

Politechnika Warszawska Wydział Matematyki i Nauk Informacyjnych



Automatyczna kategoryzacja tematyczna tekstów przy użyciu metryk w przestrzeni ciągów znaków

Natalia Potocka *Warszawa*, 21.04.2014

Plan działania

- Cel pracy
- O metrykach słów kilka
- Postęp prac
- Co dalej?

CEL PRACY

Celem pracy jest skategoryzowanie tekstów z polskiej Wikipedii pod względem tematu na podstawie liczności słów występujących w tekście. Można się spodziewać, że jeśli w dwóch tekstach występuje dużo podobnych do siebie słów, to pochodzą one z tej samej kategorii tematycznej.

CEL PRACY

Celem pracy jest skategoryzowanie tekstów z polskiej Wikipedii pod względem tematu na podstawie liczności słów występujących w tekście. Można się spodziewać, że jeśli w dwóch tekstach występuje dużo podobnych do siebie słów, to pochodzą one z tej samej kategorii tematycznei.

Α		В		С		D	
całka	10	całka	5	niewłaściwe	3	ułamek	4
pochodna	5	pochodna	15	powieść	7	mianownik	5
niewłaściwa	4	granica	7	granica	15	niewłaściwy	6

CEL PRACY

Co ze słowami podobnymi?

Przykładowo słowa niewłaściwy i niewłaściwa mają ten sam temat, różnią się tylko rodzajem (męski / żeński). W tekstach mogą też występować błędy ortograficzne, błędy spowodowane brakami znaków diaktrycznych (ą, ę, ł, ...) itd. Takie słowa również chcielibyśmy traktować jako "podobne". W celu określenia jak bardzo dwa słowa są do siebie podobne, posłużą metryki określone na napisach.

OPERACJE EDYTOWANIA

Metryki oparte na operacjach edytowania zliczają liczbę opercji potrzebnych do przetworzenia jednego napisu w drugi. Najczęściej wymieniamymi operacjami są:

- zamiana znaku, np. $'ala' \rightarrow 'ela'$
- usunięcie znaku, np. $'ala' \rightarrow 'aa'$
- wstawienie znaku, np. $'ala' \rightarrow 'alka'$
- ullet transpozycja dwóch przylegających znaków, np. 'ala'
 ightarrow 'laa'

OPERACJE EDYTOWANIA

Metryki oparte na operacjach edytowania zliczają liczbę opercji potrzebnych do przetworzenia jednego napisu w drugi. Najczęściej wymieniamymi operacjami są:

- zamiana znaku, np. $'ala' \rightarrow 'ela'$
- usunięcie znaku, np. $'ala' \rightarrow 'aa'$
- wstawienie znaku, np. $'ala' \rightarrow 'alka'$
- transpozycja dwóch przylegających znaków, np. $'ala' \rightarrow 'laa'$

Przykładowe metryki: Hamminga, najdłuższego wspólnego podnapisu (longest common substring), Levenshteina, optymalnego dopasowania napisów (optimal string alignment), Damareu-Levenshteina.

Metryka najdłuższego wspólnego podnapisu, ozn. d_{lcs} , zlicza liczbę usunięć i wstawień, potrzebnych do przetworzenia jednego napisu w drugi. Np. $d_{lsc}('leia', 'leela') = 3$, bo $leela \xrightarrow{us. e} lela \xrightarrow{us. l} lea \xrightarrow{wst. i} leia$.

Metryka najdłuższego wspólnego podnapisu, ozn. d_{lcs} , zlicza liczbę usunięć i wstawień, potrzebnych do przetworzenia jednego napisu w drugi. Np. $d_{lsc}('leia','leela') = 3$, bo $leela \xrightarrow{us.\ e} lela \xrightarrow{us.\ l} lea \xrightarrow{wst.\ i} leia$. Odległość Levenshteina, ozn. d_{lv} zlicza sumę usunięć, wstawień oraz zamian znaków, potrzebnych do przetworzenia jednego napisu w drugi.

Metryka najdłuższego wspólnego podnapisu, ozn. d_{lcs} , zlicza liczbę usunięć i wstawień, potrzebnych do przetworzenia jednego napisu w drugi. Np. $d_{lsc}('leia','leela')=3$, bo $leela \xrightarrow{us.\ e} lela \xrightarrow{us.\ l} lea \xrightarrow{wst.\ i} leia$. Odległość Levenshteina, ozn. d_{lv} zlicza sumę usunięć, wstawień oraz zamian znaków, potrzebnych do przetworzenia jednego napisu w drugi. np. $d_{lv}('leia','leela')=2$, bo $leela \xrightarrow{us.\ e} lela \xrightarrow{zm.\ l\ na\ i} leia$.

Metryka optymalnego dopasowania napisów, ozn. d_{osa} , zlicza liczbę usunięć, wstawień, zamian oraz transpozycji przylegających znaków, potrzebnych do przetworzenia jednego napisu w drugi. Np. $d_{osa}('leia', 'leela') = 2$, bo $leela \xrightarrow{us. e} lela \xrightarrow{zm. l \ na \ i} leia$.

Co zostało zrobione?

 \bullet wczytano 1075~568artykułów z polskiej Wikipedii

Co zostało zrobione?

- wczytano 1 075 568 artykułów z polskiej Wikipedii
- razem to $2\ 806\ 765$ różnych słów...

Co zostało zrobione?

- wczytano 1 075 568 artykułów z polskiej Wikipedii
- razem to 2 806 765 różnych słów...
- ullet ... z czego 49% wystąpiło tylko w **jednym** tekście
- ullet ... a 44% wystąpiło tylko **jeden raz** we wszystkich tekstach

Co zostało zrobione?

- wczytano 1 075 568 artykułów z polskiej Wikipedii
- razem to 2 806 765 różnych słów...
- ullet ... z czego 49% wystąpiło tylko w **jednym** tekście
- ullet ... a 44% wystąpiło tylko **jeden raz** we wszystkich tekstach

Po usunięciu tzw. *stopwords*, czyli słów nieistotnych w kontekście analizy, jak np. *a, bo, co, jak, to, w, z, że*, słów jednoliterowych oraz słów w językach obcych z niełacińskiego alfabetu, pozostało $2\ 805\ 858$ słów do analizy.

UŻyWająca używający użyła zużywające zużywające UżyWa używająca używają użył używająca używają użył używaj używają używał używają używała używała używała używano używającym

```
dodając
poddania
wznając
uznając
```

uznawany uznawane o o uznawani w o uznawani w o uznawali uznawanej z n uznawanej uznawanej

RYSUNEK: Przykładowe grupy utworzone przy pomocy metryki Levenshteina. Maksymalna odległość w klastrze to 7.

zużywającym
Zażywający
używające
używające
zużywające
zużywające
używające w używają
zużywającymi używającymi zużywającymi zażywającymi zażywającymi zażywającymi zażywającymi zażywającymi używającymi używającym

kamieniach kamieniach

błoto błota

płoto błota

płoto płon spony
bodo procesowa podenia pode

 $\ensuremath{\mathrm{RYSUNEK}}$: Przykładowe grupy utworzone przy pomocy metryki lcs. Maksymalna odległość w klastrze to 7.

UŻyWająca używający użyła zużywające zużywający używane używające UŻyWa używającą używają użył używającą używaj użył używaż podają używaż używaż używającą używała używano używającym

rowlandem
hollander
hofmanna polsce
roelandem
hollandem
hollande
lowlands
poddania

zrywających używających grajacych grajacych grających grających grapiących prawiących prawiących prawiących prawiących prawiących zużywających zużywających trawiących grywających

Rysunek : Przykładowe grupy utworzone przy pomocy metryki osa. Maksymalna odległość w klastrze to 7.

Z powodu słabej jakości grupowania oraz braku możliwości obliczeniowych dokonano grupowania przy pomocy tzw. *stemmingu*. Polega on na przyporządkowaniu do słowa jego rdzenia, a więc takiej jego części, która jest odporna na odmiany przez rodzaje, przyimki, przypadki itd. Przykładowo dla słowa *używająca* rdzeniem jest *żyw*.

Z powodu słabej jakości grupowania oraz braku możliwości obliczeniowych dokonano grupowania przy pomocy tzw. *stemmingu*. Polega on na przyporządkowaniu do słowa jego rdzenia, a więc takiej jego części, która jest odporna na odmiany przez rodzaje, przyimki, przypadki itd. Przykładowo dla słowa *używająca* rdzeniem jest *żyw*. Do stemmingu użyto narzędzia Hunspell, które sprawdza pisownię dla wielu programów, takich jak: OpenOffice, Mozilla Firefox, Thunderbird czy Google Chrome.

Dzięki niemu udało się poklastrować 733 828 słów ($\approx 26\%$ wszystkich) z czego 89% stanowiły polskie słowa 5,5% - słowa angielskie, a po ponad 2% - słowa francuskie i niemieckie. Innych języków nie sprawdzano. Liczba uzyskanów grup (klastrów) to 186~942.

Co z pozostałymi słowami?

Słowa, które wystąpiły więcej niż raz we wszystkich tekstach, dołączono do już istniejących grup przy pomocy metryk. Takich słów było 973~855, co dało łącznie pogrupowanych słów w liczbie 1~707~683. Co więcej, grupy, które miało mało słów (między 1~a~5) połączono lub dołączono do innych zbiorów. W ten sposób uzyskano 13~różnych zbiorów grup słów.

użyjeszUżyłem

sużyłośż użyłem

sużyłośż użyliśmy

użyją użyła użyli

użyj użyc sużylję

użyt sużyt sużyliście

użyto użyto użyty

użyto użyto użyty

użyto użyto użyty

użyto użyto użyty

yokoyamie białykamień akiyamie kamieniem kamieniem grakamieńcze kamieńcze kamień kageyamie kamień gryzikamień kamień kamadiem kamień kamieniami kamieniami kamieniami kamieniami kamieniami

hotem 5 nabłoto bołotowem błotowij

 $rac{ ext{RYSUNEK}}{ ext{lcs}}$: Przykładowe klastry utworzone przy pomocy Hunspella oraz metryki lcs.

Następnie dla próbki tekstów z trzech kategorii: matematyka, historia sztuki oraz wojny, dokonano grupowania artykułów. Kryterium była liczność **grup słów** występujących w danym tekście. Do grupowania użyto metody sferycznych k-średnich.

Opierając się na kategoriach z Wikipedii, poprawnie sklasyfikowanych zostało 61% z 59~403 artykułów.

tytuł	kat	id_kat	kl
kościół św. rocha w poznaniu	szt	1	1
portret	szt	1	2
quantum of solace (gra komputerowa)	szt	1	2
kurka wodna (seria gier)	szt	1	2
technika macierzy rzadkich	mat	2	2
kryterium walda	mat	2	2
generalized markup language	mat	2	2
czesław falkiewicz		3	3
william goodenough		3	3
kazimierz gallas	woj	3	3
wacław krzywiec	woj	3	3
fabian aleksandrowicz	woj	3	3

Metoda ta dała dość dobre rezultaty dla małej próbki i małej liczby tematów. Dla większej liczby tematów i większej próbki, R miał problemy z pamięcią...

Metoda ta dała dość dobre rezultaty dla małej próbki i małej liczby tematów. Dla większej liczby tematów i większej próbki, R miał problemy z pamięcią...

Stąd pomysł użycia metody mini-batch kmeans w pythonie. Dla 20~482artykułów i 981 tematów, obliczenia trwały ok. 70s. Wyniki dokładności były trochę gorsze, tzn. ok. 59%artykułów z tego samego tematu znalazło się w tej samej grupie.

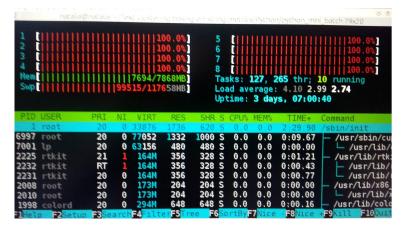
Metoda ta dała dość dobre rezultaty dla małej próbki i małej liczby tematów. Dla większej liczby tematów i większej próbki, R miał problemy z pamięcią...

Stąd pomysł użycia metody mini-batch kmeans w pythonie. Dla 20~482 artykułów i 981 tematów, obliczenia trwały ok. 70s. Wyniki dokładności były trochę gorsze, tzn. ok. 59% artykułów z tego samego tematu znalazło się w tej samej grupie.

Postanowiono puścić obliczenia na jednym ze zbiorów na wszystkich artykułach (1 $075\ 464$), dla $67\ 969$ tematów. Same dane wczytywały się 3 godziny i zajęły...

Postanowiono puścić obliczenia na jednym ze zbiorów na wszystkich artykułach (1 075 464), dla 67 969 tematów. Same dane wczytywały się 3 godziny i zajęły... $100~\mathrm{GB}~\mathrm{RAMu}$ (+swap). I algorytm liczył... i liczył...

Postanowiono puścić obliczenia na jednym ze zbiorów na wszystkich artykułach (1 075 464), dla 67 969 tematów. Same dane wczytywały się 3 godziny i zajęły... $100~{\rm GB}~{\rm RAMu}~(+{\rm swap})$. I algorytm liczył... i liczył... Po 22 dniach stwierdziłam, że czas przerwać obliczenia.



Kolejne testy wykazały, że bardzo duża liczba tematów powoduje, że obliczenia przetwarzają się w nieskończoność. W tej chwili jestem na etapie zmniejszenia liczby tematów i rozproszenia obliczeń na kilka(-naście) komputerów.

BIBLIOGRAFIA

- [1] Martin Kober Christian Buchta Kurt Hornik, Ingo Feinerer. Spherical k-means clustering. *Journal of Statistical Software*, 50(10):1–22, 9 2012.
- [2] Mark P. J. van der Loo. The stringdist Package for Approximate String Matching. *The R Journal*, 6:111–122, 2014.
- [3] Stefan Wild. Seeding Non-Negative Matrix Factorizations with the Spherical K-Means Clustering. University of Colorado, Colorado, 2002.

Dziękuję za uwagę.