

# Books Clustering

Drazkowski Hubert, Wilk Marcin

Czerwiec 2021

## 1 Wstep

Pierwszym celem projektu jest zastosowanie metod analizy grupowania danych (ang. clustering).

Drugim celem projektu jest odpowiedź na pytanie jak podobne są do siebie starożytne teksty religijne.

Analiza statystyczna została wykonana w programie Python. Dla czytelności i zwiezłości raportu niektóre opisy metod lub rysunki wykorzystywane podczas analizy zostały pominięte. Można je znaleźć w notatniku jupyter. Wykresy do modelowania także w prezentacji.

Inspiracją do stworzenia danych wyjściowych był "Project Gutenberg".  
Czyszczenie danych przeprowadzili PREETI SAH i ERNEST FOKOUE w (2019)  
"WHAT DO ASIAN RELIGIONS HAVE IN COMMON? AN UNSUPERVISED  
TEXT ANALYTICS EXPLORATION".

Dostępne były dwa źródła danych:

- Macierz rzadka gdzie zmiennymi była liczba wystąpień danego słowa w danym fragmencie
- Oryginalny tekst ze wszystkimi fragmentami skonkatenowanymi

Fragmenty pochodzą z 8 tekstów z czterech religii. Stąd potencjalnie rodzą się dwie warstwy etykietek

- Hinduism : Yogasutras, Upanishads
- Buddhism : Four Noble Truth of Buddhism
- Taoism : Tao Te Ching
- Christianity : Book of Proverb, Book of Ecclesiastes, Book of Ecclesiasticus, Book of Wisdom

## 2 Eksploracyjna analiza danych i inżynieria cech

Wnioski po eksploracyjnej analizie danych macierzy zliczeń słów:

1. Jest bardzo dużo słów, które rzadko występują, ponad 50% z nich występuje 2 razy lub rzadziej.
2. Mamy do czynienia z problemem dużej ilości wymiarów. Obserwacji jest 589, a zmiennych 8266
3. Niektóre słowa pojawiają się często, rzadko zdarzają się słowa, które występują wiele ilości razy we fragmencie
4. Fragmenty różnią się ilością zawieranych słów. Zdecydowana większość jest w przedziale 0-200. Potem jest kilkadziesiąt do 400, powyżej 400 jest niewiele obserwacji.

Można zauważyć, że mamy tutaj problem przekleństwa wymiarowości. Mamy do czynienia z macierzą o wymiarach 590 na 8266, na dodatek jest to macierz rzadka, to jest pełna zer. Kolumny stanowią indykatory zliczające wystąpienie danego słowa w korpusie. Każdy jeden fragment korpusu jest obserwacją. Zobaczmy czy jakieś charakterystyczne słowa występują najczęściej w tych tekstach.

Niektóre pojawiają się naprawdę często, nie są one jednak na pierwszy rzut oka charakterystyczne dla jakichkolwiek tekstów. Może być jednak tak, że okaże się zupełnie co innego, bo tylko część z ksiąg będzie używała składni do której będą wymagane niektóre z tych słów. Przyjrzyjmy się z kolei tym rzadziej występującym słowom, będziemy tutaj szukać czy da się tutaj jakoś zmniejszyć wymiar naszych danych przez usunięcie słów, które występują bardzo rzadko albo w ogóle.

Po przeczyszczeniu oryginalnego tekstu tak aby pozbyć się dziwnych znaków i uzyskać mniej więcej oryginalny podział (z dokładnością do pewnych słów nieinformatywnych usuniętych przez autorów artykułu) tworzymy nową ramkę danych. Ilość korpusów zgadza się z wielkością macierzy. Analiza sentymentu zawiera między innymi dwa wskaźniki polarity i subjectivity. Polarity wpada w przedział  $[1,1]$ , a im niższa wartość tym tekst oceniany jest za bardziej nacechowany negatywnie zaś subjectivity w  $[0,1]$ , a interpretować należy te skale w taki sposób, że wartości bliższe zera powinny wskazywać na mniej nacechowane subiektywnie treści.

Im większe ARI tym bardziej skomplikowany tekst Im większe FRI tym tekst łatwiejszy Powinny być ze sobą ujemnie skorelowane ze względu na wspólny czynnik we wzorze

$$ARI = 4.71\left(\frac{\text{characters}}{\text{words}}\right) + 0.5\left(\frac{\text{words}}{\text{sentences}}\right) - 21.43,$$

$$FRE = 206.835 - 1.015\left(\frac{\text{words}}{\text{sentences}}\right) - 84.6\left(\frac{\text{syllables}}{\text{words}}\right).$$

Sentiment : <https://github.com/sloria/TextBlob>, leksykon zdefiniowany subiektywnie przez współczesnych ekspertów.

Nie ma sensu zliczać ilości znaków lub słów. Bardziej zależy nam na cechach języka i używanego słownictwa, przekazywanych treściach. Te sama myśl przekazują w swojej konstrukcji współczynniki ARI i FRE. Dodamy jedynie średnią długość użytego słowa w zdaniu.

Pozwoliliśmy sobie na rzucenie okiem także na oryginalny plik tekstowy, w następnych krokach rysujemy ładnie wyglądającą mapę wyrazów najczęściej się pojawiających i przykładamy do owego tekstu dwa narzędzia rodem z NLP to jest wskaźnik subjectivity i polarity

- Polarity (użycie większości słów nacechowanych negatywnie albo pozytywnie)
- Subjectivity (użycie słów uważanych za opiniotwórcze)
- FRE (indeks mierzący skomplikowanie tekstu)
- ARI (indeks mierzący skomplikowanie tekstu)
- avg\_wordlength (długość średnia słowa w fragmencie)

Naszym głównym celem w tej części była redukcja wymiarów oraz dodanie kilku cech przy pomocy NLP na podstawie oryginalnych tekstów. Pierwszym podejściem było znalezienie słów, które występują tylko w jednej księdze. Dla każdej księgi policzylismy ile takich różnych słów zawiera oraz ile w sumie. Następnie usuneliśmy te słowa z naszej ramki danych i przeprowadzimy PCA.

Na oryginalnej macierzy dokonaliśmy kilku przekształceń:

- Aby zmniejszyć liczbę wymiarów usuwamy słowa które wystąpiły co najwyżej dwa razy w całym korpusie. (zostało 4394)
- Aby jakoś zachować informacje mierzymy ile razy wystąpiło jakieś unikatowe (i2) słowo w danym fragmencie. Do tego dodajemy ile takich unikatowych słów miał fragment.
- Na tym robimy PCA, uzyskując 91% tłumaczonej wariancji zostajemy przy 200 składowych głównych

Ostatecznie przed modelowaniem

- Łączymy ze sobą ramkę danych po PCA i z oryginalnej analizy tekstu za pomocą narzędzi NLP.
- Wystandardyzowujemy niektóre zmienne, bo nowo dołączone cechy miały znacznie większe wartości od tych zwróconych przez PCA, tak aby wszystkie kolumny były mierzone w podobnej skali.

### 3 Modelowanie

Patrzyliśmy na beterie metod :

1. Kmedioids
2. Kmeans
3. Agglomerative z połączeniami
  - Single
  - Average
  - Complete
  - Ward
4. DBSCAN
5. GMM

Następnie dodaliśmy zdefiniowaną przez nas metrykę  $sum(|x - y| / |x + y|)$ . Znowu patrzyliśmy na wyżej wymienione metody. Inspiracją do jej stworzenia była większa chęć zaakcentowania różnic na małych liczbach słów. W metryce euklidesowej jeśli w jednej księdze słowo  $x$  jest 0 razy a w drugiej 4, natomiast inne słowo w pierwszej księdze jest 48 razy a w drugiej 52, to wpływają one tak samo na odległość między tymi księgami. W nowej metryce pierwsze słowo dodawałoby 1 do metryki natomiast drugie słowo tylko 0.04.

Porównywaliśmy metody za pomocą wykresu łokciowego i wykresu miary Silhouette. Silhouette score był dla własnej metryki dużo wyższy niż w przypadku metryki euklidesowej, jednak nie powinniśmy jednoznacznie stwierdzać wyższości tej metody, gdyż silhouette jest wrażliwy na metrykę użytą do jego obliczania. Na sam koniec pojeśliśmy się próby wizualizacji wyników za pomocą t-SNE i policzenia homogeniczności używając oryginalnych klas.

### 4 Wnioski

Porównanie wykonywaliśmy na dwóch etykietach. Etykietach przynależności do religii i etykietach przynależności do konkretnej księgi. DBSCAN wyłapał tutaj tylko jeden duży klaster i ewentualne pojedyncze klasy odstające. Nie ma sensu tutaj nawet rozważać homogenity. Nasza własna metryka pomimo dobrego silhouette daje bardzo słabe wyniki dla naszego zbioru danych. Nie potrzebujemy liczyć homogenity, żeby to zobaczyć, wystarczy spojrzeć na t-SNE. Wybraliśmy do ostatecznego porównania modele :

- AgglomerativeClustering(n\_clusters=6,linkage="average")
  - Homogeneity religion : 0.38
  - Homogeneity book : 0.29

- GaussianMixture(n\_components=8, covariance\_type="full")
  - Homogeneity religion :0.40
  - Homogeneity book : 0.32
- KMeans(n\_clusters=6,random\_state=0)
  - Homogeneity religion : 0.39
  - Homogeneity book : 0.30