

Streszczenie rozprawy doktorskiej:  
„Automatyczna klasyfikacja dokumentów tekstowych  
w zarządzaniu aktami spraw na użytek e-administracji”

mgr Marek Gajewski

### Sformułowanie zadania

Inspiracją podejmowanego w rozprawie zadania jest zagadnienie zarządzania *aktami spraw* w ramach obiegu dokumentów. W prawodawstwie polskim wymagane jest, aby dokumentacja gromadzona przez jednostki administracji publicznej i sposób jej organizacji odzwierciedlały przebieg załatwiania *sprawy*. Organizacja i klasyfikacja dokumentów prowadzona jest na podstawie systematyzacji rzeczowej *akt spraw* znanej pod pojęciem *Jednolitego Rzeczowego Wykazu Akt* (JRWA) i stanowiącej hierarchię klas, do których powinny być zaklasyfikowane wszystkie przetwarzane dokumenty. W ogólności oznacza to, że dla każdego nowego dokumentu należy rozpoznać jego przynależność do klasy z JRWA oraz w ramach tejże klasy przypisać go do odpowiednich *akt sprawy*. Tak sformułowany problem decyzyjny stanowi inspirację dla zadania postawionego w niniejszej rozprawie – automatycznego przypisywania dokumentu do klasy i pewnej sekwencji dokumentów (*sprawy*) w ramach tej klasy.

Zadanie postawione w niniejszej pracy odnosi się więc do problemu dwóch rodzajów klasyfikacji dokumentu, rozpatrywanych łącznie: do tradycyjnie rozumianej klasy (kategorii) dokumentów (w praktycznym scenariuszu stanowi ją klasa JRWA) i do konkretnej sprawy, wyodrębnionej w ramach tej klasy. Z jednej strony należy zaklasyfikować dokument zgodnie z jego *tematyką*, z drugiej zaś strony, zgodnie z, szeroko rozumianą, *sprawą*, której on dotyczy. Pierwsze z zadań ma bogatą literaturę, liczne zaproponowane w niej sposoby rozwiązania, i często jest rozważane pod nazwą *kategoryzacja dokumentów tekstowych* (ang. *text categorization*). Drugi z rozważanych rodzajów klasyfikacji jest bardziej złożony i wymaga przypisania do tej samej klasy (*sprawy*) dokumentów, które są ze sobą powiązane tym, że dotyczą tego samego zdarzenia, zamówienia, osoby itp. To powiązanie ma charakter wykraczający poza odniesienie się wyłącznie do tematyki dokumentów. Co więcej, zbiór klas jest w tym wypadku dynamiczny: dokument może należeć do istniejącej już klasy (*sprawy*) lub może stanowić początek nowej sprawy i wymaga utworzenia nowej klasy. Kolejną trudność stanowi fakt, że sprawy (klasy) obejmują przeciętnie niezbyt wiele dokumentów, a więc brak jest dostatecznej liczby dokumentów uczących dla skutecznego zastosowania znanych algorytmów klasyfikacji. Ten dwójaki charakter klasyfikacji stanowi o oryginalności postawionego zadania.

**Formalnie, zadanie postawione w rozprawie można przedstawić następująco.** Przyjmijmy w tym celu następującą terminologię i notację:

$D = \{d_1, \dots, d_n\}$  - korpus dokumentów, z których każdy przypisany jest do jednej z klas (kategorii) i występuje w jednej ze spraw (sekwencji dokumentów) związanych z tą klasą,

$d^*$  - dokument zadany na wejściu, który należy przypisać do jednej z kategorii i jednej ze spraw należących do tej kategorii,

$C = \{c_1, \dots, c_m\}$  - zbiór rozważanych klas (kategorii) dokumentów,

$\Sigma = \{\sigma_1, \dots, \sigma_p\}$  - zbiór spraw (sekwencji dokumentów), każda należąca do jednej z klas  $c_i$ ,

$\sigma_k = \langle d_{s_1}, d_{s_2}, \dots, d_{s_l} \rangle$ , sprawa (sekwencja dokumentów), gdzie wszystkie dokumenty  $d_i$  należą do jednej klasy  $c_j$ ; wyróżniamy sprawy „w toku”, do których należy klasyfikować nowy dokument oraz sprawy już „zamknięte”.

Ostatecznym celem w rozważanym zadaniu jest zaklasyfikowanie dokumentu  $d^*$  do jednej ze spraw „w toku”  $\sigma_k$  lub do nowej sprawy, a przez to również zaklasyfikowanie do jednej z kategorii  $c_j$ .

## Metoda rozwiązania

### Reprezentacja danych

#### Dokument

W zastosowanym w niniejszej pracy podejściu, jako punkt wyjścia do opracowania reprezentacji dokumentów przyjęto klasyczny model wektorowy. Każdy z dokumentów  $d \in D$  jest reprezentowany za pomocą wektora  $[w_1, \dots, w_M]$  w przestrzeni  $T = \{t_i\}_{i \in \{1, \dots, M\}}$ , gdzie  $T$  jest zbiorem wszystkich słów kluczowych używanych do indeksowania dokumentów w kolekcji, przy czym ostatecznie używa się tylko  $K$  słów kluczowych, odpowiadających współrzędnym o największych wartościach w tym wektorze. Wartość parametru  $K$  przyjmuje się dostatecznie małą, na przykład równą 5, w celu uzyskania bardziej zwartej reprezentacji dokumentów.

Formalnie reprezentacja dokumentu jest następująca:

$$d \rightarrow [w_1^d, \dots, w_M^d] \rightarrow [w_{k_1}^d, \dots, w_{k_K}^d] \quad (1)$$

gdzie  $w_i^d$  oznacza stopień ważności słowa kluczowego  $t_i$  dla reprezentacji dokumentu  $d \in D$ , zaś  $k_1, \dots, k_K$  są indeksami takimi, że:

$$w_{k_1}^d \geq w_{k_2}^d \geq \dots \geq w_{k_K}^d \geq w_{k_{K+1}}^d \geq \dots \geq w_{k_M}^d$$

#### Sprawa (sekwencja dokumentów)

Sekwencje dokumentów  $\sigma = \langle d_{s_1}, \dots, d_{s_l} \rangle$  tworzących sprawy reprezentowane są przez połączenie wektorów reprezentujących kolejne składowe dokumenty. Formalnie:

$$\sigma \rightarrow [w_{k_1}^{d_{s_1}}, \dots, w_{k_K}^{d_{s_1}}, \dots, w_{k_1}^{d_{s_l}}, \dots, w_{k_K}^{d_{s_l}}] \quad (2)$$

Przy obliczaniu dopasowania dokumentu  $d^* \rightarrow [w_{k_1}^{d^*}, \dots, w_{k_K}^{d^*}]$  do sekwencji (2), jej reprezentacja ulega zmianie i przyjmuje następującą postać:

$$\sigma \rightarrow [w_{k_1}^\sigma, \dots, w_{k_K}^\sigma] \quad (3)$$

gdzie  $k_i$  oznacza indeks  $i$ -tego słowa kluczowego,  $t_{k_i}$ , uwzględnionego w reprezentacji dokumentu  $d^*$ ,  $i \in [1, K]$ , zaś jego stopień ważności  $w_{k_i}^\sigma$  w reprezentacji sekwencji (2) jest obliczany następująco:

$$w_{k_i} = \begin{cases} 0 & \text{jeśli } t_{k_i} \text{ nie występuje w reprezentacji żadnego z dokumentów } \sigma \\ w_{k_j}^{d_{s_m}} & \text{w przeciwnym przypadku} \end{cases} \quad (4)$$

gdzie  $k_j = k_i$  oraz  $m \in [1, l]$  jest największym indeksem dokumentu w sekwencji  $\sigma$  (indeksem dokumentu występującego w sekwencji najpóźniej), w którego reprezentacji występuje słowo  $k_i$ .

#### Kategoria

Reprezentacja kategorii  $c \in C$  przybiera postać średniej wszystkich wektorów reprezentujących poszczególne dokumenty aktualnie przypisane do kategorii  $c$ . Formalnie:

$$c \rightarrow [w_1^c, \dots, w_{M_c}^c] \quad (5)$$

gdzie

$$w_{t_i}^c = \frac{\sum_{d \in c} w_{t_i}^d}{|\{d: d \in c\}|} \quad (6)$$

gdzie  $d \in c$  oznacza, że dokument  $d$  przypisany jest do kategorii  $c$ , a  $|\cdot|$  oznacza licznosc zbioru. Reprezentacja dokumentow uzywana w (6) jest okreslona zgodnie z (1). Tak wiec, slowo kluczowe  $t_i$  pojawi sie w reprezentacji kategorii  $c$  o ile tylko wystepuje ono w reprezentacji przynajmniej jednego z dokumentow przypisanych do tej kategorii. We wzorze (6) brane sa pod uwage wszystkie dokumenty znajdujace sie aktualnie w kolekcji.

### Klasyfikacja

Klasyfikacja dokumentu odbywa sie na podstawie oceny jego podobienstwa do poszczegolnych spraw i kategorii z nimi zwiazanych. Dla kazdego nowego dokumentu  $d^*$  oraz kazdej sprawy „w toku”  $\sigma$  obliczane sa dwa stopnie podobienstwa, odpowiednio:

$$\text{sim}_{seq}(d^*, \sigma) \quad (7)$$

$$\text{sim}_{cat}(d^*, c) \quad (8)$$

przy czym pierwszy z nich dotyczy danej sprawy  $\sigma$ , a drugi dotyczy kategorii, do ktorej sprawa ta nalezy. Nastepnie, laczne podobienstwo dokumentu  $d^*$  do sprawy  $\sigma$  obliczane jest jako suma wazona podobienstw (7) – (8):

$$w_{seq} * \text{sim}_{seq}(d^*, \sigma) + w_{cat} * \text{sim}_{cat}(d^*, c) \quad (9)$$

Dokument przypisywany jest do tej sprawy, dla ktorej wartosc podobienstwa wyliczona z uzyciem wzoru (9) jest najwieksza. Ten ogolny schemat algorytmu klasyfikacji jest realizowany poprzez dobór:

- 1) odpowiednich sposobow obliczania podobienstwa,  $\text{sim}_{seq}$  i  $\text{sim}_{cat}$ , oraz
- 2) stopni waznosci,  $w_{seq}$  i  $w_{cat}$ .

Tak wiec mamy tu do czynienia z cala rodzina algorytmow, ktore uzyskuje sie z algorytmu bazowego przez przyjecie odpowiednich sposobow obliczania podobienstwa i stopni waznosci.

W zaproponowanym w rozprawie podejsciu, sposob liczenia podobienstwa odwotuje sie do pojecia *zawierania sie zbiorow rozmytych*. Zarowno klasyfikowany dokument  $d^*$ , jak i reprezentacje spraw  $\sigma$  i kategorii  $c$ , traktuje sie jako zbior rozmyty i podobienstwa (7)-(8) oblicza sie, odpowiednio, jako stopien zawierania sie  $d^*$  w  $\sigma$  i w  $c$ . To rozwiazanie okazuje sie byc skuteczne i efektywne.

Rozwiazanie zadania opiera sie wiec na zastosowaniu wskaznika zawierania sie zbiorow rozmytych,  $\text{sub}(A, B)$ :

$$\text{sub}: [0,1]^X \times [0,1]^X \rightarrow [0,1] \quad (10)$$

gdzie  $X$  oznacza przestrzen rozważań.

Okresla on stopien, w jakim zbior rozmyty  $A$  jest podzbiorem zbioru rozmytego  $B$ . W rozprawie rozwaza sie zastosowanie wielu roznych wskaznikow, w tym wskaznik zawierania Kosko i wskazniki bazujace na operatorze implikacji rozmytej i ich warianty.

W zadaniu archiwizacji, stanowiacym inspiracje dla zadania rozważanego w rozprawie, wystepuje dodatkowy aspekt, ktory zostal uwzględniony w proponowanym rozwiazaniu. Przyjmuje sie mianowicie, ze kategorie  $c \in C$ , do ktorych nalezy przypisac dokumenty, stanowia liście pewnej struktury hierarchicznej (drzewiastej)  $H$ . Kategorie te sa wiec ze soba powiazane i to ich powiazanie moze zostac wykorzystane w procesie klasyfikacji dokumentow. Zasadnicza idea takiego wykorzystania

hierarchii polega na założeniu, że jeśli dokonuje się wyboru pomiędzy kategoriami (liściami hierarchii) posiadającymi tego samego rodzica w hierarchii  $H$ , to przy ustalaniu wspomnianego stopnia zawierania się, rola słów kluczowych wspólnych dla tych kategorii jest mniej znacząca i wprowadza się dodatkowy stopień ważności słów kluczowych, który ma tę redukcję znaczenia reprezentować.

W związku z tym, w rozprawie proponuje się rozszerzenie pojęcia zawierania się zbiorów rozmytych z uwzględnieniem stopnia ważności elementów przestrzeni rozważań:

$$sub^W: [0,1]^X \times [0,1]^X \times [0,1]^X \rightarrow [0,1] \quad (11)$$

i bada się sposoby zaadoptowania znanych wskaźników zawierania się zbiorów rozmytych, w szczególności wskaźnika Kosko, tak aby uwzględniały one taki dodatkowy stopień ważności.

## Eksperymenty obliczeniowe

Sformułowany problem jest nowy i nie istnieją jeszcze standardowe korpusy danych testowych dla niego. W trakcie prac nad rozprawą podjęto próby stworzenia takich korpusów. Wykorzystano w tym celu dwa popularne zbiory danych: kolekcję publikacji z zakresu przetwarzania języka naturalnego i lingwistyki obliczeniowej, znaną jako ACL Anthology Reference Corpus (ACL ARC) oraz zestaw dokumentów znany w literaturze jako korpus Browna. W obydwu wypadkach, dokumenty składające się na poszczególne korpusy utożsamiono ze sprawami, a wydzielone z nich części jako dokumenty składające się na daną sprawę.

Przeprowadzono liczne eksperymenty obliczeniowe z użyciem zaproponowanego algorytmu i różnych konfiguracji jego parametrów (konkretnych wskaźników zawierania się zbiorów rozmytych (7)-(8) i stopni ważności  $w_{seq}$  i  $w_{cat}$  (we wzorze (9)). Eksperymenty obliczeniowe miały na celu porównanie zaproponowanej metody z bazowym algorytmem klasyfikacji typu  $k$ -nn oraz zweryfikowanie efektywności różnych parametryzacji zaproponowanego algorytmu.

## Wnioski

W rozprawie zaproponowano nowe podejście do analizy i rozwiązania konkretnego problemu kategoryzacji dokumentów tekstowych, w którym każdy dokument, oprócz przynależności do właściwej kategorii, należy do określonej sekwencji dokumentów, rozumianej tutaj jako sprawa, zawierającej dokumenty z tej samej kategorii. Zaproponowany algorytm łączy dwa wskaźniki względem których dokument jest klasyfikowany: jeden odpowiada za podobieństwo do sprawy a drugi odzwierciedla podobieństwo do kategorii. Wskaźniki te są oparte na wskaźnikach zawierania się zbiorów rozmytych. Zbadano skuteczność i efektywność zaproponowanego algorytmu dla różnych kombinacji wag obu wskaźników oraz zastosowanych wskaźników zawierania się zbiorów. Uzyskane wyniki są bardzo zachęcające. Jednak dalsze eksperymenty na różnych zbiorach dokumentów są potrzebne, aby możliwe było wyciągnięcie dalej idących wniosków. Otrzymane wyniki pokazują, że zaproponowana dość prosta metoda może być skuteczna i wydajna w radzeniu sobie z trudnym problemem klasyfikacji dokumentów w przypadku niewielkiej reprezentacji klas (spraw).

Ważnym założeniem było zastosowanie niskowymiarowej reprezentacji dokumentów i innych typów danych niezbędnych do realizacji zadania. Osiągnięto to przez przyjęcie wektorowej reprezentacji dokumentów z agresywną redukcją liczby uwzględnianych słów kluczowych.

Wprowadzono również pojęcie zawierania się zbiorów rozmytych z uwzględnieniem stopni ważności poszczególnych elementów przestrzeni rozważań. Pojęcie to wynikło w trakcie rozwiązywania zadania postawionego w rozprawie, a dokładniej jego wariantu, w którym dla rozróżnienia pokrewnych hierarchicznie kategorii lub spraw należących do tej samej kategorii

interesujące może być uwzględnienie zróżnicowanego znaczenia poszczególnych słów kluczowych przy klasyfikowaniu dokumentu. Pojęcie to ma jednak ogólniejszy charakter i jest interesujące samo w sobie. W rozprawie zaproponowano rozszerzenie poszczególnych wskaźników zawierania się zbiorów o uwzględnienie takich dodatkowych stopni ważności poszczególnych elementów. W szczególności przeprowadzono pogłębioną analizę takiego rozszerzenia w wypadku wskaźnika Kosko. Pokazano, że naiwne rozszerzenie tego wskaźnika może prowadzić do teoretycznie nieintuicyjnego jego zachowania, a użycie innej interpretacji bazowego wskaźnika jako punktu wyjścia pozwala przewyciężyć tę trudność.

Marek Gojewski