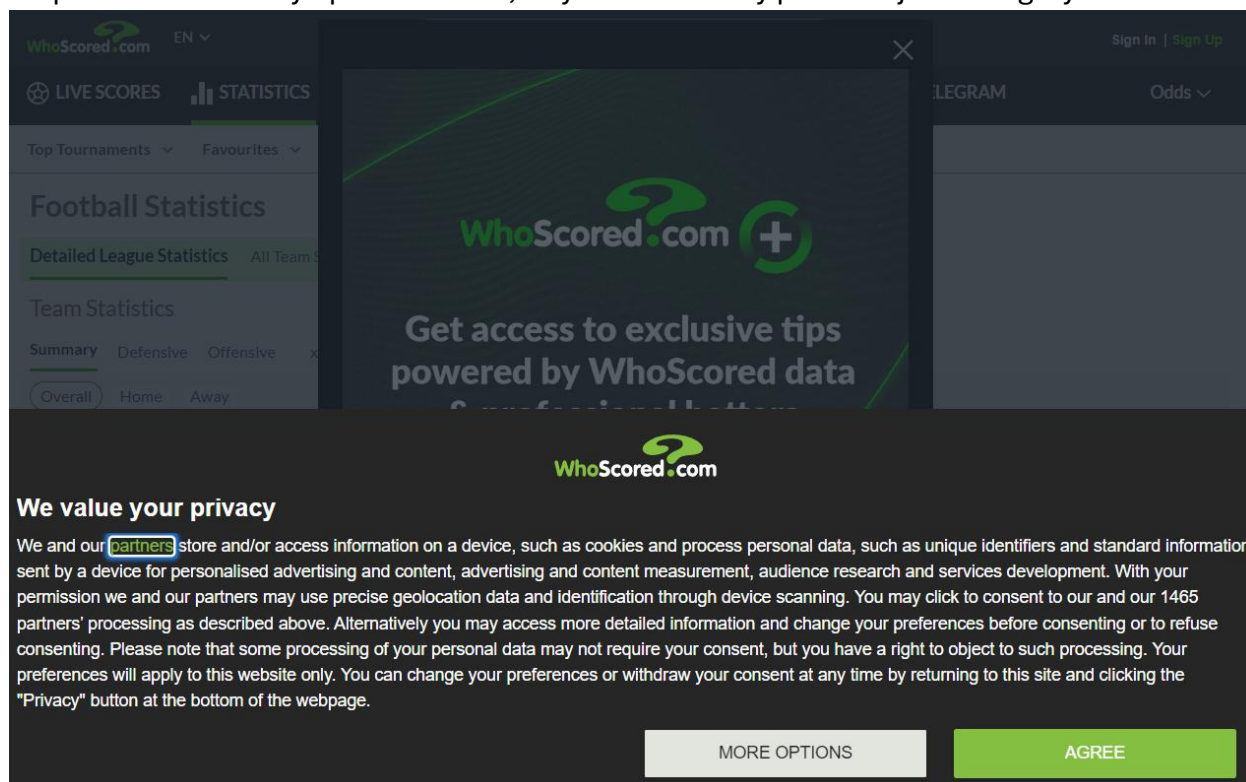
Style.py – plik zawierający style css.

Stare_kluby.csv – plik zawierające „stare dane” – plik jest potrzebny jedynie do demonstracji aktualizacji danych po zescrapowaniu.

Clubs.csv – plik z danymi zescrapowanymi ze strony.

Scrapowanie danych – ważne informacje

- Przed scrapowaniem należy szybko zaakceptować warunki strony i szybko usunąć niepotrzebne reklamy sprzed widoku, aby nie zastaniały pierwszej tabeli z góry.



Następnie zescrolować stronę do poziomu, aż przycisk „next” będzie widoczny, najlepiej mniej więcej do tego momentu:

Chrome is being controlled by automated test software.

6. Barcelona	LaLiga	51	16.8	31	2	67.2	87.9	9.2	6.83
7. Arsenal	Premier League	39	13.8	41	3	55.4	86.5	13.7	6.79
8. Eintracht Frankfurt	Bundesliga	36	15.3	26	1	48.9	83.6	14.1	6.78
9. Bayer Leverkusen	Bundesliga	40	16.4	27	0	60.5	88.1	12.8	6.77
10. Napoli	Serie A	30	14	23	0	52.5	86.4	13.7	6.77
11. Atalanta	Serie A	43	14.9	36	0	56.4	85.7	16	6.77
12. Nice	Ligue 1	35	13.9	33	2	47.3	83.9	11.6	6.76
13. Monaco	Ligue 1	28	14.4	29	3	56.2	82.7	13.1	6.75
14. Chelsea	Premier League	39	15.8	58	1	58.3	87.0	9.8	6.74
15. Mainz 05	Bundesliga	30	12.4	32	3	48.1	77.6	18.1	6.74
16. Newcastle	Premier League	34	14.1	43	1	50.9	83.6	12.5	6.73
17. Tottenham	Premier League	42	14.9	40	1	57.8	85.7	11	6.73
18. Atletico Madrid	LaLiga	33	11.9	39	1	51.6	83.9	12.9	6.73
19. Manchester City	Premier League	36	17.1	40	1	61.4	90.3	8	6.73
20. Nottingham Forest	Premier League	29	12.9	45	2	39.5	78.5	15.1	6.72

© WhoScored

first prev next last

* Only teams from English Premier League, French Ligue 1, German Bundesliga, Italian Serie A and Spanish La Liga are displayed

Privacy

Gdy nie będzie chciało przejść do następnej strony tabelki, należy poruszyć zescrolować lekko w górę/dół.

Widok strony i jej funkcjonalność

Na pierwszej zakładce, jak wspomniałem wcześniej widnieją informacje o podstawowych informacjach ramki danych:

Choose data source:

☐ Use existing data

☒ Scrape new data

Football Clubs Analysis

Overview

Plots

K-means

Overview

Data types

Statistics

	Team	Tournament	Goals	Shots pg	Possession%	Pass%	AerialsWo
Dtypes	object	category	int64	float64	float64	float64	float64

	Goals	Shots pg	Possession%	Pass%	AerialsWon	Rating
count	96.000000	96.000000	96.000000	96.000000	96.000000	96.000000
mean	25.875000	12.566667	49.858333	82.362500	13.663542	6.603438
std	9.149346	2.269786	7.150512	4.198402	2.797635	0.136527
min	11.000000	8.100000	38.600000	70.400000	7.900000	6.320000
25%	19.000000	10.900000	43.675000	79.375000	11.775000	6.510000
50%	24.000000	12.400000	48.950000	82.400000	13.600000	6.590000

The whole dataframe

Number of rows to display

5

96

	Team	Tournament	Goals	Shots pg	Possession%	Pass%	AerialsWon	Rating	Yellow Cards	Red Cards
0	Paris Saint-Germain	Ligue 1	44	19.0	69.6	91.1	8.0	7.02	20	0
1	Bayern Munich	Bundesliga	48	19.6	72.0	90.2	12.0	6.99	24	0
2	Liverpool	PL	47	16.6	56.9	86.4	10.9	6.89	41	1
3	Real Madrid	LaLiga	43	15.2	60.6	89.5	8.4	6.87	32	2
4	Inter	Serie A	45	15.9	60.3	88.9	15.3	6.84	26	0

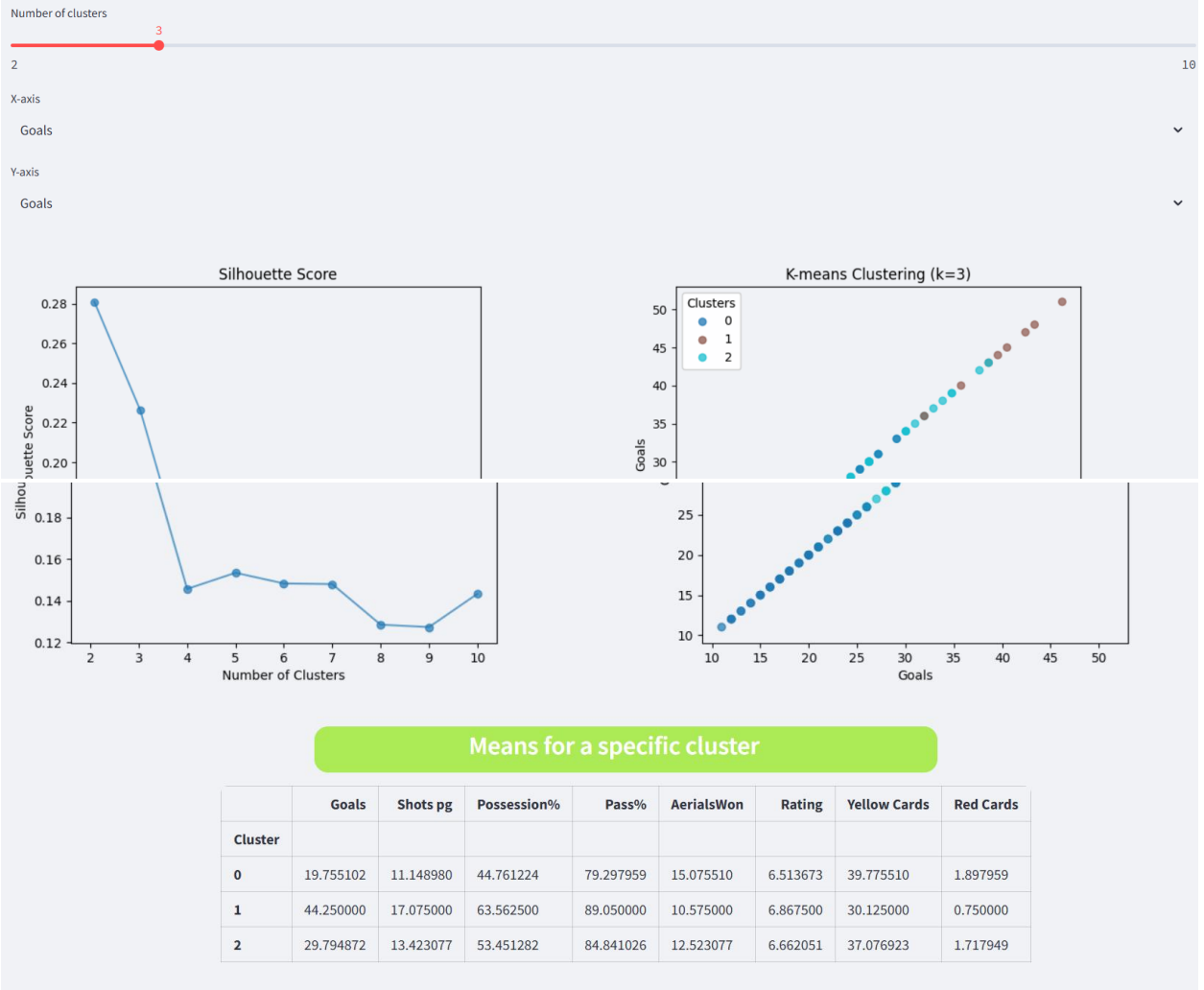
Na tej zakładce użytkownik może rozszerzyć liczbę wierszów do pokazania. Na drugiej zakładce natomiast jest już o wiele więcej funkcjonalności. Wygląda ona tak:



Użytkownik może wybrać jeden spośród trzech wykresów (countplot, boxplot i scatterplot), wybrać zmienne dla poszczególnych osi. Jest również możliwość na np. narysowanie wykresów pudełkowych z podziałem na poszczególłą ligę.

Trzecia zakładka przedstawia analizę skupień, użytkownik może wybrać liczbę skupień na podstawie np. wykresu sylwetkowego oraz zobaczyć, jak kształtują się grupy na wykresie punktowym. Dodatkowo, na samym spodzie jest tabelka ze średnimi dla poszczególnej grupy.

K-means analysis



Na ostatniej zakładce występuje porządkowanie liniowe metodą Hellwiga. Można wybrać, czy zmienna jest stymulantą, czy też nie oraz wybrać wagę danej zmiennej:

Linear Ordering

Goals			Shots per Game			Possession %			Pass %			Aerials Won			Yellow Cards			Red Cards		
Weight			Weight			Weight			Weight			Weight			Weight			Weight		
0.15 - +			0.15 - +			0.15 - +			0.15 - +			0.15 - +			0.15 - +			0.15 - +		
Type			Type			Type			Type			Type			Type			Type		
Stimulant ▾			Stimulant ▾			Stimulant ▾			Stimulant ▾			Stimulant ▾			Destimulant ▾			Destimulant ▾		

Result

	Team	Tournament	Goals	Shots pg	Possession%	Pass%	AerialsWon	Rating	Yellow Cards	Red Cards	HellwigIndex
0	Inter	Serie A	51	15.9	60.1	88.4	15.7	6.82	27	0	0.631946
1	Bayern Munich	Bundesliga	56	20.2	71.5	90.3	12.0	7.02	27	0	0.605296
2	Bayer Leverkusen	Bundesliga	44	16.3	60.5	88.2	12.6	6.79	29	0	0.536761
3	Atalanta	Serie A	46	15.0	56.6	86.1	16.0	6.75	41	0	0.529956
4	Barcelona	LaLiga	52	17.1	67.7	87.8	9.8	6.83	33	2	0.467123
5	VfB Stuttgart	Bundesliga	36	14.7	59.1	85.5	13.8	6.67	35	2	0.460119
6	Liverpool	PL	50	17.9	57.6	86.4	11.3	6.89	45	1	0.448365

Możliwa rozbudowa projektu:

Informatyk może poszerzyć statystyczną analizę danych o np. metodę k-medoid, przeprowadzić analizę wrażliwości obydwu metod, przeprowadzić analizę głównych składowych (PCA), jak i analizę czynnikową (FA). Ponadto, można sprowadzić informację o piłkarzach i zrobić analizę w tym obszarze.