

Eksploracja zasobów internetowych

Wykład 3

Wyszukiwanie dokumentów WWW bazujące na słowach kluczowych

Wstęp

Wyszukiwanie dokumentów za pomocą słów kluczowych bazujące na regułach boolowskich jest proste i szybkie, jednak posiada dużą wadę.

Nie pozwala na sortowanie zwróconych wyników wyszukiwania pod względem istotności treści.

W jaki sposób poradzić sobie z tym problemem?

Poprzez definiowanie zapytań w sposób precyzyjny lub sortowanie zwróconych treści bazując na ilości wystąpień poszczególnych termów w korpusie dokumentu.

Wstęp

Najbardziej powszechne struktury danych do przechowywania treści pobranych ze stron WWW:

- macierz term-dokument typu boolowskiego,
- macierz term-dokument typu ilościowego,
- macierz term-dokument typu pozycyjnego.

Z wykorzystaniem tych struktur można zwracać wyniki bazujące na słowach kluczowych stosując reguły boolowskie.

Wstęp

Przy wyszukiwaniu danych zgodnych z zapytaniem złożonym ze słów kluczowych ważne jest sortowanie pod względem istotności treści zboru wynikowego.

Osiągnięcie tego celu wymaga zastosowania odpowiednich algorytmów oraz zmodyfikowanych struktur do przechowywania danych.

Powstałe struktury bazują na poznanych do tej pory metodach przechowywania pobranych danych.

Dane przykładowe

Strona internetowa www.artsci.ccsu.edu:



The screenshot shows a web browser window displaying the website for the School of Arts & Sciences at Central Connecticut State University. The browser's address bar shows the URL <http://www.artsci.ccsu.edu/Departments.htm>. The website header features the CCSU logo, the university name, and the motto "Start with a Dream. Finish with a Future." Below the header, the page is titled "The School of Arts & Sciences" and "Departments". A link for "Department Chairs, Locations, Phone Numbers" is provided. A list of department names is displayed in two columns, with "Mathematical Sciences" highlighted. At the bottom, there are navigation links for "A&S Home", "A-Z Directory", "Departments", and "About the School", along with a page update date of 10/27/04.

File Edit View Favorites Tools Help

Back Forward Stop Refresh Home Search Favorites

Address <http://www.artsci.ccsu.edu/Departments.htm> Go

CCSU Central Connecticut State University *Start with a Dream. Finish with a Future.*

The School of Arts & Sciences

Departments

[Department Chairs, Locations, Phone Numbers](#)

Anthropology	History
Art	Mathematical Sciences
Biological Sciences	Modern Languages
Chemistry	Music
Communication	Philosophy
Computer Science	Physics/Earth Sciences
Criminal Justice	Political Science
Design	Psychology
Economics	Sociology
English	Theatre
Geography	

[[A&S Home](#)] [[A-Z Directory](#)] [[Departments](#)] [[About the School](#)]

page last updated: 10/27/04

Dane przykładowe

Przykładowy zbiór danych ze strony www.artsci.ccsu.edu:

Document ID	Document name	Words	Terms
d ₁	Anthropology	114	86
d ₂	Art	153	105
d ₃	Biology	123	91
d ₄	Chemistry	87	58
d ₅	Communication	124	88
d ₆	Computer Science	101	77
d ₇	Criminal Justice	85	60
d ₈	Economics	107	76
d ₉	English	116	80
d ₁₀	Geography	95	68
d ₁₁	History	108	78
d ₁₂	Mathematics	89	66
d ₁₃	Modern Languages	110	75
d ₁₄	Music	137	91
d ₁₅	Philosophy	85	54
d ₁₆	Physics	130	100
d ₁₇	Political Science	120	86
d ₁₈	Psychology	96	60
d ₁₉	Sociology	99	66
d ₂₀	Theatre	116	80
Total number of words/terms		2195	1545
Number of different words/terms		744	671

Analiza danych

Pod względem możliwości analizy danych pobranych ze stron WWW najbardziej odpowiednia jest struktura typu pozycyjnego:

- względna prostota przechowywania danych,
- łatwość wyszukiwania informacji,
- możliwość zliczenia ilości wystąpień termów w dokumencie,
- możliwość wyszukiwania słów leżących w określonych odległościach względem siebie.

Macierz typu pozycyjnego

Przykładowa macierz term-dokument typu pozycyjnego:

DID	lab	laboratory	programming	computer	program
d ₁	0	0	0	0	[71]
d ₂	0	0	0	0	[7]
d ₃	0	[65,69]	0	[68]	0
d ₄	0	0	0	[26]	[30,43]
d ₅	0	0	0	0	0
d ₆	0	0	[40,42]	[1,3,7,13,26,34]	[11,18,61]
d ₇	0	0	0	0	[9,42]
d ₈	0	0	0	0	[57]
d ₉	0	0	0	0	0
d ₁₀	0	0	0	0	0
d ₁₁	0	0	0	0	0
d ₁₂	0	0	0	[17]	0
d ₁₃	0	0	0	0	0
d ₁₄	[42]	0	0	[41]	[71]
d ₁₅	0	0	0	0	[37,38]
d ₁₆	0	0	0	0	[81]
d ₁₇	0	0	0	0	[68]
d ₁₈	0	0	0	0	0
d ₁₉	0	0	0	0	[51]
d ₂₀	0	0	0	0	0

Model wektorowy

Model wektorowy jest bezpośrednio powiązany z macierzowymi strukturami danych i wynikach z ich innej interpretacji logicznej.

Struktury modelu wektorowego pozwalają klasyfikować zbiór wyników pod kątem istotności traktując dokumenty jako wektory wielowymiarowe.

Każdy z takich wektorów posiada ilość współrzędnych równą ilości termów we wszystkich zbiorach dokumentów.

Model wektorowy

Rodzaje powszechnie stosowanych modeli wektorowych w systemach typu *Information Retrieval*:

- boolowski,
- *Term-Frequency* (TF),
- *Inverse Document Frequency* (IDF),
- *Term Frequency – Inverse Document Frequency* (TFIDF).

Model wektorowy – oznaczenia

Oznaczenia wykorzystywane na kolejnych slajdach odnoszące się do struktur logicznych:

- d_1, d_2, \dots, d_n – dokumenty
- t_1, t_2, \dots, t_m – termy
- n_{ij} – ilość termów t_i w dokumencie d_j
- m – ilość wszystkich termów
- n – ilość wszystkich dokumentów

Model wektorowy – boolowski

Pojedynczy wiersz w boolowskim modelu wektorowym opisany jest jako:

$$\vec{d}_j = (d_j^1 d_j^2 \dots d_j^n) \quad d_i^j = \begin{cases} 0 & \text{dla } n_{ij} = 0 \\ 1 & \text{dla } n_{ij} > 0 \end{cases}$$

Dla zbioru termów *lab*, *laboratory*, *programming*, *computer* oraz *program* dla dokumentu d_6 wektor będzie zadany jako:

$$\vec{d}_6 = (0 \ 0 \ 1 \ 1 \ 1)$$

Macierz typu boolowskiego

Przykładowa macierz term-dokument typu boolowskiego:

DID	lab	laboratory	programming	computer	program
d ₁	0	0	0	0	1
d ₂	0	0	0	0	1
d ₃	0	1	0	1	0
d ₄	0	0	0	1	1
d ₅	0	0	0	0	0
d ₆	0	0	1	1	1
d ₇	0	0	0	0	1
d ₈	0	0	0	0	1
d ₉	0	0	0	0	0
d ₁₀	0	0	0	0	0
d ₁₁	0	0	0	0	0
d ₁₂	0	0	0	1	0
d ₁₃	0	0	0	0	0
d ₁₄	1	0	0	1	1
d ₁₅	0	0	0	0	1
d ₁₆	0	0	0	0	1
d ₁₇	0	0	0	0	1
d ₁₈	0	0	0	0	0
d ₁₉	0	0	0	0	1
d ₂₀	0	0	0	0	0

Model wektorowy – TF

Pojedynczy wiersz w modelu wektorowym typu *Term-Frequency* opisany jest jako:

$$\vec{d}_j = (d_j^1 d_j^2 \dots d_j^n)$$

Każdy współczynnik wektora opisany jest zależnością:

$$d_j^i = TF(t_i, d_j)$$

Model wektorowy – TF

Sposoby obliczania współczynników TF:

- suma termów:
$$TF(t_i, d_j) = \begin{cases} 0 & \text{dla } n_{ij} = 0 \\ \frac{n_{ij}}{\sum_{k=1}^m n_{kj}} & \text{dla } n_{ij} > 0 \end{cases}$$

- maksimum:
$$TF(t_i, d_j) = \begin{cases} 0 & \text{dla } n_{ij} = 0 \\ \frac{n_{ij}}{\max_k n_{kj}} & \text{dla } n_{ij} > 0 \end{cases}$$

- logarytm:
$$TF(t_i, d_j) = \begin{cases} 0 & \text{dla } n_{ij} = 0 \\ 1 + \log(1 + \log n_{ij}) & \text{dla } n_{ij} > 0 \end{cases}$$

Model wektorowy – TF

Współczynniki *TF* służą do normalizacji wartości termów opisujących dokumenty WWW. Dzięki nim można przeskalować duże wartości termów związanych z dokumentami do wartości mniejszych, mieszczących się w zdefiniowanym zakresie.

Współczynniki *TF* powiązane są z każdym termem w każdym dokumencie. Do zapisania współczynników *TF* wymagana jest macierz dwuwymiarowa (w sensie logicznym!).

Model wektorowy – IDF

Założmy, że zbiór D jest zbiorem wszystkich dokumentów, zaś zbiór $D_{t_i} = \{d_j | n_{ij} > 0\}$ zbiorem dokumentów zawierających term t_i .

Sposoby obliczania współczynników *IDF*:

- ułamek:
$$IDF(t_i) = \frac{|D|}{|D_{t_i}|}$$
- logarytm:
$$IDF(t_i) = \log \frac{1 + |D|}{|D_{t_i}|}$$

Model wektorowy – IDF

Współczynniki *Inverse Document Frequency* służą do skalowania współczynników wektorów dokumentów.

Dla termów występujących często w różnych dokumentach, istotność tego termu nie może być tak duża, jak termu występującego w niewielu dokumentach.

Czy współczynniki *IDF* powiązane są tylko z termami czy również z dokumentami?

Współczynniki *IDF* powiązane są tylko z termami (bez uwzględniania dokumentów).

Model wektorowy – *TFIDF*

Pojedynczy wiersz w modelu wektorowym typu *Term-Frequency Inverse Document Frequency* opisany jest jako:

$$\vec{d}_j = (d_j^1 d_j^2 \dots d_j^n)$$

Każdy współczynnik wektora opisany jest zależnością:

$$d_j^i = TF(t_i, d_j) IDF(t_i)$$

Model wektorowy *TFIDF* łączy w sobie zalety współczynników *TF* oraz współczynników *IDF*.

TFIDF – przykład

Wektor *TF* dokumentu d_6 (strona wydziału *Computer Science*):

$$\vec{d}_6 = (0 \ 0 \ 0,026 \ 0,076 \ 0,039)$$

Współczynniki modelu *IDF* dla poszczególnych termów (logarytm):

lab	laboratory	Programming	computer	program
3.04452	3.04452	3.04452	1.43508	0.559616

Wektor *TFIDF* dokumentu d_6 :

$$\vec{d}_6 = (0 \ 0 \ 0,079 \ 0,112 \ 0,022)$$

Macierz typu pozycyjnego

Przykładowa macierz term-dokument typu pozycyjnego:

DID	lab	laboratory	programming	computer	program
d ₁	0	0	0	0	[71]
d ₂	0	0	0	0	[7]
d ₃	0	[65,69]	0	[68]	0
d ₄	0	0	0	[26]	[30,43]
d ₅	0	0	0	0	0
d ₆	0	0	[40,42]	[1,3,7,13,26,34]	[11,18,61]
d ₇	0	0	0	0	[9,42]
d ₈	0	0	0	0	[57]
d ₉	0	0	0	0	0
d ₁₀	0	0	0	0	0
d ₁₁	0	0	0	0	0
d ₁₂	0	0	0	[17]	0
d ₁₃	0	0	0	0	0
d ₁₄	[42]	0	0	[41]	[71]
d ₁₅	0	0	0	0	[37,38]
d ₁₆	0	0	0	0	[81]
d ₁₇	0	0	0	0	[68]
d ₁₈	0	0	0	0	0
d ₁₉	0	0	0	0	[51]
d ₂₀	0	0	0	0	0

Słowa kluczowe – wyszukiwanie

Korzystając z modelu wektorowego można przeprowadzić wyszukiwanie dokumentów bazując na obliczaniu odległości pomiędzy wektorami.

W jaki sposób przedstawić zapytanie w formie wektora?

Słowa kluczowe zapytania można przekształcić do termów i bazując na zbiorze wszystkich termów stworzyć wektor reprezentujący to zapytanie.

Słowa kluczowe – wyszukiwanie

Zwracane wyniki będzie można posortować pod względem istotności bazującej na termach poprzez obliczenie odległości pomiędzy wektorem reprezentującym zapytanie, a pozostałymi wektorami reprezentującymi dokumenty.

Odległości pomiędzy wektorami muszą być obliczane z wykorzystaniem wybranej normy metrycznej.

Wyszukiwanie – normy

Najczęściej w technice IR wykorzystywane są następujące normy metryczne:

- norma Euklidesowa:

$$\|\vec{q} - \vec{d}_j\| = \sqrt{\sum_{i=1}^m (q^i - d_j^i)^2}$$

- podobieństwo cosinusowe:

$$\vec{q} \cdot \vec{d}_j = \sum_{i=1}^m q^i d_j^i$$

Wyszukiwanie – normy

Które ze zwróconych dokumentów dla normy Euklidesowej są bardziej istotne?

Dla podobieństwa euklidesowego dokument jest tym bardziej istotny, im mniejsza jest wartość wynikowa normy.

Które ze zwróconych dokumentów dla normy cosinusowej są bardziej istotne?

Dla podobieństwa cosinusowego dokument jest tym bardziej istotny, im większa jest wartość wynikowa normy.

Wyszukiwanie – przykład

Założmy, że korzystając z poznanych wcześniej metod chcemy wyszukać te dokumenty, które zawierają termy *computer* oraz *program*.

Wektor q odpowiadający zapytaniu jest określony współrzędnymi:

$$\vec{q} = (0 \ 0 \ 0 \ 0,5 \ 0,5)$$

Po nałożeniu współczynników *IDF* na wektor q otrzymamy:

$$\vec{q} = (0 \ 0 \ 0 \ 0,718 \ 0,28)$$

Wyszukiwanie – przykład

Doc	TFIDF Coordinates (normalized)					$\bar{q} \cdot \bar{d}_j$ (rank)	$ \bar{q} - \bar{d}_j $ (rank)
d ₁	0	0	0	0	1	0.363	1.129
d ₂	0	0	0	0	1	0.363	1.129
d ₃	0	0.972	0	0.234	0	0.218	1.250
d ₄	0	0	0	0.783	0.622	0.956 (1)	0.298 (1)
d ₅	0	0	0	0	1	0.363	1.129
d ₆	0	0	0.559	0.811	0.172	0.819 (2)	0.603 (2)
d ₇	0	0	0	0	1	0.363	1.129
d ₈	0	0	0	0	1	0.363	1.129
d ₉	0	0	0	0	0	0	1
d ₁₀	0	0	0	0	0	0	1
d ₁₁	0	0	0	0	0	0	1
d ₁₂	0	0	0	1	0	0.932	0.369
d ₁₃	0	0	0	0	0	0	1
d ₁₄	0.890	0	0	0.424	0.167	0.456 (3)	1.043 (3)
d ₁₅	0	0	0	0	1	0.363	1.129
d ₁₆	0	0	0	0	1	0.363	1.129
d ₁₇	0	0	0	0	1	0.363	1.129
d ₁₈	0	0	0	0	0	0	1
d ₁₉	0	0	0	0	1	0.363	1.129
D ₂₀	0	0	0	0	0	0	1

Wyniki wyszukiwania

Zarówno norma euklidesowa, jak i norma cosinusowa pozwalają zwracać zbiory dokumentów posortowane pod względem ich istotności. Obliczana istotność zależy od częstości występujących w ich korpusach termów. Jednak normy te nie uwzględniają faktu występowania wyszukiwanego termu w treści dokumentu.

Należy zatem pamiętać, aby brać pod uwagę tylko te dokumenty, które zawierają wszystkie termy z zapytania opartego o słowa kluczowe.

Wyszukiwanie z operatorami

Silniki wyszukiwania mają możliwość przetwarzania zapytań złożonych ze słów kluczowych wraz z operatorami boolowskimi, np. *AND*, *OR* lub *NOT*.

Wprowadzając domyślne zapytanie złożone ze słów kluczowych, wykorzystywany jest operator *AND*.

W jaki sposób zrealizować operatory *OR* oraz *NOT*?

Przed obliczeniem odległości (lub w trakcie obliczeń) można wybierać tylko, te dokumenty, które spełniają podane zależności.

Wyszukiwanie słów z błędami

Problemem w definiowaniu zapytań złożonych ze słów kluczowych są możliwe błędy w zapisie poszczególnych słów. Jedną z metod radzenia sobie z tym problemem jest dekompozycja termów na *n-gramy*.

W przypadku popełnienia błędu w zapisie termu, porównanie fragmentów termów pozwoli na znalezienie podobieństwa i zwrócenie odpowiednich wyników.

Średnia długość stosowanych *n-gramów* waha się w zakresie od 2 do 4.

Wyszukiwanie słów z błędami

Przykładowo, term *program* może być rozłożony na 2-gramy:

$\{pr, ro, og, gr, ra, am\}$

Term *program* zapisany z błędem, np. *prorgam* zostanie rozłożony na następujące 2-gramy:

$\{pr, ro, or, rg, ga, am\}$

Porównanie dwóch sekwencji pokazuje, że 2-gramy pokrywają się w 3 na 6 przypadków, w związku z tym można podejrzewać, że termy te są takie same.

Sprzężenie zwrotne oceny wyników

W procesie sortowania wyszukanych dokumentów pod względem istotności, można wykorzystać ocenę zwróconych wyników przez użytkownika na zasadzie sprzężenia zwrotnego.

Użytkownik przydziela wyniki do dwóch zbiorów:

- D_+ – dokumenty istotne,
- D_- – dokumenty nieważne.

Problemem w tym rozwiązaniu jest jednak czas, który użytkownik musi poświęcić na klasyfikację dokumentów.

Pseudo-ocena wyników: np. 5 wyników do D_+ , reszta do D_- .

Sprzężenie zwrotne oceny wyników

Po dokonaniu oceny przez użytkownika wektor zapytania jest przeliczany z wykorzystaniem metody *Rocchio*. Dokumenty istotne zwiększają współczynniki wektora zapytania, zaś dokumenty nieistotne współczynniki te osłabiają.

Zależność modyfikująca wektor zapytania q :

$$\vec{q}' = \alpha \vec{q} + \beta \sum_{d_j \in D_+} \vec{d}_j - \gamma \sum_{d_j \in D_-} \vec{d}_j$$

Ocena wyników – przykład

Założmy, że $\alpha = 1$, $\beta = 0.5$, zaś $\gamma = 0$. Zapytanie q zostanie zmodyfikowane przez trzy istotne dokumenty zwrócone przez oryginalne zapytanie. Dodatkowo, wybrane zostaną trzy termy o najwyższych współczynnikach *IDF*: *lab*, *laboratory*, *programming*. Modyfikacja wektora zapytania q będzie wyglądała następująco:

$$\vec{q}' = \vec{q} + 0.5 (\vec{d}_4 + \vec{d}_6 + \vec{d}_{14})$$

$$\vec{q}' = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0.932 \\ 0.363 \end{pmatrix} + \left[\begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix} + \begin{pmatrix} 0 \\ 0 \\ 0.559 \\ 0 \\ 0 \end{pmatrix} + \begin{pmatrix} 0.89 \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix} \right] = \begin{pmatrix} 0.445 \\ 0 \\ 0.28 \\ 0.932 \\ 0.363 \end{pmatrix}$$

Ocena wyników – przykład

Po modyfikacji wektora zapytania q i ponownym zwróceniu wyników zgodnie z podobieństwem liczonym za pomocą miary cosinusów, otrzymane zostały następujące wartości dla dokumentów:

- d_6 (*Computer Science*): 0.863,
- d_4 (*Chemistry*): 0.846,
- d_{14} (*Music*): 0.754.

Wynik jest lepszy, gdyż zapytanie złożone z termów *computer* oraz *program*, powinno na pierwszym miejscu zwrócić stronę WWW dotyczącą wydziału *Computer Science*.

Miary oceny jakości wyników

W systemach typu *IR* istotne są mechanizmy pomiaru jakości działania algorytmów zwracających zbiory wynikowe. W tym celu opracowany został system *precision-recall*.

Model ten operuje na dwóch zbiorach:

- R_q – zbiór dokumentów zwróconych przez algorytm zgodnie z zapytaniem q ,
- D_q – zbiór istotnych dokumentów zgodnych z zapytaniem q , stworzony przez eksperta.

Miary oceny jakości wyników

Ilość zwróconych i istotnych wyników w stosunku do wszystkich istotnych dokumentów jest określany jako *recall*:

$$recall = \frac{|D_q \cap R_q|}{|D_q|}$$

Wartości parametru *recall* mogą się zmieniać w zakresie od 0 do 1.

Jaki jest najgorszy i najlepszy przypadek?

Najgorszy przypadek to wartość 0, kiedy algorytm nie zwrócił żadnego istotnego dokumentu. Najlepszy przypadek to 1, jednak wcale nie oznacza to, że algorytm zwrócił poprawny zbiór wyników.

Miary oceny jakości wyników

Ilość zwróconych i istotnych wyników w stosunku do wszystkich zwróconych dokumentów jest nazywany precyzją (ang. *precision*):

$$precision = \frac{|D_q \cap R_q|}{|R_q|}$$

Wartości precyzji mogą się zmieniać w zakresie od 0 do 1.

Jaki jest najgorszy i najlepszy przypadek?

Najgorszy przypadek to wartość 0, kiedy algorytm nie zwrócił żadnego istotnego dokumentu. Najlepszy przypadek to 1, kiedy zbiór wynikowy zwrócony przez algorytm zawiera wszystkie istotne dokumenty.

Miary oceny jakości wyników

Każdy system *IR* dąży do tego, aby jednoczesna wartość obydwu współczynników była równa 1, jednak jest to praktycznie niemożliwe.

Modyfikacja zapytań pod kątem uzyskania maksymalnej wartości (ale nie jednocześnie!) *precision* **lub** *recall* jest banalna.

Aby *recall* wyniósł 1, należy tworzyć ogólne zapytania, np. złożone z 1 słowa kluczowego.

Aby *precision* wyniósł 1, należy tworzyć tak szczegółowe zapytania, że będą dotyczyły tylko jednego dokumentu.

Miary oceny jakości wyników

W rzeczywistych systemach *IR* występuje zawsze jedna z zależności:

- $D_q \cap R_q \subset D_q$ – w tym przypadku zbiór wynikowy trzeba powiększyć,
- $D_q \subset R_q$ – w tym przypadku, zbiór wynikowy trzeba zmniejszyć.

Działanie skutecznego systemu *IR* musi być zawsze kompromisem pomiędzy wartością współczynnika *precision* i wartością współczynnika *recall*.

Dziękuję za uwagę!

