

**Bag of Words, wynik F1-score: 0.5686**

- tworzy macierz wystąpień słów – dla każdego dokumentu liczy, ile razy dane słowo z ustalonego słownika się w nim pojawia. Nie uwzględnia kontekstu, tylko czystą częstotliwość.
- Prosta i szybka implementacja
- Nie rozpoznaje synonimów i kontekstu
- Im więcej cech tym rzadsze dane

**TF-IDF, wynik F1-score: 0.5580**

- TF-IDF nadaje wagę słowom, biorąc pod uwagę częstość ich występowania w dokumencie (jak w BoW) oraz rzadkość występowania w całym zbiorze danych. Słowa częste w całym korpusie są mniej istotne.
- Lepsza filtracja „szumu” (np. bardzo częstych, ale mało znaczących słów).
- Nadal brak kontekstu semantycznego (nie rozpoznaje znaczenia słów).

**Word2Vec (medium model), wynik F1-score: 0.5094**

- Umożliwia uchwycenie semantyki i kontekstu.
- Znacząco wolniejsza – wymaga przetwarzania tekstów z pomocą modelu "Medium" od spaCy, ponieważ model "sm" nie posiada wyuczonych "vector'ów".
- Osiągnął najgorszy wynik przy najdłuższym wykonywaniu, potencjalnie ponieważ test był robiony na 25k danych i reprezentacja przestrzenna nie oddała wystarczająco kontekstu klasyfikacji

Najlepsze wyniki uzyskał Bag-of-Words, co może wynikać z faktu, że klasyfikacja opierała się na obecności konkretnych słów ("skin", "love"...), a nie na ich kontekście. Możliwe jest, iż wynik jest spowodowany stosunkowo niską próbką w porównaniu z dostępnymi danymi.