

REGRESJA

Projekt praktyczny





Rozpoznanie danych

Epicurious - przepisy z oceną i wartościami odżywczymi:

Oceny z zakresu 0 - 5

Wartości odżywcze:

- Kalorie,
- Tłuszcze,
- Białko,
- Sód

674 kolumn z tagami i ponad 22,000 rekordów w pliku CSV



Problematyka

Analizowane problemy:

- Wpływa wartości odżywczych na ocenę
- Wpływ występowania mięsa w daniu na ocenę
- Wpływ ilości tagów przepisu na ocenę

Pozostałe hipotezy:

- Wpływ typu posiłku (śniadanie, obiad, kolacja) na ocenę
- Najczęściej występujące słowa w nazwach dań



Zapoznanie się z bazą danych

Wstępna obróbka danych

1. Usunięcie wierszy z wartościami NaN
2. Usunięcie wierszy z oceną 0
3. Czyszczenie outlierów (quantiles)

Dostosowanie danych do analizy:

1. Stworzenie kolumny z informacją o zawartości mięsa w daniu
2. Stworzenie kolumny z liczbą tagów dla każdego dania
3. Usunięcie zbędnych kolumn



Ogólne statystyki przygotowanych danych

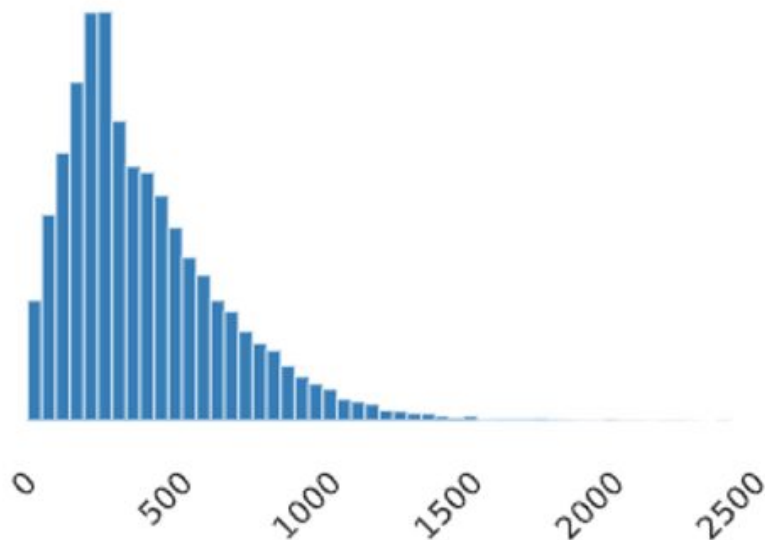
	rating	calories	protein	fat	sodium	with_meat	breakfast	lunch	dinner	num_tags
count	8858.000000	8858.000000	8858.000000	8858.000000	8858.000000	8858.000000	8858.000000	8858.000000	8858.000000	8858.000000
mean	4.081409	403.966358	15.194965	22.653534	438.542335	0.235042	0.045496	0.074622	0.119553	12.899639
std	0.628286	198.702606	12.322600	15.018741	363.788044	0.424049	0.208400	0.262795	0.324456	4.741075
min	1.250000	80.000000	3.000000	1.000000	15.000000	0.000000	0.000000	0.000000	0.000000	1.000000
25%	3.750000	250.000000	6.000000	12.000000	150.000000	0.000000	0.000000	0.000000	0.000000	9.000000
50%	4.375000	364.000000	10.000000	19.000000	332.000000	0.000000	0.000000	0.000000	0.000000	13.000000
75%	4.375000	529.000000	23.000000	30.000000	641.000000	0.000000	0.000000	0.000000	0.000000	17.000000
max	5.000000	985.000000	49.000000	85.000000	1592.000000	1.000000	1.000000	1.000000	1.000000	33.000000



Co można wynieść ze statystyk?

Według National Health Service (NHS) w Wielkiej Brytanii, przeciętny dorosły mężczyzna potrzebuje około 2500 kalorii dziennie, aby utrzymać swoją masę ciała na stałym poziomie, podczas gdy dorosłe kobiety potrzebują średnio około 2000.

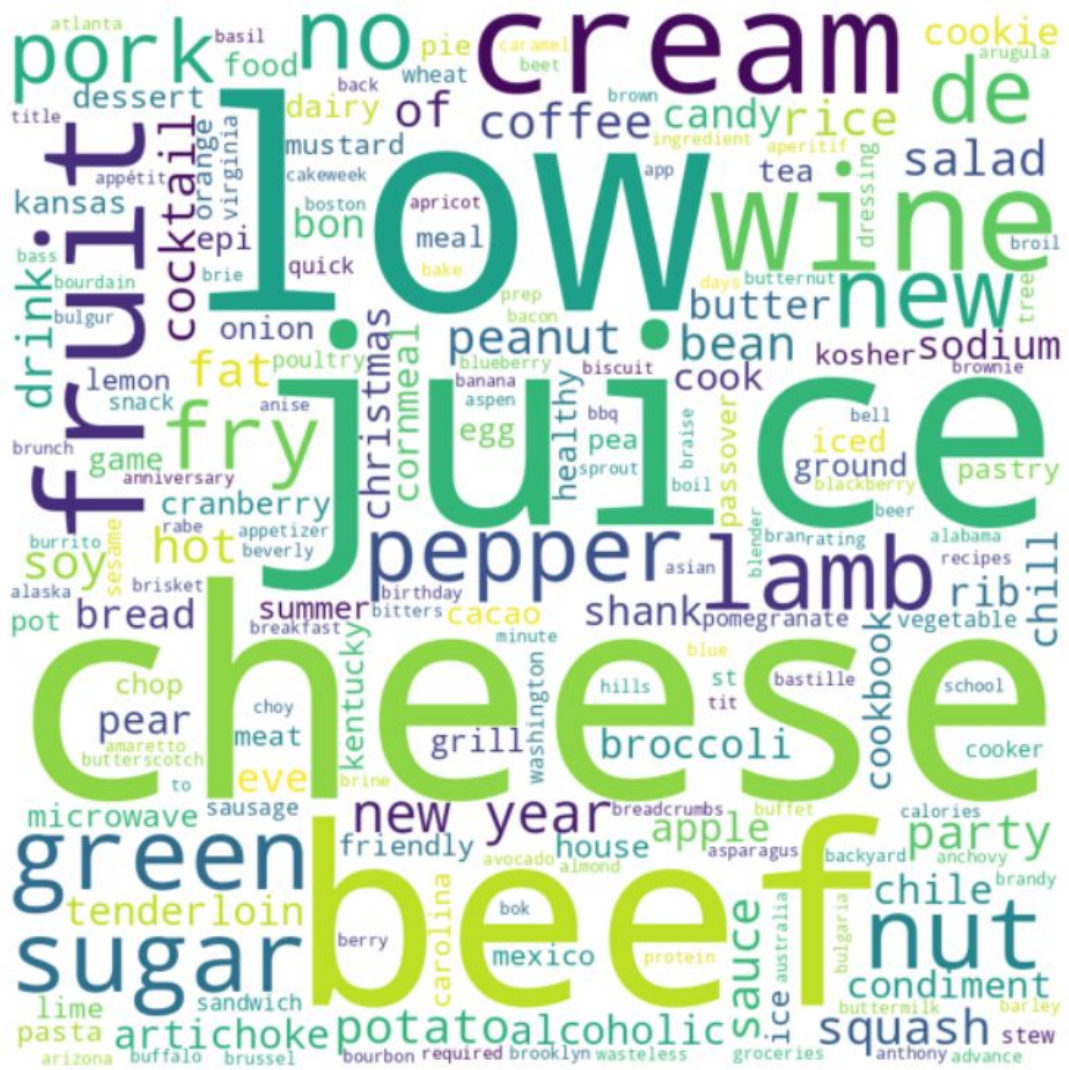
Średnie zawartość kalorii w daniach w zbiorze wyniosła 404kcal. Biorąc pod uwagę 5 posiłków dziennie bilans kaloryczny zostaje spełniony



Rozkład kalorii

Mapa słów

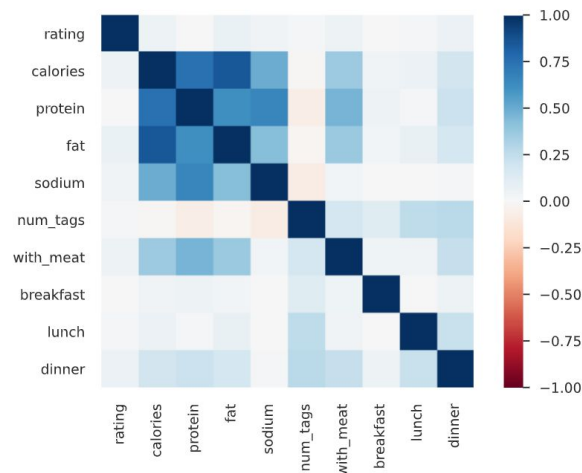
Mapa przedstawia najczęściej występujące słowa w nazwach dań w tabeli

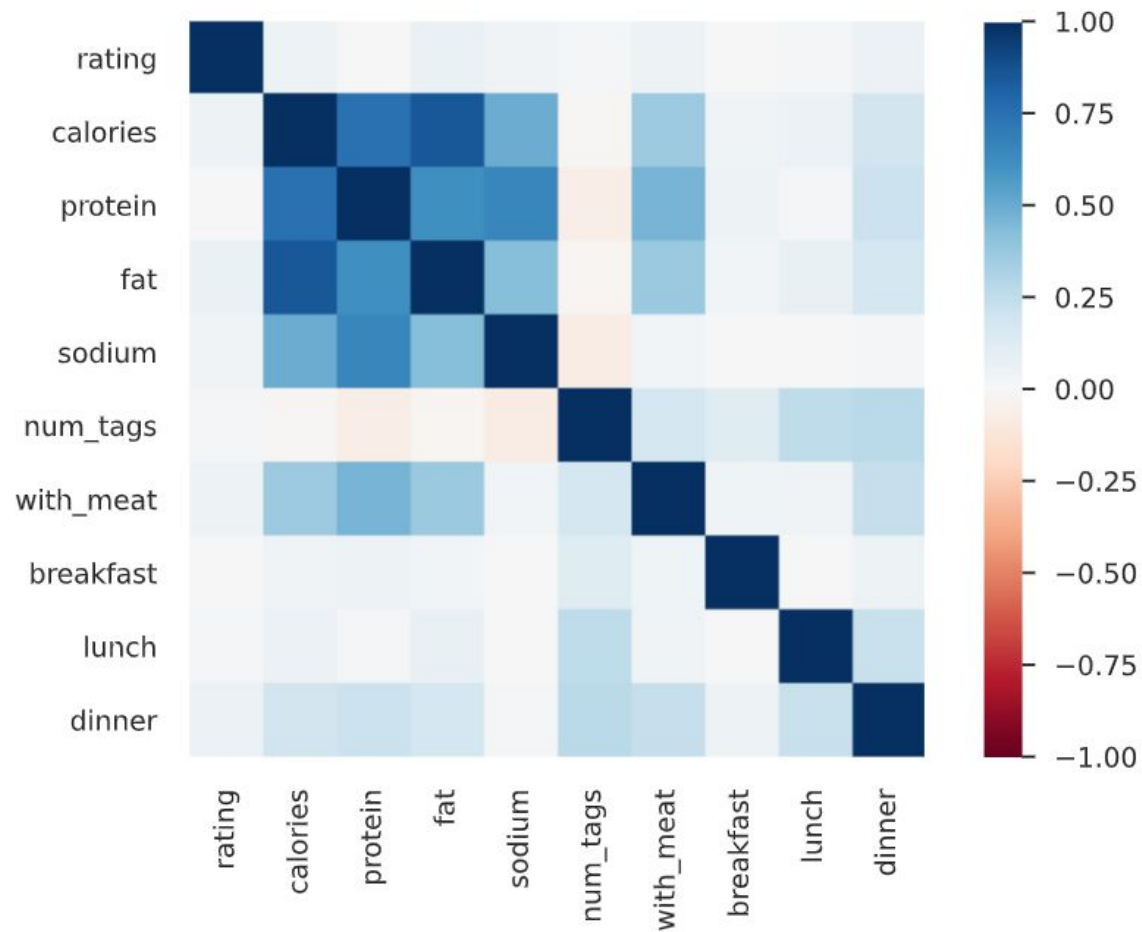


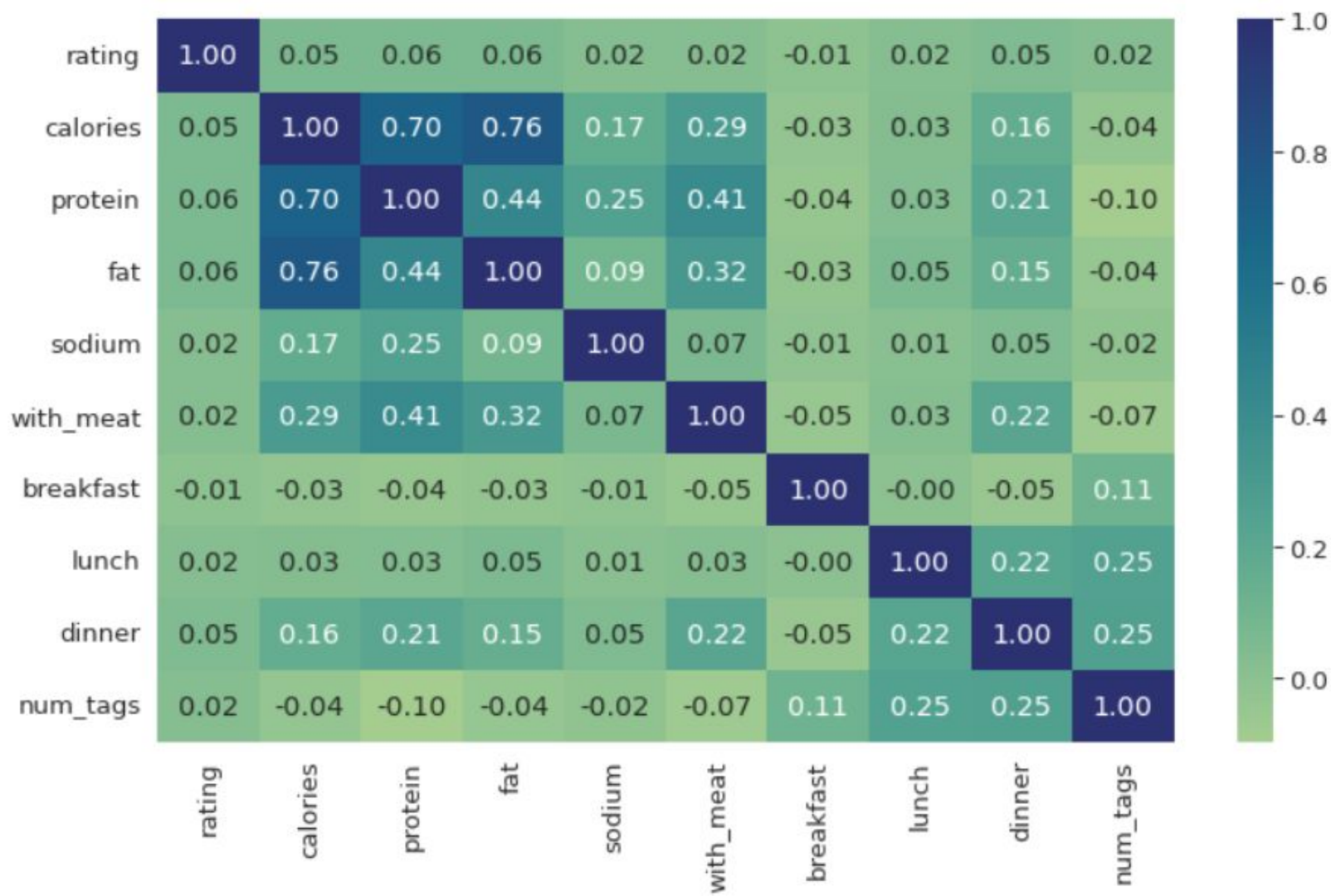


Heatmapa

Heatmapa to narzędzie do analizy danych, które w graficzny sposób (najczęściej wykorzystując kolory) reprezentuje to, jak zmienne są ze sobą skorelowane.









Korelacje zmiennych

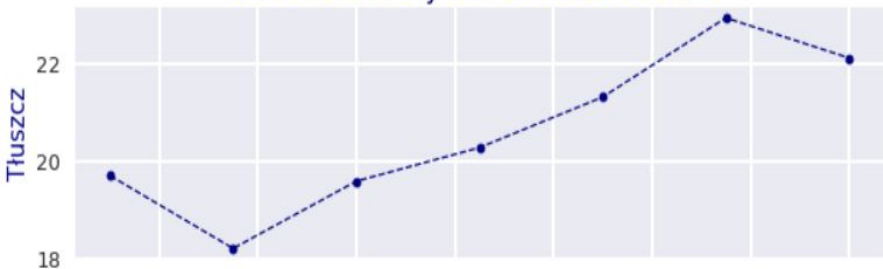
<code>title</code> has a high cardinality: 11640 distinct values	High cardinality
<code>calories</code> is highly overall correlated with <code>protein</code> and 1 other fields	High correlation
<code>protein</code> is highly overall correlated with <code>calories</code> and 2 other fields	High correlation
<code>fat</code> is highly overall correlated with <code>calories</code> and 1 other fields	High correlation
<code>sodium</code> is highly overall correlated with <code>protein</code>	High correlation
<code>breakfast</code> is highly imbalanced (77.4%)	Imbalance
<code>lunch</code> is highly imbalanced (65.4%)	Imbalance
<code>sodium</code> is highly skewed ($\gamma_1 = 30.2760962$)	Skewed
<code>title</code> is uniformly distributed	Uniform
<code>protein</code> has 222 (1.7%) zeros	Zeros



Wykresy liniowe zależności oceny

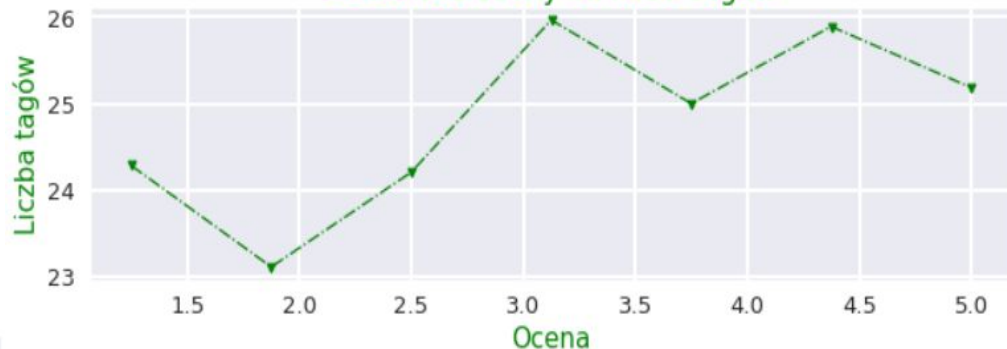
Przy ilości 22g tłuszczów najbardziej prawdopodobne jest uzyskanie maksymalnej oceny

Zależność oceny od zawartości tłuszczu

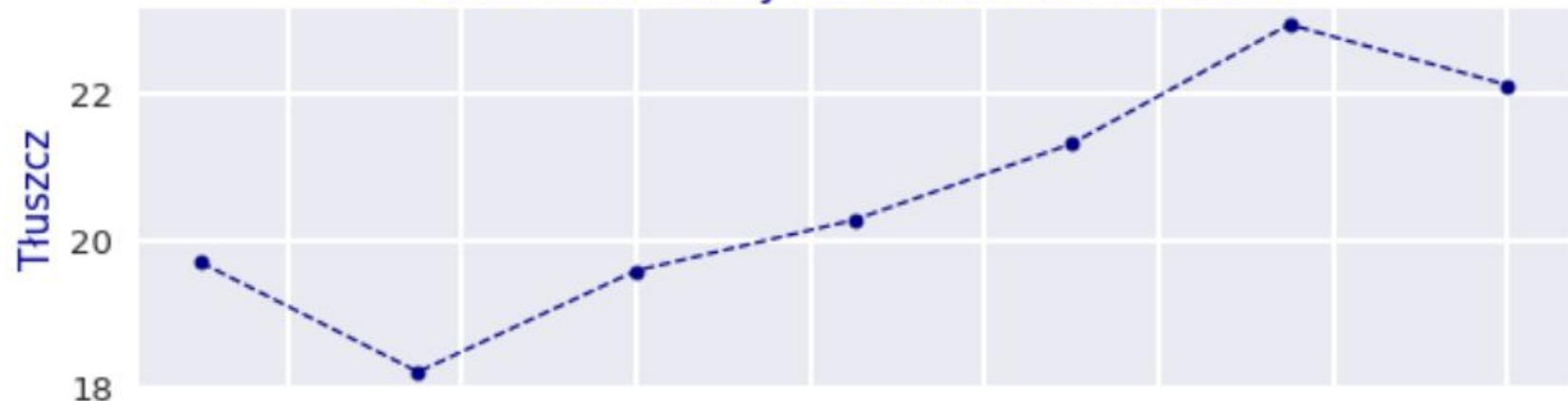


W celu uzyskania maksymalnej oceny optymalną wartością tagów jest 25

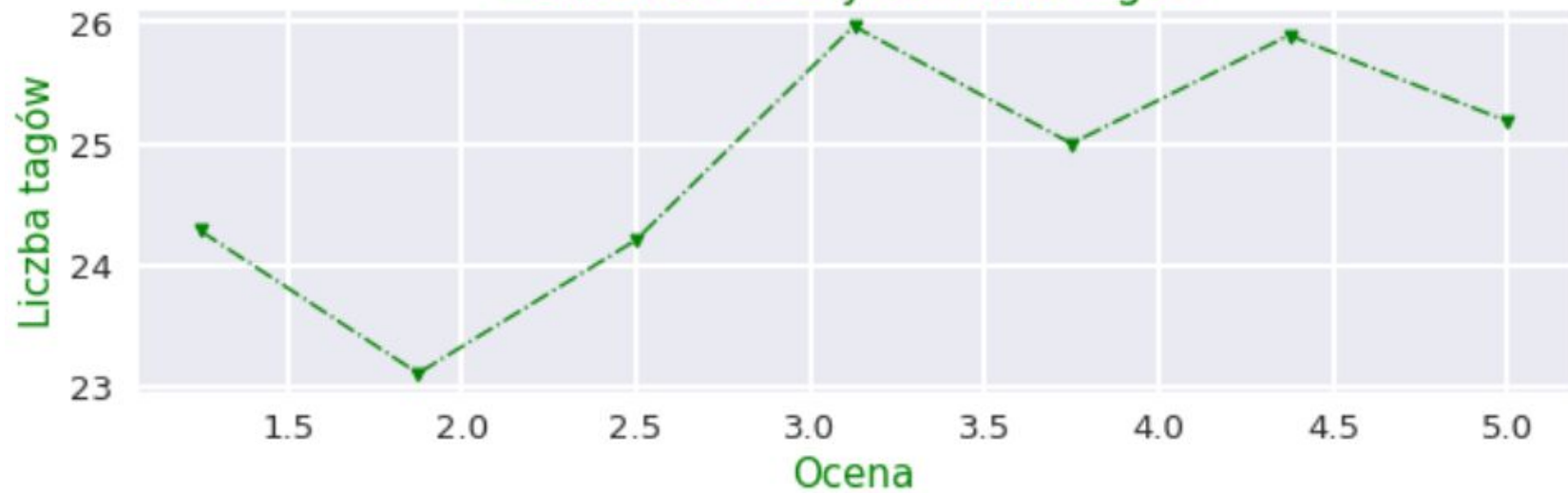
Zależność oceny od ilości tagów



Zależność oceny od zawartości tłuszczu



Zależność oceny od ilości tagów





Jak rozkładają się oceny dań?

Distinct	7
Distinct (%)	0.1%
Missing	0
Missing (%)	0.0%
Infinite	0
Infinite (%)	0.0%
Mean	4.0838572

Minimum	1.25
Maximum	5
Zeros	0
Zeros (%)	0.0%
Negative	0
Negative (%)	0.0%
Memory size	202.6 KiB



Regresja

- **Zmienne niezależne:**

"calories", "protein", "fat", "sodium", "with_meat", "breakfast", "lunch", "dinner", "num_tags"

- **Zmienna zależna:**

"rating"

- **Podział na dane testowe i treningowe:**

X_train, X_test, y_train, y_test



Regresja liniowa

Wyniki

- **Mean Absolute Error** (Średni błąd predykcji):
0.49
- **Root Mean Squared Error**:
0.64
- **R2 score** (W ilu % nasz model przewiduje dobrze):
0.01



Drzewo decyzyjne

- Najlepsze parametry z wykorzystaniem krosvalidacji:

`max_depth': 14,`

`'min_samples_split': 160`

- **Mean Absolute Error** (Średni błąd predykcji):

0.49

- **Root Mean Squared Error:**

0.65

- **R2 score** (W ilu % nasz model przewiduje dobrze):

-0.86



Wyniki

Żaden z dwóch modeli nie przewiduje oceny w zadowalającym stopniu.

Model liniowy jest zdecydowanie lepszy od drzewa decyzyjnego, dla którego R^2 Score wyszedł ujemny.

Ocena dania na podstawie wykorzystanych przez nas danych jest niemożliwa do przewidzenia

Inny sposób agregacji i przekształceń danych mogą pomóc w przewidywaniu oceny.