

News Selector

czyli jak automatycznie podsumować wiadomości z polskich serwisów informacyjnych

Jakub Kubajek

1. Wstęp
2. Pobieranie artykułów
3. Kluczowość słów
4. Grupowanie
5. Automatyczne podsumowanie

Wstęp

- Keynote speech **Lynn Cherny** podczas PyData w Warszawie:
Tl;dr: Summarisation

- Keynote speech **Lynn Cherny** podczas PyData w Warszawie:
Tl;dr: Summarisation
- Setki artykułów dziennie

Charakterystyka

- **Uczenie nienadzorowane**

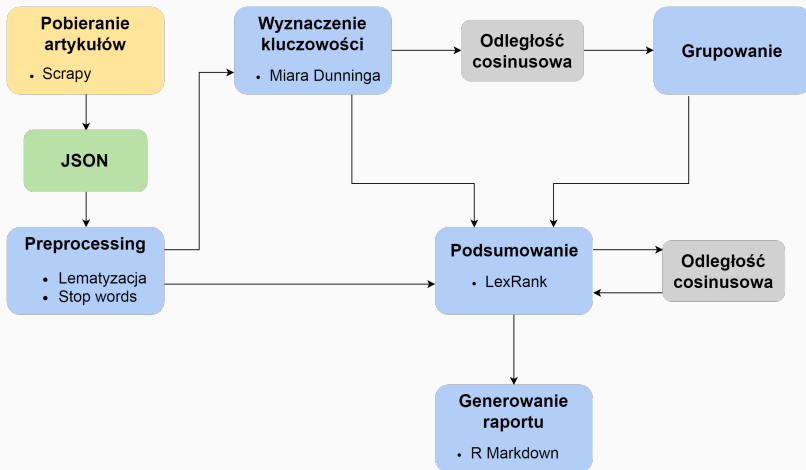
Charakterystyka

- Uczenie nienadzorowane
- Podejście słownikowe - *Morfeusz*

Charakterystyka

- Uczenie nienadzorowane
- Podejście słownikowe - *Morfeusz*
 - Sprowadzanie słów do lematu

Schemat modelu



Pobieranie artykułów

Scrapy

- Asynchroniczność
- 10x szybszy niż *rvest*
- Obsługa błędów



interia



GAZETA.PL

dziennik  **pl**



RMF **24**

ZET ^{radio}



Kluczowość słów

Informacje ogólne

- Test ilorazu wiarygodności (*Likelihood ratio test*)
- Brak założenia o normalności rozkładu
- Odpowiednie dla zjawisk o małej liczbie wystąpień

Informacje ogólne

- Test ilorazu wiarygodności (*Likelihood ratio test*)
- Brak założenia o normalności rozkładu
- Odpowiednie dla zjawisk o małej liczbie wystąpień

Wzór

$$\lambda = \frac{L(p, k_0, n_0) L(p, k_1, n_1)}{L(p_0, k_0, n_0) L(p_1, k_1, n_1)} \quad (1)$$

$$L(p, k, n) = p^k (1 - p)^{n-k} \quad (2)$$

gdzie, $p = \frac{k_0 + k_1}{n_0 + n_1}$, $p_0 = \frac{k_0}{n_0}$, zaś $-2\log\lambda$ ma rozkład χ^2 z jednym stopniem swobody. k_0 to liczba wystąpień danego słowa z wykluczeniem badanego dnia, k_1 to liczba wystąpień danego słowa w badanym dniu.

Modyfikacja miary Dunninga

- Mnożenie statystyki przez -1, gdy $p_1 < p_0$
- Brak logarytmizacji zer - $\max(\epsilon, p)$

Modyfikacja miary Dunninga

- Mnożenie statystyki przez -1, gdy $p_1 < p_0$
- Brak logarytmizacji zer - $\max(\epsilon, p)$

Kryteria wyboru słów

- Wartość statystyki χ^2 większa od 10
- Liczba wystąpień większa od wartości 90 percentyla

Grupowanie

Opis

- Miara podobieństwa dwóch wektorów
- Nie uwzględnia kolejności - *bag of words*
- Przyjmuje wartości z zakresu $[0, 1]$

Opis

- Miara podobieństwa dwóch wektorów
- Nie uwzględnia kolejności - *bag of words*
- Przyjmuje wartości z zakresu $[0, 1]$

Wzór

$$\cos(\theta) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\| \|\mathbf{B}\|} = \frac{\sum_{i=1}^n A_i B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}}$$

Podział

- Aglomeracyjny - zaczynamy od n grup
- Deglomeracyjny - zaczynamy o 1 grupy

Podział

- Aglomeracyjny - zaczynamy od n grup
- Deglomeracyjny - zaczynamy o 1 grupy

Metody łączenia w grupy

- Pojedyncze połączenie - $D(X, Y) = \min \{d(x, y) : x \in X, y \in Y\}$
- Kompletne połączenie - $D(X, Y) = \max \{d(x, y) : x \in X, y \in Y\}$
- Średnie połączenie - $D(X, Y) = \text{mean} \{d(x, y) : x \in X, y \in Y\}$

Grupowania w modelu

- Aglomeracyjne

Grupowania w modelu

- Aglomeracyjne
- Macierz podobieństwa między słowami (grupami) - wspólne występowanie w akapitach

Grupowania w modelu

- Aglomeracyjne
- Macierz podobieństwa między słowami (grupami) - wspólne występowanie w akapitach
- Aktualizacja macierzy podobieństwa po stworzeniu nowej grupy:

Grupowania w modelu

- Aglomeracyjne
- Macierz podobieństwa między słowami (grupami) - wspólne występowanie w akapitach
- Aktualizacja macierzy podobieństwa po stworzeniu nowej grupy:
 - Sumowanie wektorów dystrybucji połączonych grup

Grupowania w modelu

- Aglomeracyjne
- Macierz podobieństwa między słowami (grupami) - wspólne występowanie w akapitach
- Aktualizacja macierzy podobieństwa po stworzeniu nowej grupy:
 - Sumowanie wektorów dystrybucji połączonych grup
 - Wyznaczenie podobieństwa między nową grupą, a pozostałymi grupami

Grupowania w modelu

- Aglomeracyjne
- Macierz podobieństwa między słowami (grupami) - wspólne występowanie w akapitach
- Aktualizacja macierzy podobieństwa po stworzeniu nowej grupy:
 - Sumowanie wektorów dystrybucji połączonych grup
 - Wyznaczenie podobieństwa między nową grupą, a pozostałymi grupami
- Zakończenie grupowania

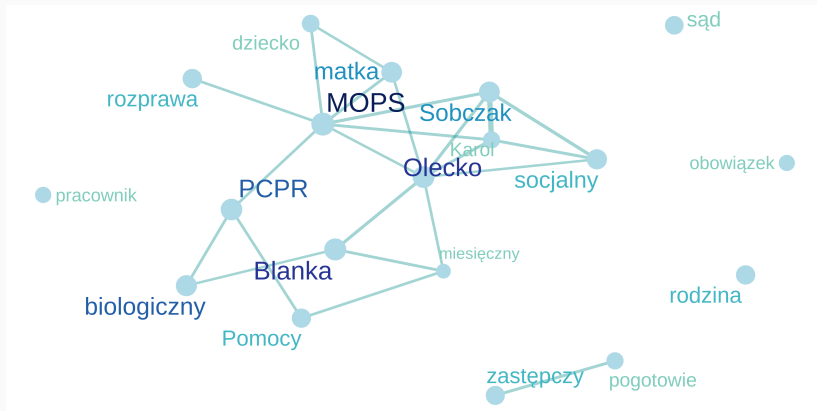
Grupowania w modelu

- Aglomeracyjne
- Macierz podobieństwa między słowami (grupami) - wspólne występowanie w akapitach
- Aktualizacja macierzy podobieństwa po stworzeniu nowej grupy:
 - Sumowanie wektorów dystrybucji połączonych grup
 - Wyznaczenie podobieństwa między nową grupą, a pozostałymi grupami
- Zakończenie grupowania
 - Minimalna wartość podobieństwa

Grupowania w modelu

- Aglomeracyjne
- Macierz podobieństwa między słowami (grupami) - wspólne występowanie w akapitach
- Aktualizacja macierzy podobieństwa po stworzeniu nowej grupy:
 - Sumowanie wektorów dystrybucji połączonych grup
 - Wyznaczenie podobieństwa między nową grupą, a pozostałymi grupami
- Zakończenie grupowania
 - Minimalna wartość podobieństwa
 - Algorytm Silhouette

Przykład



Automatyczne podsumowanie

LexRank

- Algorytm grafowy, bazujący na idei algorytmu *PageRank*

LexRank

- Algorytm grafowy, bazujący na idei algorytmu *PageRank*
- Wybiera te zdania, które zawierają najwięcej odniesień (słów) do innych zdań

LexRank

- Algorytm grafowy, bazujący na idei algorytmu *PageRank*
- Wybiera te zdania, które zawierają najwięcej odniesień (słów) do innych zdań
- Połączenia między zdaniami to podobieństwo cosinusowe

LexRank

- Algorytm grafowy, bazujący na idei algorytmu *PageRank*
- Wybiera te zdania, które zawierają najwięcej odniesień (słów) do innych zdań
- Połączenia między zdaniami to podobieństwo cosinusowe
- Wagi poszczególnych słów do wyznaczenia podobieństwa - $TF \cdot IDF$

LexRank

- Algorytm grafowy, bazujący na idei algorytmu *PageRank*
- Wybiera te zdania, które zawierają najwięcej odniesień (słów) do innych zdań
- Połączenia między zdaniami to podobieństwo cosinusowe
- Wagi poszczególnych słów do wyznaczenia podobieństwa - $TF \cdot IDF$
 - *TF (Term Frequency)* - liczba wystąpień danego słowa

LexRank

- Algorytm grafowy, bazujący na idei algorytmu *PageRank*
- Wybiera te zdania, które zawierają najwięcej odniesień (słów) do innych zdań
- Połączenia między zdaniami to podobieństwo cosinusowe
- Wagi poszczególnych słów do wyznaczenia podobieństwa - $TF \cdot IDF$
 - TF (*Term Frequency*) - liczba wystąpień danego słowa
 - IDF (*Inverse Document Frequency*) - $IDF_i = \log(\frac{D}{d_i}) + 1$, gdzie D to liczba zdań, a d_i to liczba zdań w których wystąpiło i-te słowo

LexRank

- Algorytm grafowy, bazujący na idei algorytmu *PageRank*
- Wybiera te zdania, które zawierają najwięcej odniesień (słów) do innych zdań
- Połączenia między zdaniami to podobieństwo cosinusowe
- Wagi poszczególnych słów do wyznaczenia podobieństwa - $TF \cdot IDF$
 - TF (*Term Frequency*) - liczba wystąpień danego słowa
 - IDF (*Inverse Document Frequency*) - $IDF_i = \log(\frac{D}{d_i}) + 1$, gdzie D to liczba zdań, a d_i to liczba zdań w których wystąpiło i-te słowo
- Ranking to stan ustalony łańcucha Markowa zbudowanego jako ważona suma:

LexRank

- Algorytm grafowy, bazujący na idei algorytmu *PageRank*
- Wybiera te zdania, które zawierają najwięcej odniesień (słów) do innych zdań
- Połączenia między zdaniami to podobieństwo cosinusowe
- Wagi poszczególnych słów do wyznaczenia podobieństwa - $TF \cdot IDF$
 - TF (*Term Frequency*) - liczba wystąpień danego słowa
 - IDF (*Inverse Document Frequency*) - $IDF_i = \log(\frac{D}{d_i}) + 1$, gdzie D to liczba zdań, a d_i to liczba zdań w których wystąpiło i-te słowo
- Ranking to stan ustalony łańcucha Markowa zbudowanego jako ważona suma:
 - **Macierzy podobieństwa między zdaniami**

LexRank

- Algorytm grafowy, bazujący na idei algorytmu *PageRank*
- Wybiera te zdania, które zawierają najwięcej odniesień (słów) do innych zdań
- Połączenia między zdaniami to podobieństwo cosinusowe
- Wagi poszczególnych słów do wyznaczenia podobieństwa - $TF \cdot IDF$
 - TF (*Term Frequency*) - liczba wystąpień danego słowa
 - IDF (*Inverse Document Frequency*) - $IDF_i = \log(\frac{D}{d_i}) + 1$, gdzie D to liczba zdań, a d_i to liczba zdań w których wystąpiło i-te słowo
- Ranking to stan ustalony łańcucha Markowa zbudowanego jako ważona suma:
 - Macierzy podobieństwa między zdaniami
 - Macierzy losowego przejścia między zdaniami

Modyfikacje

- Wybór artykułów - minimalny odsetek słów z tematu

Modyfikacje

- Wybór artykułów - minimalny odsetek słów z tematu
- Wagi słów w podobieństwie cosinusowym

Modyfikacje

- Wybór artykułów - minimalny odsetek słów z tematu
- Wagi słów w podobieństwie cosinusowym
 - TF - liczba wystąpień danego słowa

Modyfikacje

- Wybór artykułów - minimalny odsetek słów z tematu
- Wagi słów w podobieństwie cosinusowym
 - TF - liczba wystąpień danego słowa
 - Zmodyfikowana miara Dunninga

Modyfikacje

- Wybór artykułów - minimalny odsetek słów z tematu
- Wagi słów w podobieństwie cosinusowym
 - TF - liczba wystąpień danego słowa
 - Zmodyfikowana miara Dunninga
 - $\log(D_i + 1)$ dla $p_1 \geq p_0$

Modyfikacje

- Wybór artykułów - minimalny odsetek słów z tematu
- Wagi słów w podobieństwie cosinusowym
 - TF - liczba wystąpień danego słowa
 - Zmodyfikowana miara Dunninga
 - $\log(D_i + 1)$ dla $p_1 \geq p_0$
 - 0 dla $p_1 < p_0$

Modyfikacje

- Wybór artykułów - minimalny odsetek słów z tematu
- Wagi słów w podobieństwie cosinusowym
 - TF - liczba wystąpień danego słowa
 - Zmodyfikowana miara Dunninga
 - $\log(D_i + 1)$ dla $p_1 \geq p_0$
 - 0 dla $p_1 < p_0$
 - Podobieństwo cosinusowe występowania danego słowa do występowania tematu (grupy słów)

Modyfikacje

- Wybór artykułów - minimalny odsetek słów z tematu
- Wagi słów w podobieństwie cosinusowym
 - TF - liczba wystąpień danego słowa
 - Zmodyfikowana miara Dunninga
 - $\log(D_i + 1)$ dla $p_1 \geq p_0$
 - 0 dla $p_1 < p_0$
 - Podobieństwo cosinusowe występowania danego słowa do występowania tematu (grupy słów)
- Skalowanie wartości rankingu - $F_i = \frac{\log(w_i^T)}{\log(w_i)}$

Modyfikacje

- Wybór artykułów - minimalny odsetek słów z tematu
- Wagi słów w podobieństwie cosinusowym
 - TF - liczba wystąpień danego słowa
 - Zmodyfikowana miara Dunninga
 - $\log(D_i + 1)$ dla $p_1 \geq p_0$
 - 0 dla $p_1 < p_0$
 - Podobieństwo cosinusowe występowania danego słowa do występowania tematu (grupy słów)
- Skalowanie wartości rankingu - $F_i = \frac{\log(w_i^T)}{\log(w_i)}$
 - gdzie F_i to współczynnik skalujący, w_i liczba słów z tematu w i-tym zdaniu, a W liczba słów w temacie

Modyfikacje

- Wybór artykułów - minimalny odsetek słów z tematu
- Wagi słów w podobieństwie cosinusowym
 - TF - liczba wystąpień danego słowa
 - Zmodyfikowana miara Dunninga
 - $\log(D_i + 1)$ dla $p_1 \geq p_0$
 - 0 dla $p_1 < p_0$
 - Podobieństwo cosinusowe występowania danego słowa do występowania tematu (grupy słów)
- Skalowanie wartości rankingu - $F_i = \frac{\log(w_i^T)}{\log(w_i)}$
 - gdzie F_i to współczynnik skalujący, w_i liczba słów z tematu w i-tym zdaniu, a W liczba słów w temacie
- Pominięcie podobnych zdań ($\cos(\theta) > 0,4$)



T. Dunning.

Accurate methods for the statistics of surprise and coincidence.

COMPUTATIONAL LINGUISTICS, 19(1):61–74, 1993.



G. Erkan and D. R. Radev.

Lexrank: Graph-based lexical centrality as salience in text summarization.

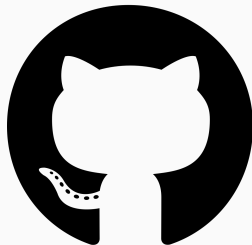
Journal of artificial intelligence research, 22:457–479, 2004.



L. Page, S. Brin, R. Motwani, and T. Winograd.

The pagerank citation ranking: Bringing order to the web.

Technical report, Stanford InfoLab, 1999.



jkubajek