

Raport z klasyfikacji nagłówków

Maria Galanty - kandydatka na stanowisko Junior Data Scientist

Luty 2020

1 Wstęp

Zadanie polegało na zbudowaniu modelu, który sklasyfikuje nagłówki jako sarkastyczne bądź nie. W przedstawionym rozwiązaniu wykorzystałam Pythona.

2 Eksploracja danych

Dane zawierały 26709 rekordów w tym 11724 nagłówków satyrycznych i 14985 w drugiej klasie, nie było brakujących zmiennych. Dwie klasy nie różnią się znacząco między sobą długością tekstu ani ilością zawartych znaków interpunkcyjnych.

3 Text Preprocessing

Przy przetwarzaniu języka naturalnego występują dwa problemy związane z formatem danych. Po pierwsze klasyfikatory potrzebują zmiennych numerycznych, a nie tekstowych. Algorytmy te, wymagają wektora liczbowego. Można to rozwiązać za pomocą podejścia *bag-of-words*, w którym każde unikalne słowo w tekście będzie reprezentowane przez jedną liczbę. W związku z czym pojawia się drugi problem, który polega na ograniczeniu ilości wyrazów poddanych analizie. W tym celu najpierw zostały usunięte *stop words* są to najbardziej popularne wyrazy w języku angielskim, które wiążą się z gramatyką języka, należą tutaj takie wyrazy jak *a*, *the*, *you*. Następnie wykonany został *stemming*, który różne formy wyrazów sprowadza do formy podstawowej - korzenia wyrazu. Przykładowo *playing*, *play*, *plays* zostaną zamienione na *play*. Kolejnym krokiem było usunięcie znaków interpunkcyjnych. Tak przygotowane dane zostały zmienione na wektory liczbowe. TF-IDF jest skrótem terminu *Term Frequency-Inverse Document Frequency* i jest algorytmem przekształcającym tekst w reprezentację liczbowa.

4 Klasyfikator

Ze względu na zbalansowany zbiór danych, do ewaluacji modelu użyłam cross-walidacji i metrykę *accuracy*. Przy cross-walidacji (umożliwiającej sprawdzenie, czy model nie jest przeuczony, co jest szczególnie ważne przy tak dużym zbiorze danych) podzieliłam zbiór na 5 części, co pozwoliło trenować model na 4 z nich i testować na pozostałej. Zbudowałam następujące dwa klasyfikatory: Random Forest oraz Multinomial Naive Bayes. Drugi z nich uzyskał lepszy średni wynik klasyfikacji na danych testowych wynoszący 0.78.