

Projekt : Clustering

H.Drażkowski, M.Wilk

08 czerwiec 2021

Wydział Matematyki i Nauk Informacyjnych
Politechnika Warszawska

Agenda

Agenda

1. Agenda
2. EDA
3. Feature Engineering
4. Modelowanie
5. Wyniki

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.

EDA

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

Wnioski po eksploracyjnej analizie danych macierzy zliczeń słów:

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

Feature Engineering

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.

Dodajemy

- 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 we fragmencie)

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 jakiegokolwiek unikatowe (<2) słowo w danym fragemencie. 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

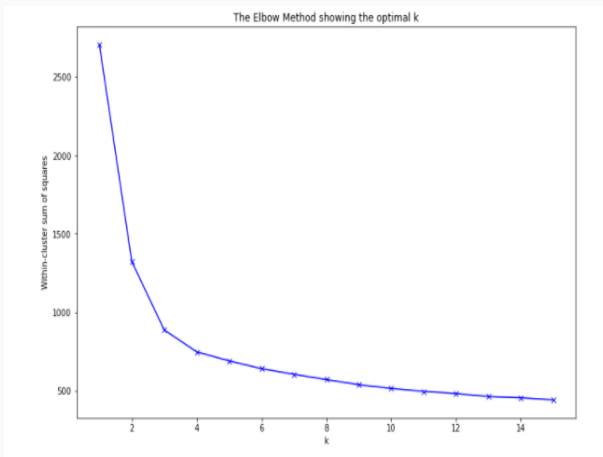
- Łączymy ze sobą ramke danych po PCA i z oryginalnej analizy tekstu za pomocą narzędzi NLP.
- Wystandaryzowujemy niektóre zmienne, bo nowo dołączone cechy miały znacznie większe wartości od tych zwróconych przez PCA, tak aby wszystkie kolumny były mierzone w podobnej skali.

Modelowanie

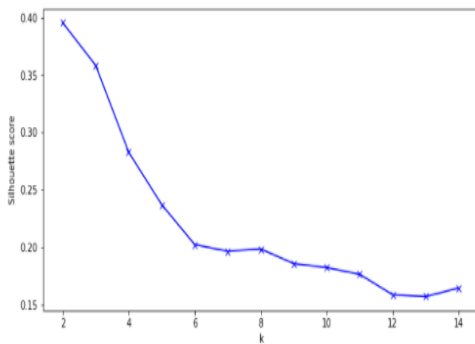
Patrzyliśmy na baterię metod :

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

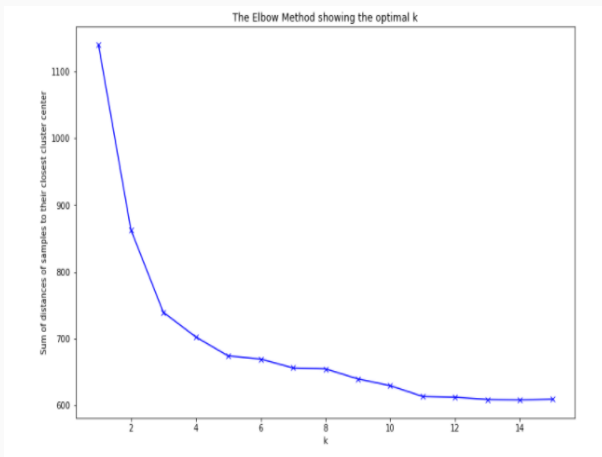
K-Means



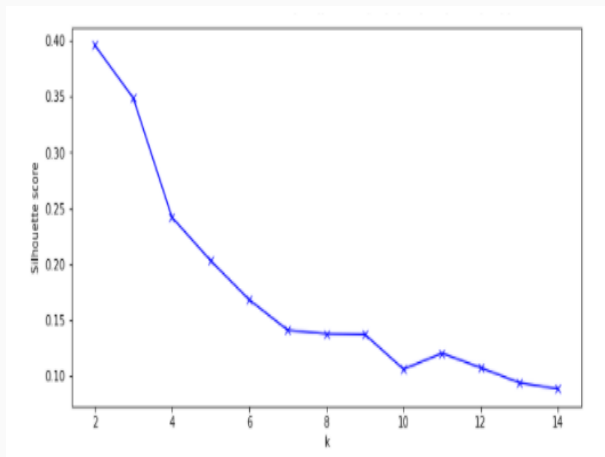
K-Means



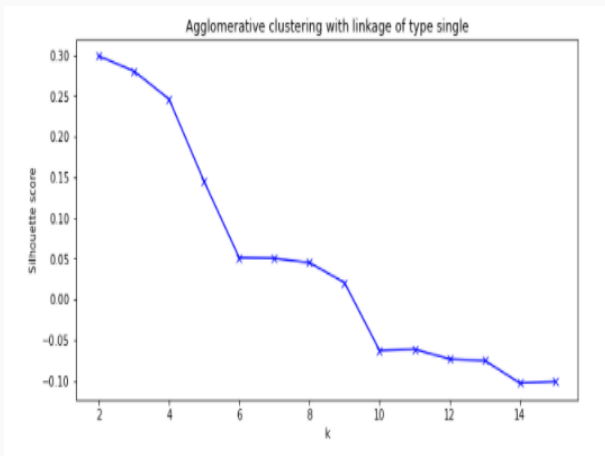
K-Medoids



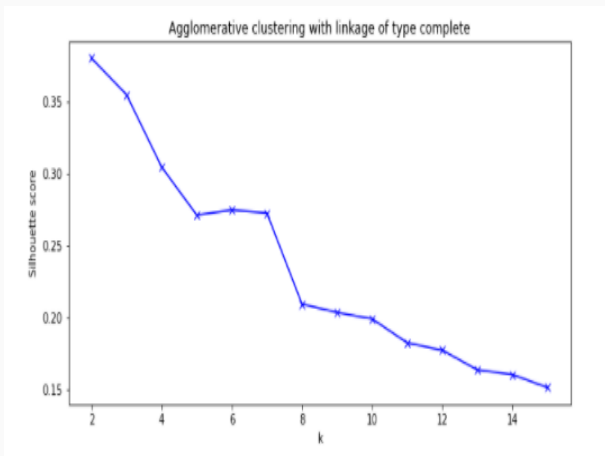
K-Medoids



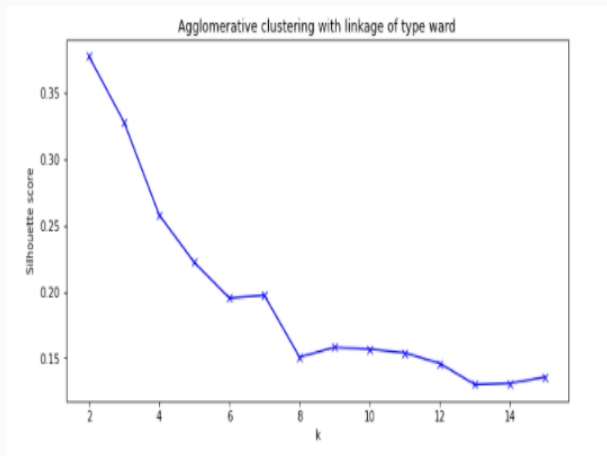
Agglomerative



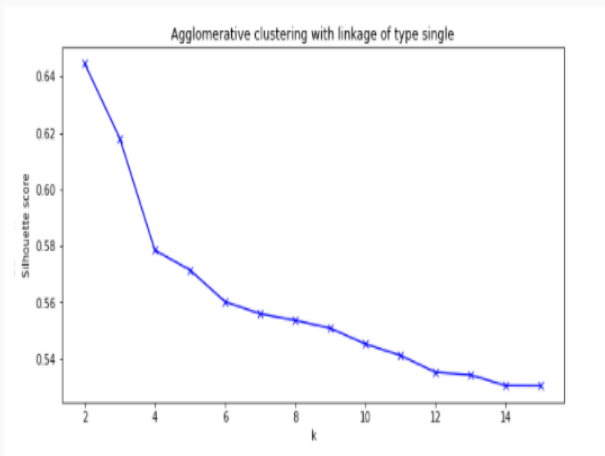
Agglomerative



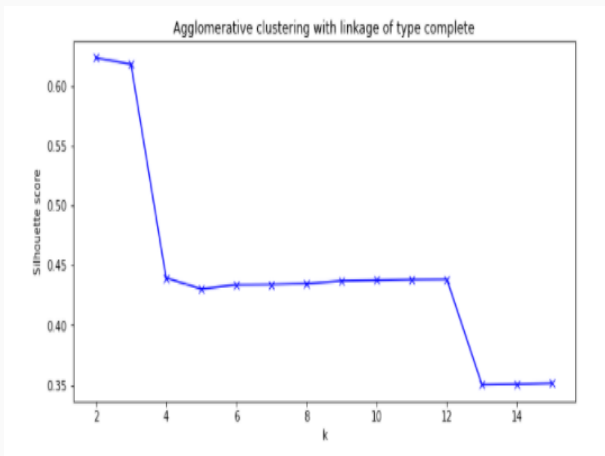
Agglomerative



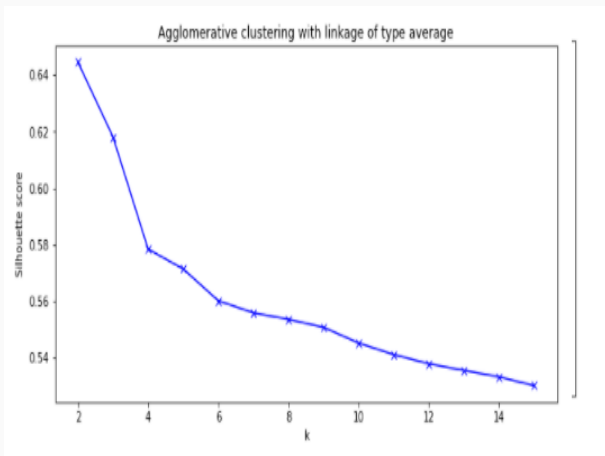
Agglomerative



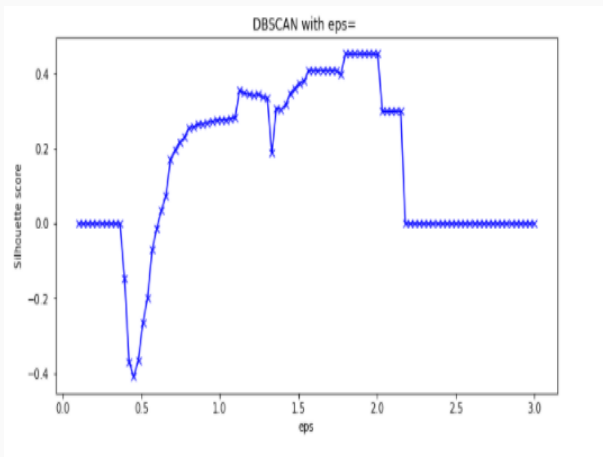
Agglomerative



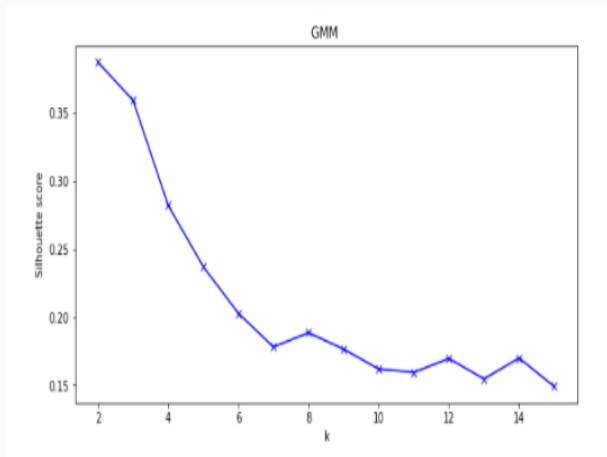
Agglomerative



DBSCAN



GMM



Wyniki

Porównywaliśmy metody za pomocą wykresu łokciowego i wykresu miary Silhouette.

Silhouette score by 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.

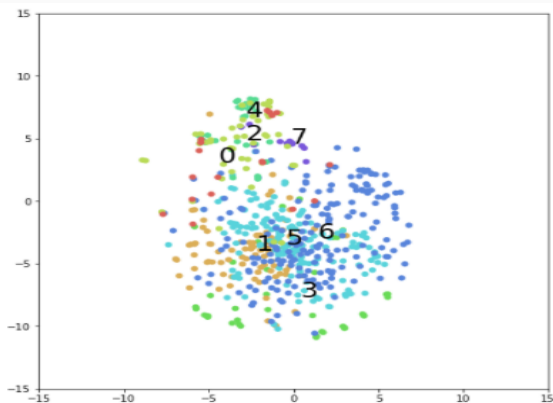
Na sam koniec pojeśliśmy się próby wizualizacji wyników za pomocą t-SNE i policzenia homogeniczności używając oryginalnych klas.

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 ostatecznie trzy modele do porównania:

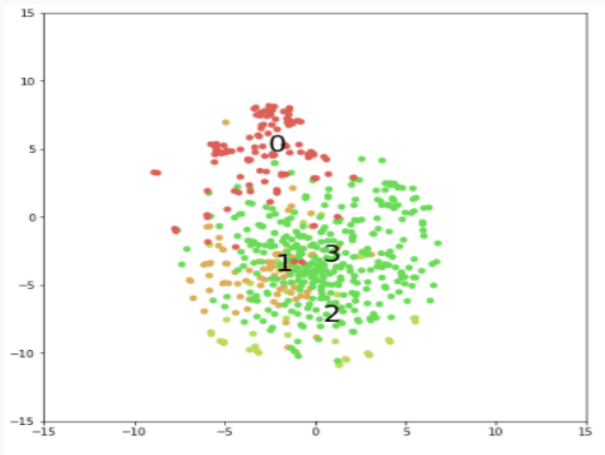
- `AgglomerativeClustering(n_clusters=6,linkage='average')`
 - Homogeneity religion : 0.38
 - Homogeneity book : 0.29
- `GaussianMixture(n_components=8, covariance_type='full')`
 - Homogeneity religion : 0.40
 - Homogeneity book : 0.32
- `KMeans(n_clusters=6,random_state=0)`
 - Homogeneity religion : 0.39
 - Homogeneity book : 0.30

Książki na płaszczyźnie 2D



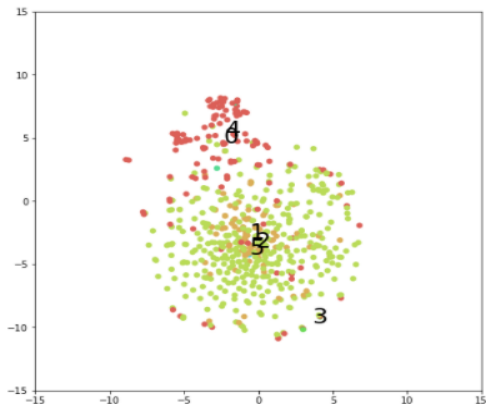
Tu widzimy jak prezentują się religie na płaszczyźnie 2D.

Religie na płaszczyźnie 2D



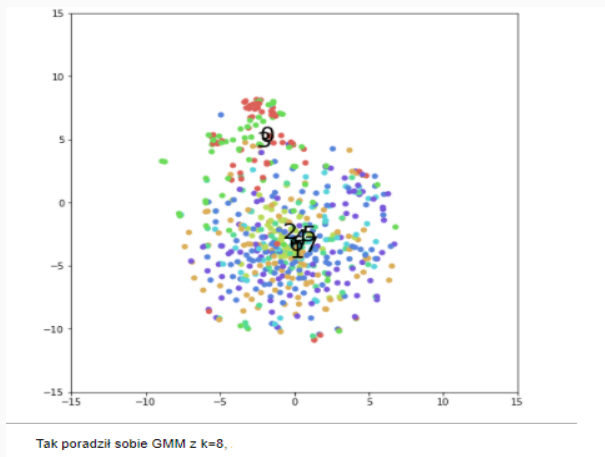
Podsumowanie

`AgglomerativeClustering(n_clusters=6,linkage='average')`



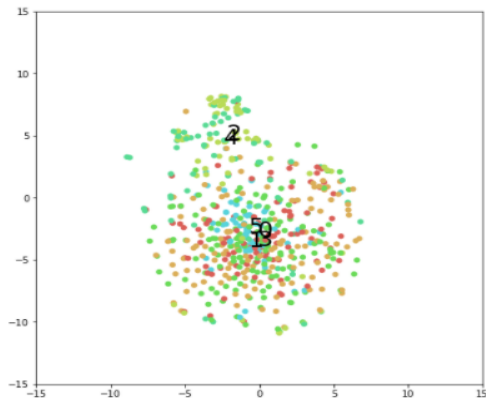
Tak prezentuje się agglomerative clustering z metryką euklidesową, `linkage='average'` i `k=6`.

GaussianMixture(n_components=8, covariance_type="full")



Podsumowanie

KMeans(n_clusters=6,random_state=0)



Tak wygląda podział przewidziany przez KMeans.

DBSCAN

