

Spis treści

1	Wstęp [TODO]	3
2	State of the art	4
3	Algorytmy	6
3.1	Folksonomia	6
3.2	SocialPageRank	6
3.2.1	Opis	6
3.2.2	Wyniki algorytmu dla przykładowych danych	7
3.2.3	Implementacja	8
3.2.4	Baza danych	8
3.2.5	Tworzenie i mnożenie macierzy	9
3.2.6	Wyniki - czesciowe: TODO	9
3.2.7	Problemy : TODO	10
3.2.8	Przykładowe wyniki wyszukiwarki : TODO	10
4	Adapted PageRank	11
4.1	Opis	11
4.1.1	Algorytm FolkRank	12
4.1.2	Przykładowe wyniki	12
4.2	Implementacja	13
4.3	Wyniki	14
4.4	Porównanie algorytmu SocialPageRank i algorytmu Adapted PageRank : TODO	14
4.4.1	Porównanie na małym przykładzie	14
4.4.2	Porównanie w działającym systemie	14
5	Opis aplikacji [TODO]	15
5.1	Folksonomia	15
5.2	delicious.com : TODO	15
5.3	Architektura	15
5.3.1	Wątki	15
5.3.2	Dane	16
5.3.3	Dane testowe TODO	18

5.4	Baza danych	18
5.5	Interface użytkownika TODO	19
5.6	Lucene	20
5.6.1	Pobieranie stron	20
5.6.2	Wyszukiwanie	21
5.7	Wyszukiwanie TODO	22
6	Wyniki [TODO]	23

Rozdział 1

Wstęp [TODO]

TODO

Rozdział 2

State of the art

Obecnie prowadzonych jest wiele badań nad sieciami społecznymi używającymi etykiet (tzw. tagów, czy adnotacji) do porządkowania zasobów użytkownika. Część z nich za cel stawia sobie ulepszenie algorytmów wyszukiwania dokumentów, inne używają tych zasobów dla personalizacji wyników dla użytkowników, rekomendacji dokumentów, użytkowników, czy tagów w czasie procesu etykietowania.

Autorzy prac [9] i [5] analizują użyteczność etykietowania dokumentów przy wyszukiwaniu. Ich analizy bazują głównie na danych systemu delicious. Wnioski z tych prac są pozytywne. Zwracają one uwagę na korelacje między popularnością stron w serwisie, a ich rankingiem według algorytmu PageRank. Dodatkowo wskazują na problemy z nowymi dokumentami dodanymi do sieci, z którymi nie radzą sobie algorytmy bazujących na odnośnikach pomiędzy dokumentami, a które mogą być rozwiązane przez strony takie jak delicious. Również rozwiązanie problemu nowych danych jest zaproponowane w pracy [2]. W tym artykule, jako rozwiązania, autorzy proponują użycie informacji pobranych z systemu twitter, który służącego do mikro-blogowania. Pod uwagę brane są krótkie informacje umieszczane przez użytkowników zawierające odnośniki do dokumentów jak również inne informacje dostępne w profilu użytkownika.

Autorzy pracy [3] zaprezentowali formalny model folksonomi i algorytmy AdaptedPageRank i jego bardziej spersonalizowana wersję algorytm FolkRank. Oba te algorytmy bazują na metodzie PageRank. Algorytm ten pozwala na rekomendacje tagów dla użytkowników jak również na ustalenie rankingu w wyszukiwanych elementach.

W pracy [1] przedstawione zostały algorytmy SocialSimRank i SocialPageRank. Algorytm SocialPageRank jest statycznym rankingiem zasobów, opartym również o ideę algorytmu PageRank. Algorytm ten bierze pod uwagę ilość tagów wskazujących na dany dokument jak również różną wagę opisujących dokument etykiet. Drugim proponowanym w tej pracy algorytmem jest SocialSimRank, który estymuje podobieństwo między używanymi

tagami, a następnie wyniki te są używane dla wyliczanie podobieństwa między zapytaniem użytkownika a tagami przypisanymi do danego zasobu.

W pracach [8] i [6] autorzy biorą pod uwagę całą sieć społecznościową użytkownika, dodając dodatkowo znajomości użytkowników i wykorzystują te dane dla przedstawienia spersonalizowanych wyników. W artykule [8] opisywany jest algorytm ContextMerge, który wykonuje wyszukiwanie biorąc pod uwagę zależności między użytkownikami. Dodatkowo pozwala na dynamiczne dodawanie nowych wyników do odpowiedzi dla uzyskania k . W pracy [6] autorzy zaproponowali system, który łączy wyniki wyszukiwarki z danymi użytkownikami aplikacji takich jak delicious. Jako bazowa wyszukiwarka może być użyta dowolna aplikacja, która zwraca dokumenty wraz z przypisanym do nich rankingiem. Wyniki te są następnie przeliczane za pomocą danych uzyskanych od użytkownika, czyli ze zbiorem dokumentów i tagów opisanych przez niego i innych użytkowników. Jako rezultat otrzymujemy nowy, bardziej spersonalizowany ranking dokumentów.

W pracy [7] celem autorów było spersonalizowanie wyników wyszukiwania w dokumentach z tagami. Do tego celu stworzone zostały dwa modele: użytkownik-tag, który ukazuje jak użytkownik użył tagi podobne do wybranego tagu. Drugim modelem jest model tag-zasób opisujący w jaki sposób dokumenty podobne do danego zasobu zostały opisane etykietami.

Probabilistyczne podejście do rankingu dokumentów przy użyciu tagów zostało zaprezentowane w pracy [10] i [4]. W [10] autorzy zaprezentowali model probabilistyczny dla generowania etykiet i zależnych im dokumentów i wyszukiwania ich tematów. W [4] autorzy prezentują metodę rankingu tagów w zależności od ich tematu i konstruują graf przejścia między tagami należącymi do różnych tematów. Metoda ta jest wykorzystana następnie dla rekomendacji tagów dla użytkownika w czasie opisywania dokumentów.

Rozdział 3

Algorytmy

3.1 Folksonomia

Folksonomią nazywamy społeczne klasyfikowanie, czy tagowanie różnych zasobów. Formalnie model folksonomii został przedstawiony w pracy [3]. Można ją przedstawić jako krotkę $F := (U, T, R, Y)$, gdzie: U, T, R to zbiory skończone, których elementy składają się odpowiednio z użytkowników, tagów i dokumentów. Y jest relacją “przypisania tagu” pomiędzy tymi elementami $Y \subseteq U \times T \times R$

Użytkownicy i tagi identyfikowani są na podstawie ich unikalnych nazw własnych. Dokumenty mogą być różnymi danymi: stronami www, zdjęciami, plikami np: pliki pdf. Ta praca bazuje na danych pobranych z witryny delicious, które w zdecydowanej większości są stronami www. Dane które nie są stroną www nie są brane pod uwagę w tej pracy.

3.2 SocialPageRank

3.2.1 Opis

SocialPageRank jest algorytmem wyliczającym statyczny ranking stron z perspektywy użytkownika sieci. Algorytm bazuje na obserwacji relacji między popularnymi stronami, tagami i aktywnymi użytkownikami. Popularne strony są dodawane przez aktywnych się użytkowników, które są opisywane popularnymi tagami. Aktywni się użytkownicy używają popularnych tagów dla popularnych stron. Popularne tagi używane są do annotacji popularnych stron przez ważnych użytkowników.

Bazując na powyższych założeniach algorytm propaguje i wzmacnia zależności między popularnymi tagami, użytkownikami i dokumentami.

Dane wejściowe:

N_T : ilość tagów

N_U : ilość użytkowników
 N_D : ilość dokumentów
 M_{DU} : macierz $N_D \times N_D$ asocjacyjna między dokumentami a użytkownikami
 M_{UT} : macierz $N_U \times N_T$ asocjacyjna między użytkownikami a tagami
 M_{TD} : macierz $N_T \times N_D$ asocjacyjna między tagami a dokumentami
 P_0 : wektor, o długości N_D ,

Inicjalizacja

W komórce macierzy $M_{DU}(d_n, u_k)$ znajduje się wartość będąca ilością adnotacji przypisanych do dokumentu d_n przez użytkownika u_k . Podobnie dla pozostałych macierzy, elementy $M_{UT}(u_k, t_n)$ to ilość dokumentów opisanych tagiem t_n przez użytkownika u_k , elementy $M_{TD}(t_n, d_k)$: ile użytkowników dodawało dokument d_k i oznaczyło go anotacją t_n .

Wektor P_0 zainicjalizowany został losowymi wartościami z przedziału $[0, 1]$. Jest on pierwszym przybliżeniem rank dokumentów.

repeat

$$U_i = M_{DU}^T * P_i$$

$$T_i = M_{UT}^T * U_i$$

$$P'_i = M_{TD}^T * T_i$$

$$T'_i = M_{TD} * P'_i$$

$$U'_i = M_{UT} * T'_i$$

$$P_{(i+1)} = M_{DU} * U'_i$$

until wartości wektora P_n nie zbiegną

Złożoność

Złożoność czasowa każdej iteracji wynosi $O(N_u * N_d + N_t * N_d + N_t * N_u)$.

3.2.2 Wyniki algorytmu dla przykładowych danych

W poniższej tabelce znajdują się dane, dla których zostało sprawdzone działanie algorytmu Social PageRank. Dane są nie duże i składają się z trzech różnych dokumentów, dwóch użytkowników i trzech tagów.

Strony www	użytkownicy	
	użytkownik 1	użytkownik 2
http://www.ted.com/	inspiration	
http://www.colourlovers.com/	design	inspiration
http://www.behance.net/	portfolio, design	portfolio, inspiration

Dla takich danych macierz $M_{d,u}$ mówiąca o zależności dokumentów z użytkownikami ma postać:

$$M_{d,u} = \begin{pmatrix} 1 & 0 \\ 1 & 1 \\ 1 & 2 \end{pmatrix}$$

Macierz użytkowników i tagów, $M_{u,t}$:

$$M_{u,t} = \begin{pmatrix} 1 & 1 & 1 \\ 2 & 0 & 1 \end{pmatrix}$$

Macierz tagów i dokumentów, $M_{t,d}$:

$$M_{t,d} = \begin{pmatrix} 1 & 1 & 1 \\ 0 & 1 & 1 \\ 0 & 0 & 2 \end{pmatrix}$$

wyniki:

Dla powyższych danych wyniki algorytmu zbiegają po czterech iteracjach z dokładnością $|P_3 - P_4| < 10^{-10}$. Wyniki zostały przedstawione w poniższej tabelce:

	Social PageRank
www.ted.com	0.2381373691295440
www.colourlovers.com	0.4343479235414989
www.behance.net	0.8686958470829979

Można zauważyć, że największy ranking ma strona behance.net która została dodana przez dwóch użytkowników i oznaczonych najpopularniejszymi tagami - 2 razy tagiem portfolio, użytym tylko dla tej strony, raz tagiem design, który użyty był 2 razy w powyższych danych i również raz tagiem inspiration, który jest najpopularniejszym tagiem, użytym w przykładzie aż 3 razy.

3.2.3 Implementacja

Algorytm pozwala na wcześniejsze wyliczenie rankingu dlatego został zaimplementowany jako osobny proces. Dodatkowo algorytm wymaga danych, które muszą być wyliczone i zapisane w bazie danych przed rozpoczęciem jego działania.

3.2.4 Baza danych

Poniżej znajduje się obrazek przedstawiający bazę danych ze zmianami wymaganymi dla sprawnego działania algorytmu. Dodatkowe tabele dodane zostały dla przyspieszenia generowania danych wejściowych. Dane w nich zawarte wyliczane są z danych wyliczane są z danych już istniejących w bazie danych.

TODO: OBRAZEK - BAZA DANYCH Z DODATKOWYMI TABELAMI I ZAZNACZONYMI UŻYWANYMI POLAMI W BAZIE DANYCH

Tabela DOKUMENT zawiera dodatkowo pole `soc_page_rank` służące do przechowywania wyników algorytmu. Tabele TAG_USR i TAG_DOC jak również pole w tabeli USERTAGDOC.how_much zawierają redundantne dane wykorzystywane do generowania macierzy. Pole USERTAGDOC.how_much zawiera informacje o ilości tagów użytych przez użytkownika do opisu dokumentu. TAG_USR zawiera relacje między użytkownikami i tagami, ilość annotowanych dokumentów przez tę parę znajduje się w polu how_much. Analogicznie TAG_DOC jest relacją między annotacjami a dokumentami z ilością ich wykorzystania.

3.2.5 Tworzenie i mnożenie macierzy

Algorytm wymaga w każdej iteracji wykorzystania sześć macierzy. Z powodu wielkości danych nie jesteśmy w stanie przechowywać ich wszystkich w pamięci. Dodatkowo biblioteka wykorzystana do mnożenia macierzy nakłada ograniczenia na ilość kolumn i wierszy w macierzach. Kolejnym problemem jest czas potrzebny na pobranie danych z bazy danych.

Z powodu tych ograniczeń macierze pobierane są do partiami do pamięci z bazy danych, przetwarzane i zapisywane są w struktury ułatwiające szybki do nich dostęp. Następnie, zapisywane są w plikach zawierające największe porcje danych mieszczące się jednocześnie w pamięci.

W każdej iteracji algorytm pobiera dane z dysku. Z tych danych tworzona jest macierz o maksymalnej możliwej ilości wierszy na które pozwala wykorzystana biblioteka. Każda z tych części jest osobno mnożona przez wektor. Wyniki następnie składane są w wektorze wynikowym, który przekazywany jest do kolejnego mnożenia macierzy

WYKRES POKAZUJĄCY PROCES TWORZENIA I MNOŻENIA CZĘŚCIOWEGO MACIERZY

Mimo, że wymagania pamięciowe sprawiają, że trzeba wykonać większą liczbę działań w czasie mnożenia, macierze generowane są dość rzadkie. Większość pól zawiera wartość 0, co przyspiesza mnożenie macierzy i wektorów.

TODO: ilość wykorzystanych plików przy prawdziwych danych (1 mln)

TODO: ilość mnożeń

TODO: wykorzystana biblioteka: `cern.colt`

3.2.6 Wyniki - częściowe: TODO

// wyniki dla danych 66 000 dokumentów - TODO - wyniki dla dużych danych

Strona <http://www.pythonchallenge.com/> jest jedną ze stron z największym wynikiem socialpagerank (0.00301). Dodane jest przez 785 różnych

użytkowników. Zostało użyte do tego 209 unikalnych tagów. Najpopularniejsze tagi, w kolejności od najczęściej użytego to python, programming, challage i puzzle.

Jedną ze stron o najniższym rankingu jest np: <http://djangosnippets.org/snippets/1314/>. Strona ta zawiera specyficzne rozszerzenie dla frameworku django. Można się spodziewać że nie będzie to popularna witryna. Została ona dodana przez jednego użytkownika i opisana siedmioma annotcjami.

Strony które uzyskały ranking 0 to strony które nie zostały opisane żadnymi tagami przez użytkowników.

3.2.7 Problemy : TODO

Potencjalne problemy zauważone na mniejszej ilości danych: algorytm jest podatny na cykle które mogą zostać stworzone przez dużą ilość wygenerowanych użytkowników.

3.2.8 Przykładowe wyniki wyszukiwarki : TODO

TODO: OPISAC POZNIEJ, PO ZAIMPLEMENTOWANIU WYSZUKIWANIA UZYWAJACEGO SOCIALPAGERANK

Rozdział 4

Adapted PageRank

4.1 Opis

Algorytm Adapted Page Rank jest zainspirowany algorytmem PageRank. Ideą za algorytm PageRank jest pomysł, że strona jest ważna jeśli dużo innych stron ma odnośniki wskazujące na tą stronę i te strony są również ważne.

Ponieważ algorytm page rank nie może być bezpośrednio zastosowany do zebranych danych autorzy algorytmu adapted page rank zmienili strukturę danych na nieskierowany graf trydzienny $G_f = (V, E)$.

Proces tworzenia grafu G_f :

1) zbiór wierzchołków V powstaje z sumy rozłącznej zbioru użytkowników, tagów, i dokumentów: $V = U \cup T \cup R$

2) wszystkie wystąpienia łącznie tagów, użytkowników, dokumentów stają się krawędziami grafu G_f : $E = \{\{u, t\}, \{t, r\}, \{r, u\} | (u, t, r) \in Y\}$. Wagi tych krawędzi są przydzielane w następujący sposób: każda krawędź $\{u, t\}$ ma wagę $|\{r \in R : (u, t, r) \in Y\}|$, czyli jest to ilość dokumentów, którym użytkownik u nadał anotację t . Analogicznie dla krawędzi $\{t, r\}$: $waga = |\{u \in U : \{(u, t, r) \in Y\}|$ i krawędzi $\{r, u\}$, gdzie $waga = |\{t \in T : \{(u, t, r) \in Y\}|$.

dane wejściowe:

A – prawo stochastyczna macierz sąsiedztwa grafu G_f

p – wektor preferencji

w – losowo zainicjalizowany wektor

α, β, γ – stałe, gdzie: $\alpha, \beta, \gamma \in [0, 1]$ i $\alpha + \beta + \gamma = 1$

do:

$w = \alpha * w + \beta * A * w + \gamma * p$

while: do czasu kiedy wartości wektora w zbiegają się.

wynikiem algorytmu jest wektor w .

Dane zostały uzyskane dla parametrów: $\alpha = 0.35, \beta = 0.65, \gamma = 0$ i pochodzą z pracy (cytat! / połączenie z bibliografią).

W algorytmie adapted page rank wektor preferencji i $p = 1$.

4.1.1 Algorytm FolkRank

Algorytm FolkRank jest wersją algorytmu adapted page rank, który daje różne wyniki w zależności od tematu. Temat tej jest ustalany w wektorze preferencji p . Wektor preferencji może być ustalony na wybrany zbiór tagów, użytkowników, dokumentów, albo na pojedynczy element. Wybrany temat będzie propagowany na reszcie dokumentów, tagów i użytkowników. Można dzięki temu, ustalając wagę na zapytanie albo konto użytkownika systemu uzyskać wyniki bardziej zbliżone do zainteresowań danego użytkownika.

Algorytm FolkRank może pomóc w analizie zbioru danych, albo w systemach nie działających w czasie rzeczywistym, ale nie jest użyteczny w zaprezentowanym systemie. Nie jest możliwy do wykorzystania z powodu długiego czasu obliczania wag dokumentów.

4.1.2 Przykładowe wyniki

Działanie algorytmu dla danych złożonych z 3 różnych dokumentów, 2 użytkowników i 3 tagów.

Strony www	użytkownicy	
	użytkownik 1	użytkownik 2
http://www.ted.com/ http://www.colourlovers.com/ http://www.behance.net/	inspiration design portfolio, design	inspiration portfolio, inspiration

macierz asocjacyjna powstała z powyższych danych ma wymiary 8×8 i wygląda:

$$G_f = \begin{pmatrix} 0 & 0 & 0 & 1 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 1 & 1 & 1 & 0 \\ 0 & 0 & 0 & 2 & 2 & 1 & 1 & 2 \\ 1 & 1 & 2 & 0 & 0 & 1 & 2 & 1 \\ 0 & 1 & 2 & 0 & 0 & 2 & 0 & 1 \\ 1 & 1 & 1 & 1 & 2 & 0 & 0 & 0 \\ 0 & 1 & 1 & 2 & 0 & 0 & 0 & 0 \\ 0 & 0 & 2 & 1 & 1 & 0 & 0 & 0 \end{pmatrix}$$

	Social PageRank
doc: http://www.ted.com/	0.280676409730572
doc: http://www.colourlovers.com/	0.369220441236174
doc: http://www.behance.net/	0.384551423972747
usr: użytkownik A	0.402473669662513
usr: użytkownik B	0.354578057974890
tag: inspiration	0.383291185401107
tag: design	0.321455285546253
tag: portfolio	0.314739469615726

Zbierczość wektora została uzyskana po 22 iteracjach.

Patrząc na powyższe wyniki najwyższy ranking wśród dokumentów ma strona begence.net: została ona dodana przez 2 użytkowników i przypisane jej zostały 4 tagi. Niewiele niższy ranking ma witryna colourlovers dodana przez 2 użytkowników i opisana 2 różnymi anotacjami. Można po tym wnioskować że nadanie większej ilości tagów nie ma dużego wpływu na rank strony. Za to zmniejszenie liczby użytkowników którzy tą stronę dodali, ma duże: przykład strona ted.com i colourlovers.com gdzie widoczny jest dość duży skok wartości wyniku.

4.2 Implementacja

Z powodu długiego czasu działania i dużej ilości wymaganych danych nie mogą być one pobierane bezpośrednio z bazy danych. Przed rozpoczęciem działania algorytmu są one pobierane, zapisywane w struktury pozwalające na lepsze i szybsze ich przeglądanie i serializowane do pliku. Dodatkowo żeby przyspieszyć tworzenie danych korzystamy z dodatkowych tabel zawierających już wyliczone dane np: o ilości dokumentów dodanych i opisanych tym samym tagiem przez użytkowników.

W czasie każdej iteracji algorytmu są one pobierane z pliku, zamieniane na macierze i poddawane dalszym operacjom. Po zakończeniu działania algorytmu wyniki zapisywane są w bazie danych. Z powodu tego, że dane wymagane w algorytmie adapted pagerank są zbliżone do danych wymaganych w algorytmie socialpagerank wykorzystywane są te same zserializowane dane. Dokładniej macierz wygląda następująco:

$$G_f = \begin{pmatrix} 0 & M_{DU} & M_{TD}^T \\ M_{DU}^T & 0 & M_{UT} \\ M_{TD} & M_{UT}^T & 0 \end{pmatrix}$$

gdzie tworzenie macierzy M_{UD} , M_{TD} , M_{UT} jest opisana w rozdziale opisującym algorytm social page rank.

Dodatkowo podobnie jak w algorytmie socialPageRank z powodu wielkości macierzy i ograniczeń pamięciowych są one mnożone częściowo przez wektor wag.

4.3 Wyniki

// TODO: po zebraniu wystarczającej ilości danych
// opis kilku przykładowych dokumentów

4.4 Porównanie algorytmu SocialPageRank i algorytmu Adapted PageRank : TODO

4.4.1 Porównanie na małym przykładzie

Szybkość działania:

Zbierność wektora uzyskano szybciej - bo już po pięciu iteracjach przy algorytmie SocialPageRank. W algorytmie Adapted Page Rank wymagało to aż 22 iteracji. Zbierność można przyspieszać przez zmianę parametru α w algorytmie page rank.

Wyniki:

Kolejność wag dla dokumentów jest taka sama w przypadku jednego i drugiego algorytmu, ale ich wartości są zdecydowanie inne. W przypadku dokumentu o największej randze: behence.net i kolejnego colourlovers.com różnica dla algorytmie adapted page rank jest niewielka proporcjonalnie do wagi, a przy social page rank waga behence.net jest prawie dwa razy większa od poprzednika. Widoczna podobieństwo jest dopiero dla dokumentów colourlovers.com i ted.com gdzie różnica wyników algorytmów dla jednego i dokumentu jest znaczna. Jest to prawdopodobnie spowodowane, że na wynik jednego i drugiego algorytmu wpływa mała ilość użytkowników, którzy dodali dany dokument i mała ilość tagów przypisanych temu dokumentowi.

Dodatkowe dane:

Adapted Page Rank daje nam dodatkową wiedzę w postaci rank dla tagów i użytkowników. O ile nie jest to przydatne przy wyszukiwaniu, ale daje dodatkowe informacje o posiadanych danych.

4.4.2 Porównanie w działającym systemie

// TODO: zrobić po zebraniu dużej ilości danych
// wykresy ??
// porównanie i opis dla kilku dokumentów

Rozdział 5

Opis aplikacji [TODO]

5.1 Folksonomia

TA DEFINICJA MA BYC W INNY MIEJSCU:

Folksonomia nazywamy krotkę $F := (U, T, R, Y)$, gdzie: U, T, R to zbiory skończone, których elementy składają się odpowiednio z użytkowników, tagów i dokumentów. Y jest relacją “przypisania tagu” pomiędzy tymi elementami $Y \subseteq U \times T \times R$

Użytkownicy i tagi identyfikowani są na podstawie ich unikalnych nazw własnych. Dokumenty mogą być różnymi danymi: stronami www, zdjęciami, plikami np: pliki pdf. Ta praca bazuje na danych pobranych z witryny delicious, które w zdecydowanej większości są stronami www. Dane które nie są stroną www nie są brane pod uwagę w tej pracy.

5.2 delicious.com : TODO

TODO: opis strony delicious

5.3 Architektura

5.3.1 Wątki

Główne zadania w serwisie wykonywane są przez odpowiednie wątki. Część z nich działa cały czas, część z nich reaguje na odpowiednie zmiany w bazie danych, a inne uruchamiane są na żądanie. Poniżej jest lista głównych zadań wykonywanych przez te procesy. Dokładne opisy znajdują się w kolejnych rozdziałach:

- RecentCrawler: przegląda strone delicious.com i zapisuje nowo dodane strony, użytkowników i tagi

- URLCrawler: sprawdza czy nie nastąpiła jakaś zmiana w dokumentach zapisanych już w naszej bazie danych np: czy nie zostały one dodane przez większą liczbę użytkowników
- UserCrawler: zadaniem tego crawlera jest sprawdzenie czy użytkownicy, o których informacje posiadamy nie dodali innych stron do swoich kolekcji dokumentów.
- SocialNetowrkCrawler: ten wątek nie pobiera danych z serwisu delicious, ale sprawdza czy URL'e zapisane w bazie danych zostały dodane przez użytkowników innych serwisów, takich jak: digg, twitter, facebook.
- Procesy odpowiedzialne za algorytmy: każdy z algorytmów (social pagerank, adapted pagerank) posiada własny wątek, który przygotowuje dla niego dane, a następnie uruchamia działanie samego algorytmu.
- cache: wątki te odpowiadają są za wyliczanie danych, dzięki czemu później jest do nich szybszy dostęp. Na przykład w tabeli DOCUMENT dodawane są również dodatkowe informacje, które są wypisywane w momencie kiedy wyszukiwarka zwróci wynik dla użytkownika. Wyliczanie tych informacji w trakcie wyszukiwania byłoby zbyt czasochłonne. Dane zapisywane tam to np: kilka najczęściej używanych tagów dla tego dokumentu wraz z ich ilością i ilość użytkowników którzy dodali ten dokument.

5.3.2 Dane

Dane pobierane są na kilka sposobów. Głównym źródłem nowych informacji jest kanał RSS strony delicious. Serwis delicious daje użytkownikom dostęp do kilku kanałów RSS. Kanał 'recent' zawierający najnowsze dokumenty, dodawane właśnie przez użytkowników. Dodatkowo swój kanał RSS zawierają też użytkownicy i adnotacje - znajdują się w nich ostatnio dodane przez danego użytkownika dokumenty, a w kanale RSS adnotacji: ostatnio opisane danym tagiem dane. Każdy dokument zawiera także własny kanał RSS, znajdują się w nim informacje o użytkownikach i tagach którym został opisany dany dokument.

Nowe dane

Kanał delicious zawierający najnowsze dokumenty wykorzystywany jest do pobierania informacji o ostatnio dodanych dokumentach. Kilka wątków w aplikacji sprawdza na bieżąco stronę. Dla każdego dokumentu znajdującego się w pobranych danych, sprawdzany jest jej kanał RSS i stamtąd pobierane są dane o użytkownikach i tagach którymi dany dokument został opisany. Następnie wszystkie te informacje zapisywane są w bazie danych.

Odświeżanie danych

W czasie działania wątków pobierających nowe dane działają również wątki sprawdzające czy nie zostały dodane nowe wpisy przez zapisanych już użytkowników, albo czy informacje posiadane o dokumencie nie zmieniły się: np dodany został również przez innego użytkownika i opisany innymi tagami. Wykorzystywane do tego są kanały RSS użytkowników i dokumentów.

Czyszczenie danych

Dane pobierane z serwisu delicious nie zawsze są w postaci wymaganej przez aplikację. Spowodowane to jest błędami użytkowników, albo specyficznym stylem zapisywania tagów przez nich.

Niektóre tagi mają różne znaczenie w zależności od kontekstu, na przykład tag 'design' ma inne znaczenie w kontekście strony o programowaniu, a inne w kontekście strony o sztuce. Część użytkowników żeby poradzić sobie z tym problemem dodają do tagów informacje mówiące o ich domenie. Często domena ma wygląd 'programming@design' czy 'art#design'. Żeby rozwiązać ten problem adnotację przed dodaniem do bazy danych są dzielone ze względu na najpopularniejsze znaki specjalne. Dodane kontekstu do tagu mogłoby być przydatne w aplikacji, ale z powodu tego że każdy użytkownik ma swój specyficzny sposób opisywania dokumentów np: design@art i art#design, trudno jest je zunifikować. Dodatkowym problemem jest to, że nie jest to sposób opisu używany przez wszystkich użytkowników.

Adnotacje przypisywane przez użytkowników często kończą się lub zaczynają od znaków specjalnych. Jest to spowodowane np: błędami (dodatkowe przecinki) albo specyficznym stylem opisywania danych przez użytkownika. Wszystkie znaki specjalne z końca i początku dokumentu są usuwane przed dodaniem do bazy danych.

Przykłady danych przed i po ich oczyszczeniu:

- '@java' : 'java'
- '@@java' : 'java'
- '#java6@' : 'java6'
- 'design!\$%@art' : 'design', 'art'
- 'art!#,' : 'art'

Dane z innych serwisów społecznościowych

W systemie wykorzystywane są również dane z innych serwisów społecznościowych: twitter, facebook i digg. Dostarczają one API, które pozwala sprawdzić ile użytkowników udostępniło dana stronę w tych serwisach. Dla

każdego z tych serwisów tworzony jest osobny wątek, który sprawdza czy w bazie danych nie ma nowo dodanych dokumentów. Dla tych dokumentów sprawdzany jest odpowiedni serwis, wynik zapisywany jest w bazie danych w tabeli DOCUMENT. Co pewien okres czasu przeprowadzane jest sprawdzanie wszystkich dokumentów i informacje są nadpisywane.

Przykładowe dane z bazy danych TODO

TODO

5.3.3 Dane testowe TODO

TODO - więcej/poprawić

Żeby zapewnić unikalność danych testowych dane te pochodzą ze strony arXiv.org. Strona ta zawiera publikację z dziedzin fizyki, metamatyki, biologii, ... Z powodu specyfiki tych dokumentów i raczej małej ich popularności w serwisach społecznościowych można stworzyć z nich unikalne dane testowe.

Jako tekst który jest przekazywany do frameworku lucene używany jest tekst publikacji wyciągnięty z pliku pdf.

Tagi: są sumą słów kluczowych podanych w publikacji, tytułu i nazw działów do których dany dokument został przypisany.

TODO - więcej/poprawić

Opis: TODO

TODO

Przykładowe dane testowe: TODO

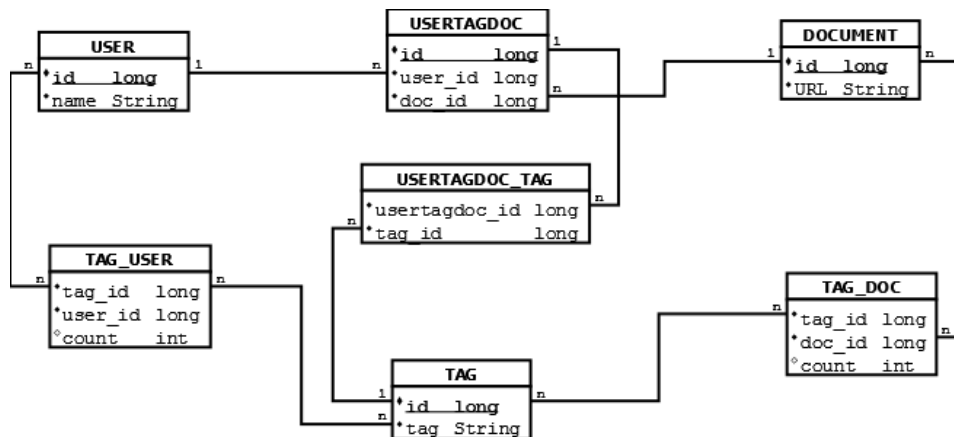
TODO

5.4 Baza danych

Jako serwer bazy używany jest MySQL 5.1. Komunikacja między aplikacją a bazą danych odbywa się za pomocą frameworku Hibernate. Framework ten zapewnia translację danych z relacyjnej bazy danych na obiekty używane w aplikacji.

Opis tabel:

- USER: tabela zawiera dane użytkowników, ich *id* i unikalną nazwę pobraną z serwisu delicious.
- DOCUMENT: tabela zawiera dane o dokumentach. *URL* który reprezentuje dokument jest unikalny. Ponieważ adresy stron po pominięciu



Rysunek 5.1: Baza danych

ostatniego slash/backslash prowadzi do tych samych witryn, przy sprawdzaniu unikalności linku brane pod uwagę są wszystkie kombinacje adresów.

- TAG: tabela zawierająca adnotacje stron.
- USERTAGDOC i USERTAGDOC_TAG: są to tabele które służą do zapisania w bazie danych relacji nadania k tagów: $t_n, t_{n+1}, \dots, t_{n+k}$ przez użytkownika $user_j$ dokumentowi doc_m . Dodatkowo wartość k : ilość przypisanych tagów, zapisana jest w polu *count* tabeli USERTAGDOC.
- TAG_USER i TAG_DOC: tabele zawierające redundantne dane przyspieszające wykonywanie innych operacji. TAG_USER zawiera informacje o ilości dokumentów opisanym tagiem tag_k przez użytkownika $user_n$. Druga tabela zawiera informacje o liczbie przypisań tagu tag_k do dokumentu doc_m . Tabele te są wykorzystywane głównie dla szybszego zbierania danych dla algorytmów Adapter PageRank i SocialPageRank.

5.5 Interface użytkownika TODO

TODO: więcej

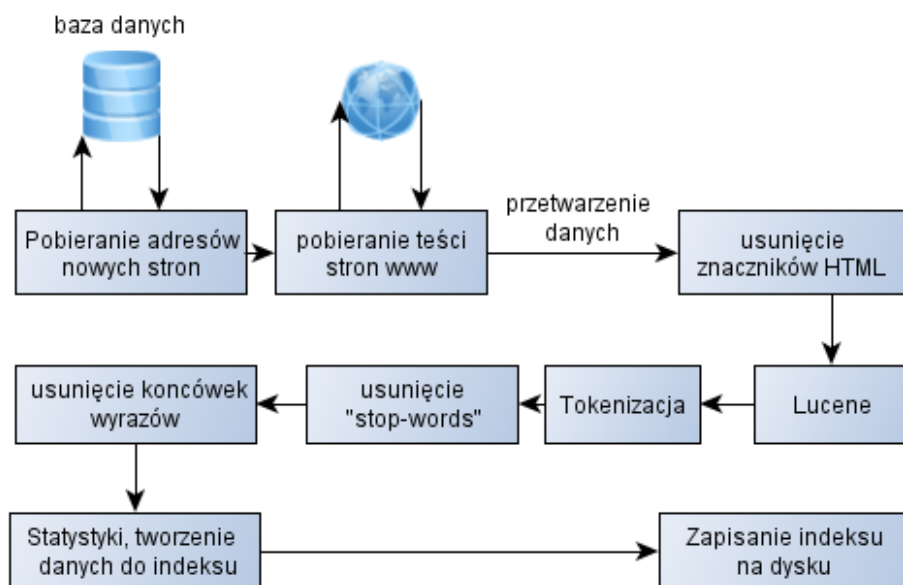
Interface użytkownika napisany jest przy pomocy technologii google-web-toolkit. Użytkownik po wczytaniu zapytania i naciśnięciu ENTER otrzymuje wyniki zapytania.

5.6 Lucene

Lucene jest biblioteką napisaną w Javie. Biblioteka ta jest w stanie indeksować dużą ilość dokumentów z różnych źródeł i przeprowadzać wyszukiwania w tych tekstach. W opisywanej aplikacji, framework Lucene przechowuje źródła stron, o których informacje zostały pobrane z serwisu delicious i zapisane w bazie danych.

5.6.1 Pobieranie stron

W pewnych odstępach czasu, wątek odpowiedzialny za indeksowanie stron sprawdza czy w bazie danych tabela DOCUMENT nie ma informacji o nowych wpisach. Z bazy danych pobrane są informacje o adresach tych stron. Następnie dla każdego adresu URL zostaje pobrana treść strony na którą wskazuje. Strona WWW następnie zostaje oczyszczona ze znaczników HTML, i przekazane do frameworku lucene do zaindeksowania. Jeśli wszystkie czynności zakończą się powodzeniem, w bazie danych zostaje odnotowana informacja o posiadaniu na dysku danego dokumentu. Rysunek 5.2 przedstawiony jest cały proces pobierania i przetwarzania danej strony.



Rysunek 5.2: Lucene: pobieranie danych i indeksowanie

Do Lucene zapisywane są następujące informacje: identyfikator *id* dokumentu z bazy danych, oraz przetworzony tekst strony WWW. Przechowywanie identyfikatora dokumentu w danych Lucene pozwala późniejsze powiązanie wyników wyszukiwania z odpowiednim rekordem w bazie danych.

W czasie indeksowania biblioteka Lucene wykonuje wiele czynności które pozwalają jej później szybko wyszukiwać informację. Główne z nich to:

- Tekst zostanie przetworzony na ciąg termów,
- usunięcie końcówek wyrazów,
- usunięcie 'stop-words' z tekstu, czyli słów nie mających dużego znaczenia przy wynikach wyszukiwania,
- obliczenie statystyk, np: wystąpienia słów w dokumencie, odległości od siebie

5.6.2 Wyszukiwanie

Czas wyszukiwania zapytania w dokumentach przechowywanych Lucene jest szybkie. Przy małej, poniżej 1GB danych, wyszukiwanie następuje praktycznie w czasie rzeczywistym. Zapytanie jest przekazywane do frameworku, w którym jest ono przekształcane na termy. Dla zapytania q i dla każdego dokumentu d wyliczana jest wartość funkcji $score(q, d)$. Wynikiem są dokumenty posortowane wg. wyniku tej funkcji.

$$score(q, d) = coord(q, d) * queryNorm(q) * \sum_{t \text{ in } q} \left(tf(t \text{ in } d) * idf(t)^2 * getBoost(t) * norm(t, d) \right)$$

gdzie:

- $coord(q, d)$: funkcja zwraca wartości zależne od miejsca występowania i odległości od siebie szukanych termów w dokumencie.
- $queryNorm(q)$: funkcja normalizująca wyniki zapytania
- $tf(t \text{ in } d)$: funkcja wyliczająca częstość występowania danego termu w dokumencie
- $idf(t)$: funkcja wyliczająca częstość występowania termu we wszystkich dokumentach.
- $getBoost(t)$ - Lucene pozwala na zwiększenie wagi niektórych termów. Nieużywane w tej aplikacji.

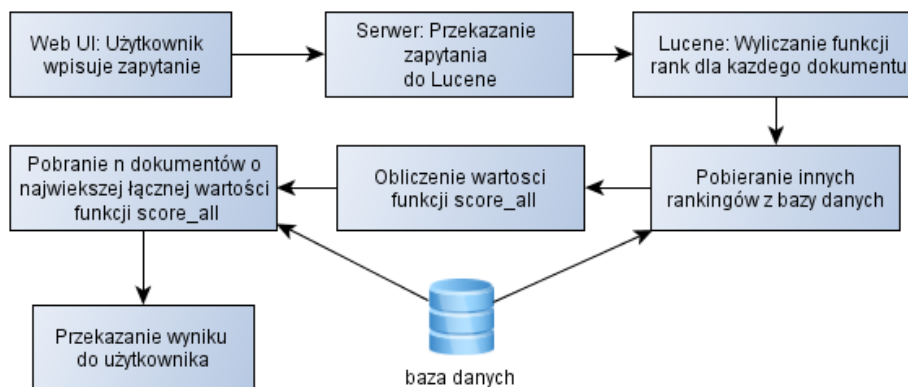
Lucene ocenia dokumenty głównie na podstawie funkcji TF-IDF. Każdy dokument reprezentowany jest przez wektor, składający się z wag słów występujących w tym dokumencie. TFIDF informuje o częstości wystąpienia termów uwzględniając jednocześnie wyważenie znaczenia lokalnego termu i jego znaczenia w kontekście pełnej kolekcji dokumentów.

5.7 Wyszukiwanie TODO

$$score_all(doc, q) = \alpha * lucene_score(doc, q) + \beta * social_page_rank_score(doc) + \gamma * adapt_page_rank(doc) + \eta * social_rank(doc)$$

gdzie:

- $\alpha, \beta, \gamma, \eta$ współczynniki wag dla poszczególnych algorytmów
- *lucene_score* - rank dokumentu *doc* dla zapytania *q*. Zwracane przez framework lucene.
- *social_page_rank* wynik algorytmu social page rank dla dokumentu *doc*
- *adapt_page_rank* wynik algorytmu adapted page rank
- *social_rank* wynik algorytmu social rank, działającego na danych pobranych z serwisu digg/twitter/facebook



Rysunek 5.3: Proces wyszukiwania dokumentów i wyliczania ich rank

Rozdział 6

Wyniki [TODO]

TODO

Bibliografia

- [1] Shenghua Bao, Guirong Xue, Xiaoyuan Wu, Yong Yu, Ben Fei, and Zhong Su. Optimizing web search using social annotations. In *Proc. WWW '07*, pages 501–510, Banff, Canada, 2007.
- [2] Anlei Dong, Ruiqiang Zhang, Pranam Kolari, Jing Bai, Fernando Diaz, Yi Chang, Zhaohui Zheng, and Hongyuan Zha. Time is of the essence: improving recency ranking using twitter data. In *Proceedings of the 19th international conference on World wide web, WWW '10*, pages 331–340, New York, NY, USA, 2010. ACM.
- [3] Andreas Hotho, Robert Jäschke, Christoph Schmitz, and Gerd Stumme. Information retrieval in folksonomies: Search and ranking. In York Sure and John Domingue, editors, *The Semantic Web: Research and Applications*, volume 4011 of *Lecture Notes in Computer Science*, pages 411–426, Berlin/Heidelberg, June 2006. Springer.
- [4] Yan'an Jin, Ruixuan Li, Kunmei Wen, Xiwu Gu, and Fei Xiao. Topic-based ranking in folksonomy via probabilistic model. *Artificial Intelligence Review*, pages 1–13, February 2011.
- [5] Makoto Kato, Hiroaki Ohshima, Satoshi Oyama, and Katsumi Tanaka. Can social tagging improve web image search? In James Bailey, David Maier, Klaus-Dieter Schewe, Bernhard Thalheim, and Xiaoyang Wang, editors, *Web Information Systems Engineering - WISE 2008*, volume 5175 of *Lecture Notes in Computer Science*, chapter 19, pages 235–249. Springer Berlin / Heidelberg, Berlin, Heidelberg, 2008.
- [6] Michael Noll and Christoph Meinel. Web search personalization via social bookmarking and tagging. pages 367–380. 2008.
- [7] Majdi Rawashdeh, Heung-Nam Kim, and Abdulmotaleb El-Saddik. Folksonomy-boosted social media search and ranking. In Francesco G. B. De Natale, Alberto Del Bimbo, Alan Hanjalic, B. S. Manjunath, and Shin'ichi Satoh, editors, *ICMR*, page 27. ACM, 2011.
- [8] Ralf Schenkel, Tom Crecelius, Mouna Kacimi, Sebastian Michel, Thomas Neumann, Josiane X. Parreira, and Gerhard Weikum. Efficient

- top-k querying over social-tagging networks. In *Proceedings of the 31st annual international ACM SIGIR conference on Research and development in information retrieval*, SIGIR '08, pages 523–530, New York, NY, USA, 2008. ACM.
- [9] Y. Yanbe, A. Jatowt, S. Nakamura, and K. Tanaka. Can social bookmarking enhance search in the web? In *JCDL '07: Proceedings of the 7th ACM/IEEE-CS joint conference on Digital libraries*, pages 107–116, New York, NY, USA, 2007. ACM.
- [10] Ding Zhou, Jiang Bian, Shuyi Zheng, Hongyuan Zha, and C. Lee Giles. Exploring social annotations for information retrieval. In *WWW '08: Proceeding of the 17th international conference on World Wide Web*, pages 715–724, New York, NY, USA, 2008. ACM.