



Narzędzie do komputerowej analizy dyskusji na forum

Projekt inżynierski

5 stycznia 2021

Karolin Izabel Boczoń Marcin Armacki
Michał Romaszkin



Wszystko zaczęło się od danych

Możliwe do pozyskania lokalnie specyficzne dane językowe:
dyskusje studentów na forach prowadzonych w ramach zajęć na wydziale.

Pomysły:

- ▶ Topic modeling
- ▶ “Lanie wody”
- ▶ Polaryzacja opinii
- ▶ Rodzaje argumentów



Dane tekstowe

- ▶ Fora z przedmiotów prowadzonych przez prof. Marciniaka o tematyce związanej ze sztuczną inteligencją i systemami informatycznymi z ostatnich kilku lat
- ▶ Około 2400 postów w sumie
- ▶ W ramach pojedynczego forum tematy podzielone na wątki merytoryczne i dyskusyjne



Polaryzacja opinii

Chcemy wskazać, które wypowiedzi zgadzają się lub nie zgadzają się z tezą postawioną w wątku.

Dzielimy każdą wypowiedź na akapity, a każdy akapit tagujemy jako:

- ▶ *pozytywny* („Zgadzam się, że”, „Masz rację”, „Doskonale powiedziane”)
- ▶ *negatywny* („Nie zgadzam się, że”, „Nie masz racji”, „To bzdura”)
- ▶ *mieszany* („Niby się zgadzam, ale”, „To zależy”)
- ▶ *fakt* („Pierwsze komputery powstały w latach 40. XX wieku”)
- ▶ *nieistotny* (nagłówki, linki, etc.)



Polaryzacja opinii c.d.

Udało się zaanotować 1500 paragrafów.

Został stworzony model regresji logistycznej, jednak osiągnął dość marne wyniki: precyzja na poziomie 59%.

Możliwe przyczyny:

- ▶ niestaranny podział na akapity,
- ▶ rozmyty charakter klasyfikacji,
- ▶ zły dobór zbioru uczącego

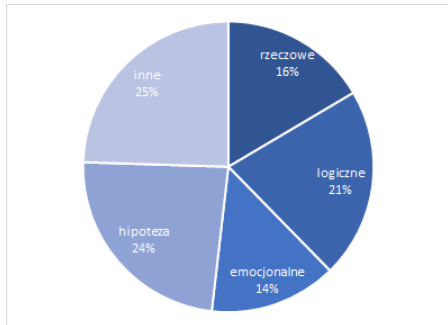
Chcemy wskazać typ argumentu użytego w dyskusji.

Znowu dzielimy każdą wypowiedź na akapity, a każdy akapit tagujemy jako:

- ▶ *argument rzeczowy* (fakty, liczby, dane statystyczne, cytaty)
- ▶ *argument logiczny* (analogia, indukcja, dedukcja, dylemat)
- ▶ *argument emocjonalny* (oparty na uczuciach: sympatii, lęku, ambicji, litości itp.)
- ▶ *inny rodzaj argumentu* (przykład z życia argumentującego, hipotetyczna sytuacja)
- ▶ *hipoteza* (stwierdzenie, które można potwierdzić lub obalić)

Lepszy zbiór uczący:

- ▶ dopracowane dzielenie na akapity
- ▶ wybrano tylko wypowiedzi z tematów „dyskusyjnych”
- ▶ bardziej jednolity udział wszystkich tagów



Niestety, udało się zaanotować tylko 600 akapitów.



Rodzaje argumentów - model

Problem jest podobny do analizy wydźwięku.

Użyto sieci neuronowej typu LSTM (*Long short-term memory*) z uwagi na osiągnięcia w podobnych zadaniach.

Wyniki okazały się fatalne, na poziomie 20% precyzji.

Możliwe przyczyny:

- ▶ zdecydowanie za mały zbiór uczący
- ▶ przynależność akapitów do różnych klas jednocześnie
- ▶ sieć uczy się szybko, ale ma problem z uogólnieniem klasyfikacji (*overfitting*)



Rodzaje argumentacji - model v2

Przy kolejnym podejściu wykorzystano schemat wykorzystywany przy klasyfikacji dokumentów do kategorii tematycznych. Stworzony model to wielomianowy naiwny klasyfikator bayesowski, który dobrze radzi sobie na małych zbiorach danych.





Precyzja tego modelu była średnio dwa razy wyższa niż modelu neuronowego. Możliwe ulepszenia:

- ▶ powiększenie zbioru uczącego
- ▶ wyrównanie dysproporcji pomiędzy klasami



The screenshot shows a web browser window with the address bar displaying "localhost:4200/my-profile/discussions/1/topic". The page title is "Mój profil". On the left, there is a sidebar with the heading "Dyskusje" and several discussion topics listed, such as "Autonomiczne pojazdy [M]" and "Systemy dialogowe [M]". The main content area displays two discussion threads. The first thread, titled "2944", contains two paragraphs of text and two dropdown menus for selecting labels ("mieszana" and "pozytywna"). The second thread, titled "1843", contains one paragraph of text and one dropdown menu for selecting a label ("mieszana"). At the top of the main content area, there are three buttons: "WRÓC DO WYBORU DYSKUSJI" (green), "ZMAPUJ DANE Z PLIKU USERS.XML" (blue), and "USUN ZMAPOWANE DANE" (red), along with a "Polaryzacja opinii" dropdown menu.

References

-  Sławomir Dadas, *A repository of polish NLP resources*, Github, 2019.
-  Sepp Hochreiter and Jürgen Schmidhuber, *Long short-term memory*, *Neural computation* **9** (1997), 1735–80.
-  Daniel Jurafsky and James H. Martin, *Speech and language processing (2nd edition)*, Prentice-Hall, Inc., USA, 2009.
-  Jeffrey Pennington, Richard Socher, and Christopher D. Manning, *Glove: Global vectors for word representation*, *Empirical Methods in Natural Language Processing (EMNLP)*, 2014, pp. 1532–1543.