

Anotacja korpusów oraz osadzenia słów i tekstów

Część II: Modele przestrzeni wektorowych

Autorzy: Oliwer Krupa, Adam Bednarski, Jan Masłowski, Łukasz Lenkiewicz

October 23, 2024

Contents

1	Word Embeddings	3
1.1	Osadzenia słów w anotowanym korpusie (Word2Vec i Fasttext)	3
1.1.1	Word2Vec	3
1.1.2	FastText	3
1.2	Porównanie k-najbardziej podobnych słów dla modeli Word2Vec i FastText	4
1.2.1	Wyniki dla Word2Vec	4
1.2.2	Wyniki dla FastText	4
1.2.3	Dyskusja wyników	5
2	Text Embeddings	5

1 Word Embeddings

1.1 Osadzenia słów w anotowanym korpusie (Word2Vec i FastText)

W celu zbadania i wizualizacji osadzeń słów w anotowanym korpusie, zastosowaliśmy dwie metody: **Word2Vec** oraz **FastText**. Dla obu modeli osadziliśmy słowa, a następnie zwizualizowaliśmy wyniki za pomocą algorytmu t-SNE, który redukuje wymiary danych, co pozwala na łatwiejszą analizę i porównanie rezultatów.

1.1.1 Word2Vec

Model **Word2Vec** został przetrenowany na korpusie tekstów polskich, a następnie osadziliśmy słowa z anotowanego korpusu, aby sprawdzić, jak różne słowa są reprezentowane w przestrzeni wektorowej. Wyniki zostały przedstawione na poniższym wykresie, gdzie kolory reprezentują różne kategorie sentymentu — niektóre słowa zostały oznaczone jako "Brak etykiety", a inne jako "Wzmacnianie".

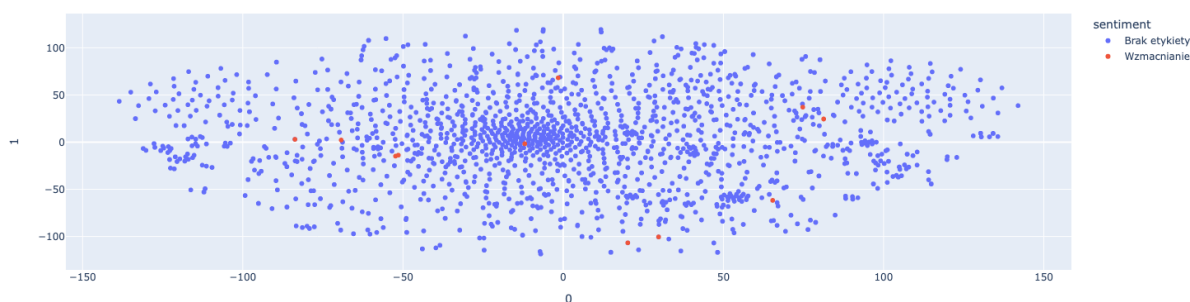


Figure 1: Wizualizacja osadzeń słów za pomocą Word2Vec

Jak widać na rysunku 1, większość słów jest zgrupowana w centralnej części wykresu, jednak niektóre słowa przypisane do kategorii "Wzmacnianie" znajdują się różnych obszarach przestrzeni wektorowej.

1.1.2 FastText

Podobną analizę przeprowadzono z użyciem modelu **FastText**, który również został przetrenowany na korpusie polskim. Wyniki osadzeń słów z tego modelu zostały przedstawione na poniższym wykresie, który podobnie jak poprzedni, został wygenerowany za pomocą t-SNE.

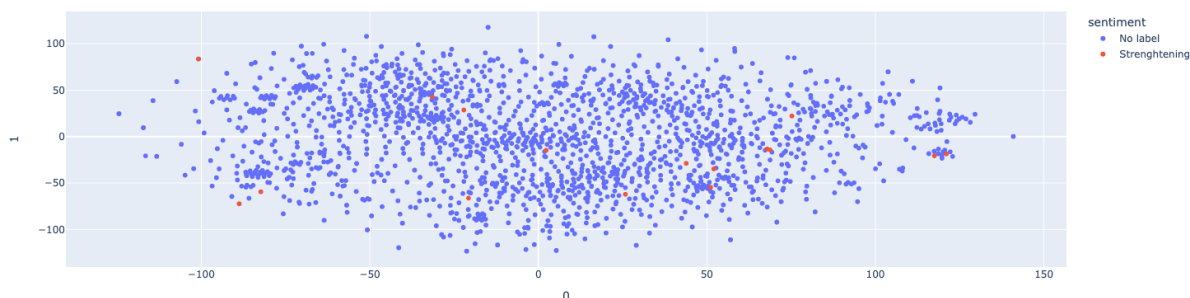


Figure 2: Wizualizacja osadzeń słów za pomocą FastText

Rysunek 2 pokazuje, że wyniki dla modelu FastText są podobne do tych uzyskanych z Word2Vec, chociaż niektóre grupy słów znajdują się w różnych częściach przestrzeni, co wynika z różnic w sposobie reprezentacji słów przez te dwa modele.

1.2 Porównanie k-najbardziej podobnych słów dla modeli Word2Vec i FastText

W ramach analizy wygenerowano listy pięciu najbardziej podobnych słów dla zaanotowanych terminów, wykorzystując modele **Word2Vec** oraz **FastText**. Wyniki pokazują różnice w sposobie modelowania podobieństw między słowami.

1.2.1 Wyniki dla Word2Vec

Dla modelu **Word2Vec**, lista najbardziej podobnych słów zawierała m.in.:

- **by**: aby (0.89), żeby (0.86), więc (0.69)
- **kto**: dlaczego (0.66), czemu (0.63), jeśli (0.59)
- **człowiek**: osobnik (0.64), osoba (0.64), tubylec (0.62)
- **seks**: sexu (0.76), masturbacja (0.71), erotyka (0.63)
- **fakt**: stwierdzenie (0.66), twierdzenie (0.65), hipoteza (0.59)

1.2.2 Wyniki dla FastText

Dla modelu **FastText**, podobne słowa to:

- **by**: By (0.77), żeby (0.61), muc (0.60)
- **kto**: Kto (0.80), ktoś (0.78), ktokolwiek (0.72)
- **człowiek**: Człowiek (0.80), czlowiek (0.76), człowiek. (0.69)
- **seks**: sex (0.80), seks. (0.78), seks- (0.76)
- **fakt**: faktem (0.68), faktu (0.66), Fakt (0.65)

1.2.3 Dyskusja wyników

Oba modele pokazują różnice w sposobie reprezentacji podobieństw. **FastText** lepiej radzi sobie z różnymi formami morfologicznymi i fleksyjnymi, np. dla słowa „człowiek” uwzględnia różne jego odmiany. Z kolei **Word2Vec** ma tendencję do generowania bardziej ogólnych podobieństw, co sprawia, że jest mniej precyzyjny w przypadku rzadkich wyrażen. FastText dzięki n-gramom lepiej odwzorowuje złożoność morfologiczną języka polskiego.

2 Text Embeddings