

# Automatyczna klasyfikacja i ekstrakcja tematu krótkich notatek w języku polskim

Paweł Obrok  
pod kierunkiem dr. Michała Korzyckiego

13 sierpnia 2012

## Spis treści

<b>1</b>	<b>Wstęp</b>	<b>3</b>
<b>2</b>	<b>Podstawy teoretyczne</b>	<b>3</b>
<b>3</b>	<b>Procedura badawcza</b>	<b>3</b>
3.1	Dobór schematu wagowego . . . . .	3
<b>4</b>	<b>Opis danych</b>	<b>3</b>
4.1	Przykładowy problem . . . . .	3
<b>5</b>	<b>Wyniki i analiza</b>	<b>6</b>
5.1	Tematy . . . . .	6
5.2	Czas działania . . . . .	9
5.3	Metryki z nadzorem . . . . .	9
5.3.1	Ranking dokumentów . . . . .	9
5.3.2	Krzywe ROC . . . . .	11
5.3.3	Przywołanie i precyzja . . . . .	15
5.4	Metryki bez nadzoru (perplexity) . . . . .	16
5.5	Wnioski . . . . .	17
<b>6</b>	<b>Podsumowanie</b>	<b>18</b>

## 1 Wstęp

## 2 Podstawy teoretyczne

## 3 Procedura badawcza

### 3.1 Dobór schematu wagowego

Aby przekształcić dokument tekstowy w wektor, na którym mogą operować algorytmy macierzowe konieczne jest zdefiniowanie procedury tego przekształcenia. W tym przypadku zastosowano podejście bag-of-words, czyli dla każdego dokumentu obliczono jedynie liczbę wystąpień każdego wyrazu z pominięciem jego miejsca wystąpienia w tekście. Zastosowano następnie schemat wagowy log-entr, aby odzwierciedlić nie tylko wagę danego wyrazu w danym tekście, ale również jego ładunek znaczeniowy ogólnie — przykładowo słowa w rodzaju “w”, “na” czy “który” będą pojawiać się bardzo często niezależnie od tematu tekstu — stąd ich waga powinna być niska.

Obliczenie wagi danego słowa  $w_i$  odbywa się według wzoru 1.

$$w_i = 1 - \sum_{j=1}^n \frac{p_{ij} \log(p_{ij})}{\log(n)} \quad \text{Naprawić to} \quad (1)$$

Drobiazgowe omówienie różnych schematów wagowych znaleźć można w [1].

## 4 Opis danych

Działanie algorytmów LDA i LSI analizowane było na zbiorze około 50 tysięcy notatek Polskiej Agencji Prasowej. Notatki te to krótkie wiadomości tekstowe, z których większość dotyczy pojedynczego wydarzenia. Dotyczą one różnych dziedzin życia, co nadaje temu zbiorowi dodatkową różnorodność i pozwala oczekiwać, że osiągnięte wyniki będą miarodajne dla efektywności algorytmów w prawdziwych zastosowaniach.

### 4.1 Przykładowy problem

Aby przetestować działanie obu algorytmów przygotowane zostało przykładowe zapytanie do systemu wyszukiwania informacji dla zbioru notatek prasowych PAP. Problem polega na znalezieniu dokumentów podobnych do pojedynczej wybranej notatki na temat bliźniaczek syjamski. Została przygotowana modelowa odpowiedź systemu dla porównania z faktycznymi odpowiedziami.

Zapytanie:

\*\*\*\*\* #000424 \*\*\*\*\*

W sobotę polskie bliźniaczki syjamskie Weronika i Wiktor przylecą do Polski rejssem Newark-Kraków - poinformował PAP oddział LOT-u na nowojorskim lotnisku. Bliźniaczki syjamskie Weronika i Wiktor urodziły się 26 maja 1999 w szpitalu Akademii Medycznej w Lublinie. Od 16 sierpnia przebywają w Filadelfii. W listopadzie przeszły udaną operację rozdzielania. Pierwszy termin wypisania dziewczynek wyznaczono na 11 lutego. Został jednak przesunięty o tydzień z uwagi na gorączkę Weroniki.

Najbardziej podobne dokumenty wybrane ręcznie:

\*\*\*\*\* #000516 \*\*\*\*\*

Wiktoria i Weronika, siostry syjamskie rozdzielone w listopadzie w Filadelfii, przyleciały z mamą w sobotę rano do Krakowa. Na lotnisku w Balicach żonę i córki przywitał ojciec dziewczynek, Edward Paleń. Byli też trzej bracia dziewczynek. "Dziewczynki przyjechały w dobrym stanie, są zdrowe - powiedziała na lotnisku w Balicach mama rozdzielonych bliźniaczek pani Krystyna Paleń. Lekarze amerykańscy twierdzili, że gdyby dziewczynki miały po powrocie do kraju trafić do polskiego szpitala, to oni by je przytrzymali dłużej u siebie. Wypisali dziewczynki w takim stanie, że mogą iść do domu - powiedziała dziennikarzom. Dodała, że Wiktoria i Weronika będą znajdowały się pod opieką lekarza pediatry w Stalowej Woli.

\*\*\*\*\* #009928 \*\*\*\*\*

W klasztorze Ojców Kapucynów w Stalowej Woli (Podkarpatcie) ochrzczone zostały w niedzielę syjamskie bliźniaczki - Weronika i Wiktor Paleniówny, które urodziły się 26 maja ub. roku częściowo zrośnięte klatkami piersiowymi i brzuskami. Pomoc w rozdzieleniu dzieci zaoferował Szpital Dziecięcy w Filadelfii w USA. Rozdzielenia dokonano 3 listopada, po kilkunastogodzinnej operacji, w której udział wzięło 45 lekarzy.

\*\*\*\*\* #011189 \*\*\*\*\*

Powoli stabilizuje się stan zdrowia 9-miesięcznego Kamila, jednego z rozdzielonych w Krakowie braci syjamskich. Drugi z braci - Patryk - nadal jest w ciężkim stanie - poinformował w poniedziałek PAP opiekujący się braćmi docent Adam Bysiek. Stan zdrowia Kamila docent Bysiek określił jako żokujący nadzieję". Obaj bracia nadal przebywają na oddziale intensywnej

terapii Polsko-Amerykańskiego Instytutu Pediatrii Uniwersytetu Jagiellońskiego w Krakowie Prokocimiu. Bracia zostali rozdzieleni tydzień temu. O terminie operacji zdecydowało nagłe pogorszenie się stanu zdrowia jednego z nich. Początkowo operacja rozdzielenia była planowana na wrzesień. Na początku czerwca braciom wszczepiono ekspandery, mające za zadanie namnożyć tkankę skórną potrzebną po operacji. Operacja trwała 8 godzin, uczestniczyło w niej 14 lekarzy oraz zespół anestezjologów i pielęgniarek. Jej pierwszą część zajęło rozdzielenie braci, w drugiej lekarze zajęli się rekonstrukcją rozdzielonych narządów. Bracia byli zrośnięci powłokami piersiowo-brzusznymi, mieli wspólną przeponę, worek osierdziowy i wątrobę. Była to siódma operacja rozdzielania bliźniaków syjamskich dokonana w Instytucie w Prokocimiu.

\*\*\*\*\* #008779 \*\*\*\*\*

Bracia syjamscy Kamil i Patryk przeszli w poniedziałek pierwszy zabieg przygotowujący ich do operacji rozdzielania, planowanej za trzy miesiące w Polsko-Amerykańskim Instytucie Pediatrii UJ w Krakowie-Prokocimiu.

\*\*\*\*\* #005662 \*\*\*\*\*

W Polsko-Amerykańskim Instytucie Pediatrii w Krakowie-Prokocimiu zmarły siostry syjamskie, które urodziły się przed tygodniem w Wejherowie.

\*\*\*\*\* #000469 \*\*\*\*\*

Bliźniaczki syjamskie z Poznania - Małgosia i Dorota - przeszły już pierwsze badania w Polsko-Amerykańskim Instytucie Pediatrii UJ w Krakowie-Prokocimiu. Z przeprowadzonych badań wynika, że siostry mają wspólną wątrobę i przeponę oraz prawdopodobnie wspólne drogi żółciowe i serce- powiedział PAP opiekujący się bliźniaczkami dr Adam Bysiek z Instytutu. Siostry urodziły się 11 lutego w poznańskiej klinice św. Rodziny. Dziewczynki zrośnięte są brzuskami i klatkami piersiowymi.

\*\*\*\*\* #005855 \*\*\*\*\*

Liczba urodzeń bliźniąt syjamskich nie odbiega w Polsce od statystycznej normy - powiedział PAP prof. Jan Grochowski, dyrektor Polsko-Amerykańskiego Instytutu Pediatrii UJ, w którym przebywają dwie pary bliźniąt syjamskich.

\*\*\*\*\* #010677 \*\*\*\*\*

Lekarze z Polsko-Amerykańskiego Instytutu Pediatrii UJ w

Krakowie Prokocimiu rozdzielili 9-miesięcznych braci syjamskich Kamila i Patryka - poinformował PAP doc. Adam Bysiek, szef zespołu opiekującego się bliźniętami.

\*\*\*\*\* #007320 \*\*\*\*\*

Prof. Louis Gerald Keith, który w ubiegłym roku przeprowadził operację rozdzielenia polskich bliźniaczek syjamskich Weroniki i Wiktorii, został odznaczony przez prezydenta Aleksandra Kwaśniewskiego Krzyżem Oficerskim Zasługi RP.

\*\*\*\*\* #007872 \*\*\*\*\*

Komitet etyki szpitala w Palermo na Sycylii wydał zgodę na operację peruwiańskich bliźniaczek syjamskich, w wyniku której jedna z nich ma szansę na ocalenie kosztem życia drugiej - doniosła prasa włoska.

\*\*\*\*\* #012193 \*\*\*\*\*

Jeden z 9-miesięcznych braci syjamskich, rozdzielonych przez lekarzy w Krakowie pod koniec czerwca, zmarł w środę wieczorem z powodu niewydolności krążenia - poinformował PAP docent Adam Bysiek.

## 5 Wyniki i analiza

Niniejszy rozdział zawiera porównanie różnych aspektów działania algorytmów LDA i LSI. Na jego końcu znajdują się wnioski jakie można wyciągnąć z zebranych danych.

### 5.1 Tematy

Tabele 5.1 i 5.1 zawierają niektóre tematy wygenerowane przez algorytmy LSI i LDA skonfigurowane na 100 tematów (po dziesięć najbardziej znaczących słów w każdym temacie). Pojedynczy wiersz tabeli zawiera jeden temat - liczby przy tokenach oznaczają wagi poszczególnych słów w danym temacie.

Tematy uzyskane przy pomocy LDA wydają się bardziej odpowiadać postrzeganiu tekstu przez człowieka niż te wygenerowane przez LSI. Przykładowo temat numer 4 w tabeli 5.1 można interpretować jako „nie pogoda i finanse” — możliwość złożenia dwóch tematów postrzeganych przez człowieka w jeden, ale z przeciwnymi znakami powoduje powstawanie tego typu kombinacji. Tematy wygenerowane przez LDA bywają złożeniami dwóch różnych konceptów, jednak zawsze mają ten sam znak, jak na przykład temat numer 5 w tabeli 5.1, który wydaje się łączyć koncepty „muzeum” i „przestępstwo”.

Tablica 1: Tematy wyekstrahowane przez algorytm LSI

Lp.	Temat
1	0.269*" + 0.181*- + 0.171*być + 0.161*procent + 0.144*polski + 0.138*rok + 0.119*) + 0.118*złoty + 0.111*( + 0.102*a
2	-0.304*procent + -0.265*wzrósć + -0.254*punkt + -0.211*WIG + -0.192*wynieść + -0.191*spać + -0.180*złoty + -0.165*spółka + -0.158*akcja + 0.158"
3	0.482*RATIO + 0.265*mecz + 0.234*: + 0.187*pokonać + 0.182*mistrzostwo + 0.149*) + 0.149*turniej + -0.142*" + 0.119*piłkarski + 0.117*wygrać
4	-0.301*stopień + -0.250*temperatura + -0.240*maksymalny + -0.228*wiatr + -0.222*umiarkowany + -0.216*deszcz + -0.212*słaby + -0.208*opad + -0.181*południe + 0.164*złoty
5	-0.390*złoty + -0.305*grosz + -0.262*dolar + -0.246*euro + 0.210*punkt + -0.195*osiągać + -0.170*milion + 0.159*WIG + -0.147*umocnić + 0.142*procent
6	-0.355*spółka + -0.301*Akcyjna + 0.259*grosz + -0.223*milion + 0.215*zamknięcie + 0.185*euro + 0.180*osiągać + 0.170*punkt + 0.153*dolar + -0.148*bank
7	-0.435*procent + 0.300*spółka + -0.232*rok + 0.212*akcja + -0.191*proca + 0.190*Akcyjna + 0.148*giełda + -0.133*milion + 0.127*zmienić + 0.124*kurs
8	-0.227*RATIO + 0.192*sąd + -0.147*: + 0.147*( + -0.135*unia + -0.129*mecz + 0.127*policja + -0.125*spółka + 0.122*tysiąc + -0.122*AWS
9	0.313*( + 0.274*) + -0.258*RATIO + -0.165*mecz + -0.158*sąd + -0.144*: + 0.133*wyścig + 0.126*mistrzostwo + 0.120*spółka + 0.120*świat
10	-0.230*sąd + 0.220*europejski + -0.187*AWS + 0.154*unia + -0.148*procent + 0.143*UE + -0.119*wyborczy + -0.118*okręgowy + 0.114*polski + 0.111*milion

Tablica 2: Tematy wyekstrahowane przez algorytm LDA

Lp.	Temat
1	0.027*open + 0.026*powodzianin + 0.021*podlaski + 0.018*Słowenia + 0.017*cukrownia + 0.017*najstarszy + 0.013*przedstawiony + 0.012*urodziny + 0.012*rata + 0.012*zrezygnować
2	0.021*europejski + 0.021*unia + 0.018*UE + 0.012*polski + 0.011*kraj + 0.010*" + 0.009*Litwa + 0.009*unijny + 0.008*państwo + 0.008*NATO
3	0.032*palestyński + 0.031*Izrael + 0.030*izraelski + 0.023*Palestyńczyk + 0.015*Arafat + 0.013*szaron + 0.013*świętokrzyski + 0.012*zawieszenie + 0.012*autonomia + 0.011*arabski
4	0.024*sąd + 0.015*aresztować + 0.015*podejrzany + 0.014*rejonowy + 0.013*okręgowy + 0.013*śledczy + 0.013*akt + 0.012*Gdynia + 0.012*oskarżenie + 0.012*Radom
5	0.016*wierzyciel + 0.013*muzeum + 0.013*wystawa + 0.011*zbiór + 0.011*śląski + 0.011*Brazylijczyk + 0.010*łączy + 0.010*zajmujący + 0.009*przestępczy + 0.009*łódzki
6	0.032*festiwal + 0.022*woj + 0.017*Białystok + 0.017*letni + 0.014*wielkopolski + 0.014*wrzesień + 0.014*kupno + 0.012*ogólnopolski + 0.012*usowanie + 0.012* impreza
7	0.019*siatkarz + 0.011*obniżka + 0.009*noc + 0.007*postać + 0.006*Gorzów + 0.006*artystyczny + 0.006*bóg + 0.006*bandyta + 0.005*nieznany + 0.005*ZSRR
8	0.012*" + 0.011*general + 0.010*motors + 0.008*Jedwabne + 0.008* kardynał + 0.007*film + 0.007*weekend + 0.007*Józef + 0.007*rocznica + 0.006*odbyć
9	0.017*świat + 0.016*klasa + 0.016*TP + 0.015*metr + 0.015*( + 0.015*) + 0.014*mistrzostwo + 0.014*zająć + 0.013*AZS + 0.013*bieg
10	0.044*obligacja + 0.021*Artur + 0.018*włosek + 0.017*pomnik + 0.016*politechnika + 0.016*białostocki + 0.016*spółeczność + 0.013*wyeliminować + 0.012*skorzystać + 0.011*wyemitować



## 5.2 Czas działania

Czy to ma sens?

## 5.3 Metryki z nadzorem

W tym rozdziale omówiono wyniki otrzymane za pomocą algorytmów LDA i LSI dla przykładowego problemu opisanego w 4.1. Należy zauważyć, że tego rodzaju ewaluacja wymaga ręcznego przygotowania danych testowych przez człowieka, co może być niepraktyczne dla dużych zbiorów danych. Jej zaletą jest fakt, że mierzy ona faktyczne osiągi danego rozwiązania w rzeczywistych problemach.

### 5.3.1 Ranking dokumentów

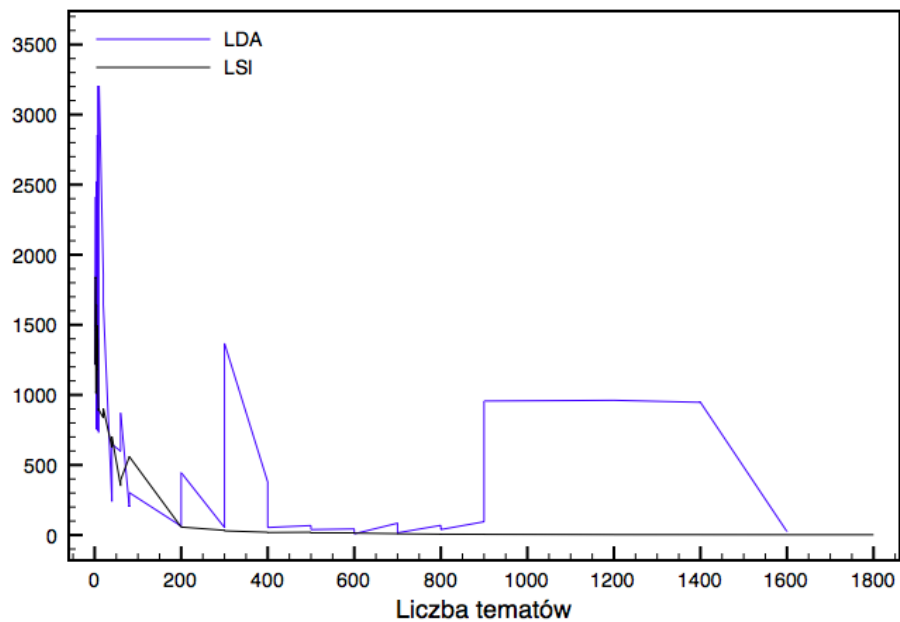
Wykresy 1 i 2 przedstawiają sumę kwadratów ranków dokumentów z wzorca przygotowanego ręcznie dla danego zapytania w wynikach działania odpowiednio algorytmów LDA i LSI dla różnej liczby tematów.

Algorytm LDA osiąga ogólnie gorsze wyniki niż LSI - poza przedziałem 50 – 100 tematów. Gorszy jest też (aczkolwiek niewiele) najlepszy wynik jaki udałooby się osiągnąć odpowiednio dobierając liczbę tematów. Na wykresie daje się także zauważyć stochastyczna natura LDA - podczas gdy dla LSI wyniki niemal monotonicznie poprawiają się wraz ze wzrostem liczby tematów dla LDA zdarza się znaczne pogorszenie wyników przy zwiększeniu tej liczby.

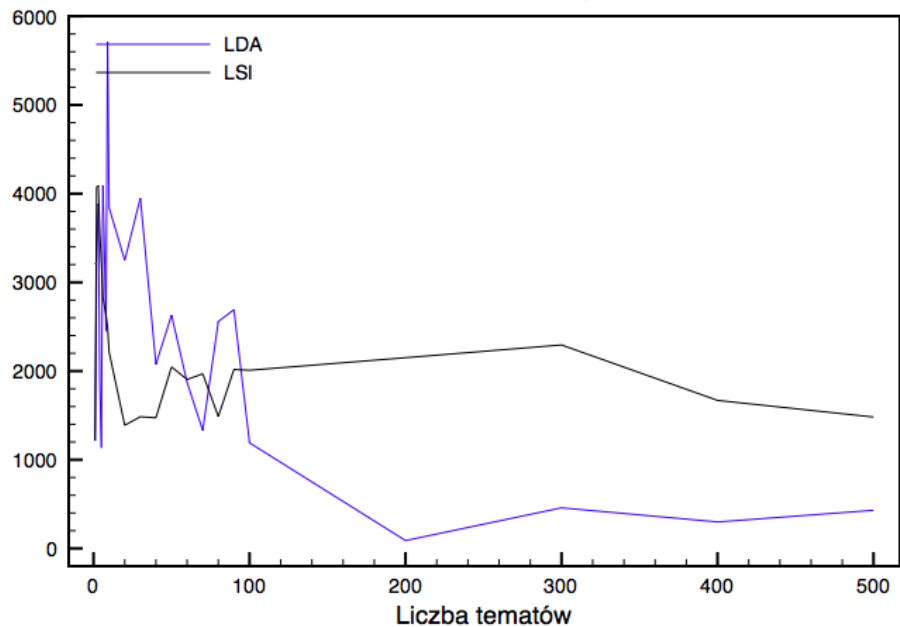
Polepszenie wyników dzięki zastosowaniu stemmingu jest widoczne na pierwszy rzut oka — polski jako język silnie fleksyjny jest znakomitą kandydatem do zastosowania tego typu techniki. W [2] zasugerowano, że ze stemmingu można zrezygnować dysponując odpowiednio dużym zbiorem danych jednak wyniki te uzyskano dla języka angielskiego, którego fleksja jest znacznie mniej rozbudowana. W tym wypadku zebranie tak dużej ilości danych może być mniej praktyczne niż skonstruowanie słownika fleksyjnego takiego jak na przykład ten opisany w [4].

Co ciekawe algorytm LDA radzi sobie znacznie lepiej od LSI bez wykorzystania stemmingu. Może to być spowodowane trudnością w przypadku LSI połączenia ze sobą słów, które różnią się formą fleksyjną i są w tym wypadku traktowane całkowicie osobno.

Rysunek 1: Suma kwadratów ranków dokumentów ze wzorca dla testowego zapytania (z wykorzystaniem stemmingu)



Rysunek 2: Suma kwadratów ranków dokumentów ze wzorca dla testowego zapytania (bez wykorzystania stemmingu)



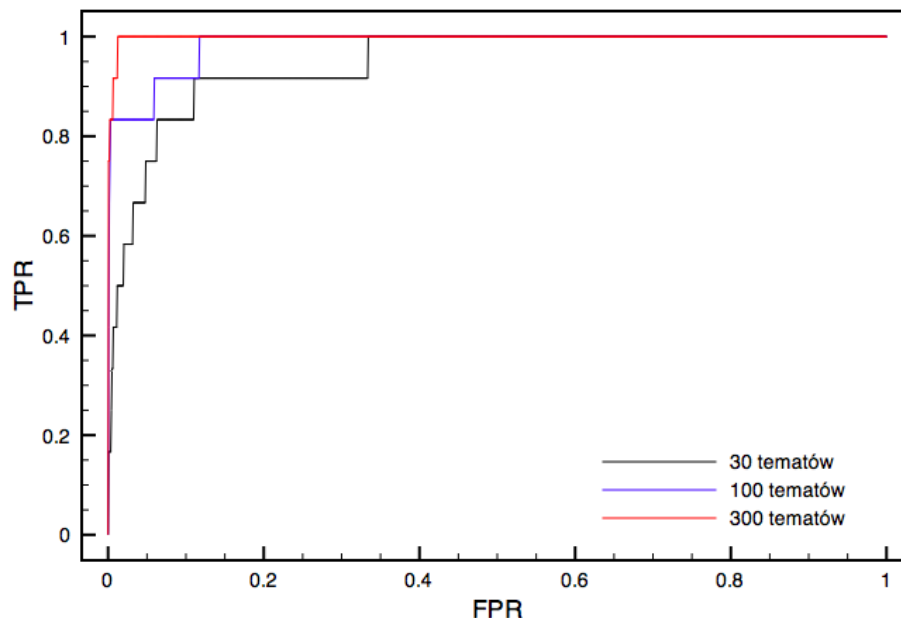
### 5.3.2 Krzywe ROC

Krzywa ROC [3] (Receiver Operation Characteristic) to wykres przedstawiający dla danego klasyfikatora zależność między stosunkiem liczby znalezionych dokumentów relewantnych do liczby wszystkich zwróconych dokumentów (TPR — True Positive Rate), a stosunkiem liczby odrzuconych dokumentów relewantnych do liczby wszystkich odrzuconych dokumentów (FPR - False Positive Rate) w miarę zmiany progu detekcji. W tym wypadku ten zmienny próg to po prostu liczba  $n$  - pierwszych  $n$  dokumentów jest traktowane jako odnalezione, a pozostałe jako odrzucone.

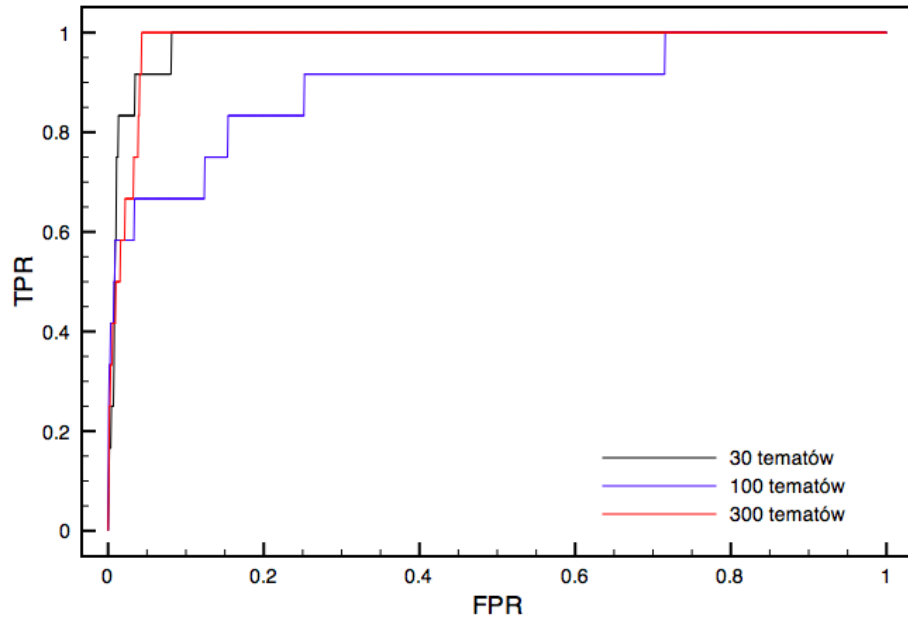
Lepsze klasyfikatory charakteryzują się krzywymi ROC położonymi dalej od linii  $x = y$ . Klasyfikatory blisko, lub na tej linii nie wykonują żadnej użytecznej pracy. Analiza odległości krzywej ROC od linii  $x = y$  w różnych miejscach wykresu może dać wskazówkę co do najlepszego dobrania progu detekcji dla danego problemu.

Wykresy 3 i 4 przedstawiają krzywe ROC dla algorytmów LDA i LSI dla różnych liczb tematów. Dla dużych liczb tematów algorytm LDA spisyje się gorzej, jednak można zauważyć, że klasyfikator uzyskany dla 30 tematów jest podobnej jakości lub lepszy jak ten uzyskany przy użyciu LSI dla 100 tematów.

Rysunek 3: Krzywe ROC dla algorytmu LSI dla wybranych liczb tematów

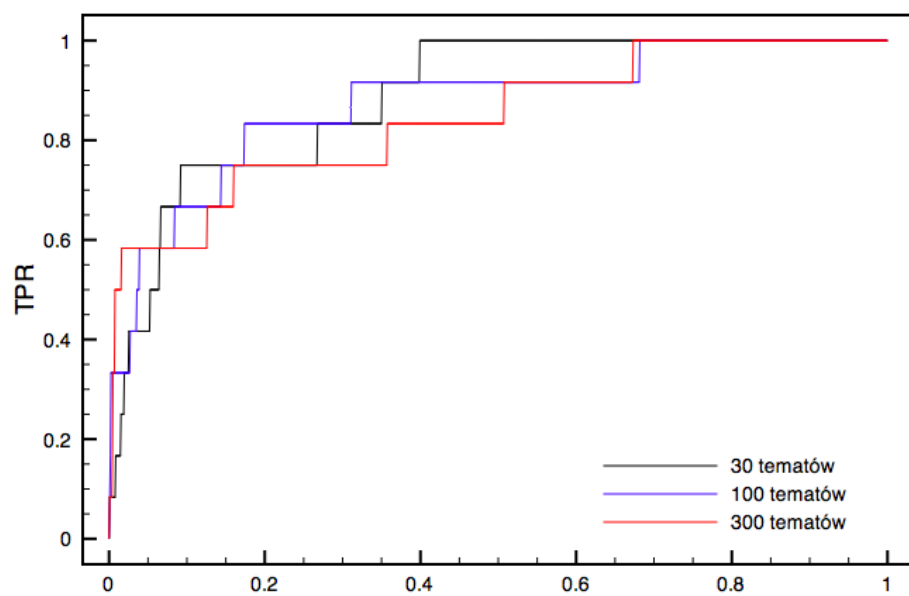


Rysunek 4: Krzywe ROC dla algorytmu LDA dla wybranych liczb tematów

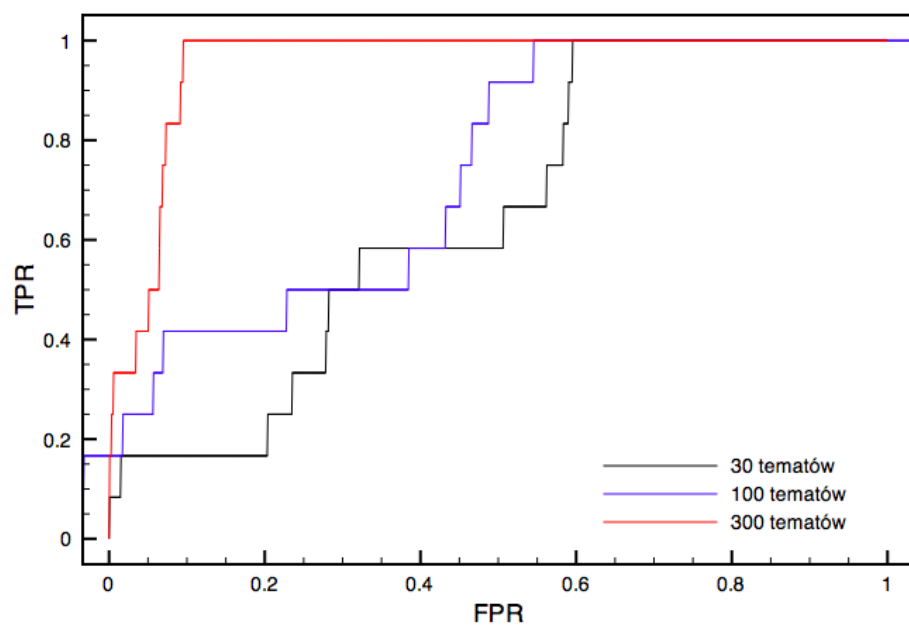


Na wykresach 5 i 6 przedstawione zostały krzywe ROC dla algorytmów LSI i LDA bez wykorzystania stemmingu. Ponownie daje się zauważyć lepsze działanie algorytmu LDA w tym wypadku - klasyfikator uzyskany dla 300 tematów jest znacznie lepszy od tego uzyskanego przy pomocy LSI.

Rysunek 5: Krzywe ROC dla algorytmu LSI dla wybranych liczb tematów bez wykorzystania stemmingu



Rysunek 6: Krzywe ROC dla algorytmu LDA dla wybranych liczb tematów bez wykorzystania stemmingu



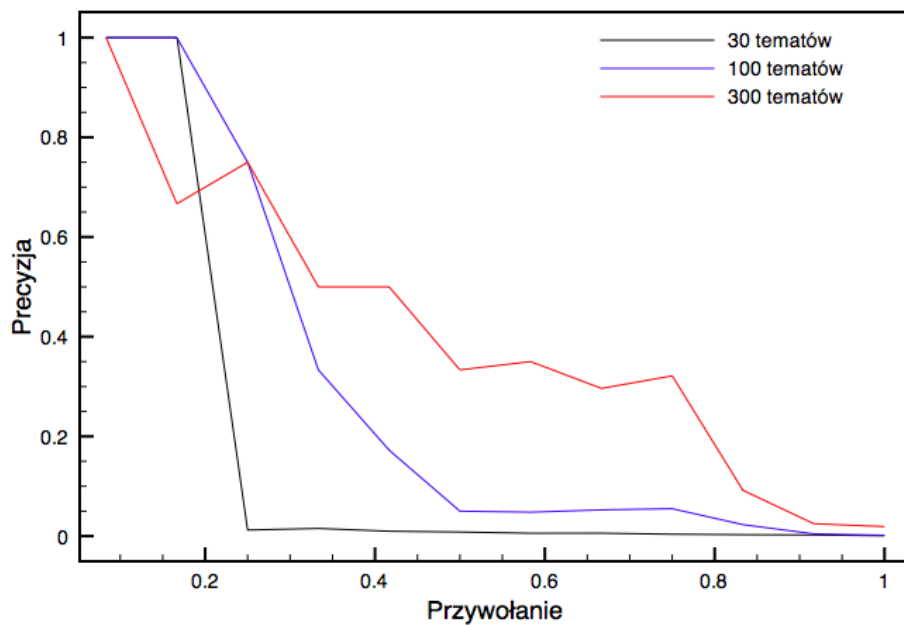
### 5.3.3 Przywołanie i precyzja

Przywołanie (stosunek liczby zwróconych relewantnych dokumentów do liczby wszystkich relewantnych dokumentów) i precyzja (stosunek liczby zwróconych relewantnych dokumentów do liczby wszystkich zwróconych dokumentów) to częste metryki w zadaniach typu information retrieval. Wybranie jakiegoś poziomu przywołania reprezentuje pewien kompromis między kompletnością zwróconych danych, a częstością występowania w nich danych relewantnych, a więc ilością czasu, które musi poświęcić operator systemu na ich dalsze przetworzenie.

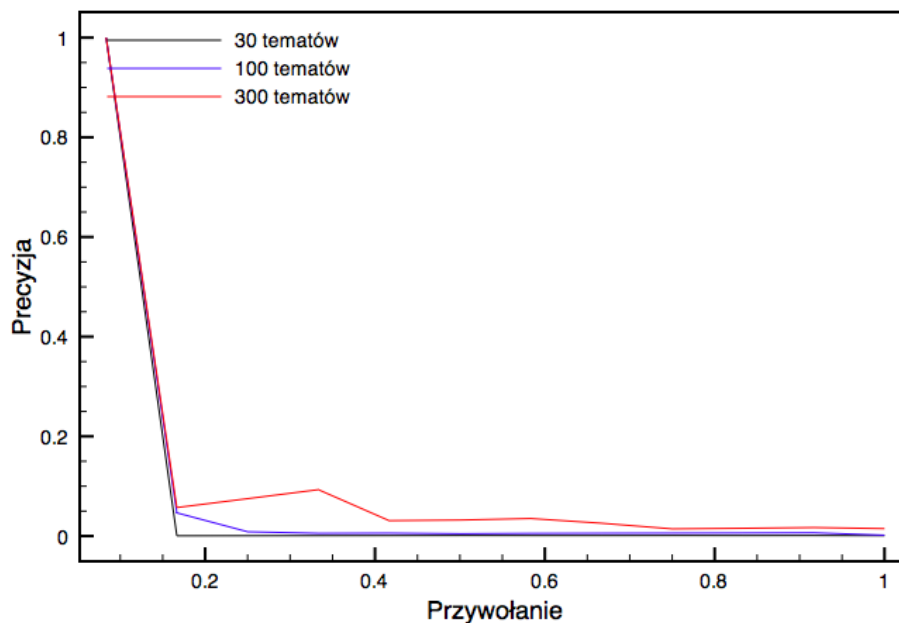
Wykresy 7 i 8 prezentują precyzję osiąganą przez algorytmy LDA i LSI na różnych poziomach przywołania dla przykładowego problemu.

Można zauważyć, że LDA daje znacznie gorszą precyzję niż LSI. Nawet najlepiej dobrana liczba tematów (w tym wypadku 300) pozwala osiągać precyzję porównywalną jedynie z LSI dla 30 tematów.

Rysunek 7: Precyzja na różnych poziomach przywołania dla algorytmu LSI



Rysunek 8: Precyzja na różnych poziomach przywołania dla algorytmu LDA



#### 5.4 Metryki bez nadzoru (perplexity)

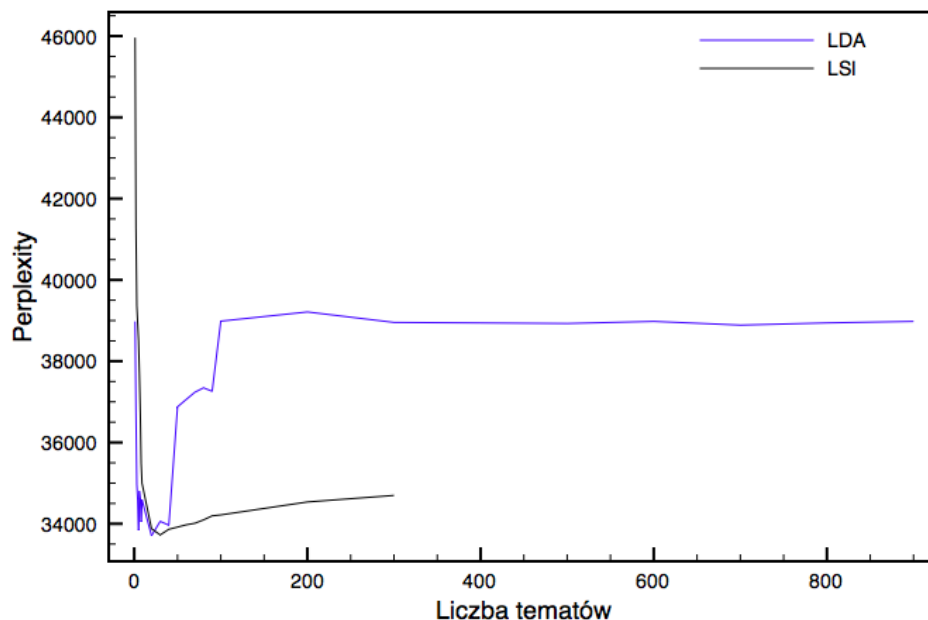
Współczynnik perplexity, którego wartości w zależności od liczby tematów są przedstawione na wykresie 9, daje pewne pojęcie o tym jak dobrze model jest w stanie przewidzieć nowe dane. Wysokie wartości współczynnika mogą wskazywać, że model jest przeuczony i będzie słabo uogólniał swoje działania na nieznane dane. Jest on dobrym wskaźnikiem jak dobrze dany model będzie sobie radził z klastrowaniem danego zbioru danych.

Wykres demonstruje, że optymalne wartości współczynnika perplexity zostają osiągnięte w okolicach 50 tematów. W tym wypadku może to sugerować, że mniej więcej na tyle właśnie grup tematycznych należałoby podzielić ten zbiór danych.

Algorytm LDA zachowuje niski współczynnik perplexity tylko stosunkowo blisko optymalnej liczby tematów. Takie zachowanie może wymagać dokładnego strojenia algorytmu do każdego zastosowania, co bywa uciążliwe i czasochłonne.



Rysunek 9: Współczynnik perplexity dla LDA i LSI w zależności od liczby tematów



## 5.5 Wnioski

Algorytm LDA daje gorszej jakości (a przynajmniej mniej stabilne) wyniki dla typowych problemów klasyfikacji i wyszukiwania informacji spotykanych w codziennej praktyce. Wydaje się za to być w stanie działać w sytuacji, gdy wiele różnych tokenów oznacza to samo (przypadek bez wykorzystania stemmingu), w odróżnieniu od LSI, którego wyniki są wtedy całkowicie nieprzydatne. To bardzo porządana cecha w przypadku braku odpowiedniego słownika fleksyjnego dla danego języka.

Przewagą LDA wydaje się być jakość generowanych tematów — przez wymuszenie dodatnich wag otrzymujemy na najbardziej znaczących pozycjach (z najwyższymi wagami) słowa opisujące dany dokument/temat, podczas gdy w przypadku LSI mogą to być słowa najodleglejsze. Takie zachowanie może okazać się korzystne w zastosowaniach typu tagowanie dokumentów lub automatyczne generowanie podsumowań czy słów kluczowych.

## 6 Podsumowanie

### Literatura

- [1] A. Figiel. Tekst jak wzorzec informacyjny — automatyczna ocena podobieństwa tematycznego tekstów za pomocą Latent Semantic Analysis. *Słowniki Komputerowe i Automatyczna Ekstrakcja Informacji z Tekstu*, 2009.
- [2] C. D. Manning and H. Schütze. *Foundations of Statistical Natural Language Processing*. MIT Press, Cambridge, MA, 1999.
- [3] D. K. McClish. Analyzing a portion of the ROC curve.
- [4] P. Pisarek. Słownik fleksyjny. *Słowniki Komputerowe i Automatyczna Ekstrakcja Informacji z Tekstu*, 2009.