

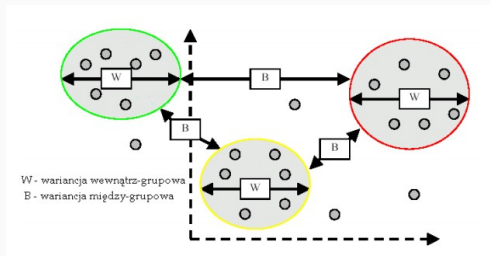
KLASTERYZACJA

Paweł Buglewicz, Krzysztof Cieśla, Justyna Olczak

20.12.2016

Politechnika Krakowska

Klasteryzacja (grupowanie) jest jedną z metod nienadzorowanej (bez dostępnej a priori wiedzy) analizy danych. Głównym celem klasteryzacji jest podział rozpatrywanego zbioru obiektów na grupy (klastry), w ten sposób, aby każda z grup była możliwie jednorodna (tzn. zawierała elementy podobne do siebie), a jednocześnie poszczególne klastry były jak najbardziej zróżnicowane między sobą (rys. 1).



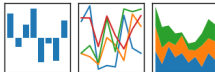
Rys. 1: Idea klasteryzacji. Dążymy do takiego podziału zbioru danych aby wariancja wewnątrz-grupowa (w każdym z klastrów była możliwie mała) a jednocześnie wariancja między-grupowa możliwie duża.

Klasteryzacji używa się w celu zredukowania złożoności wielowymiarowego problemu. Pozwala ona zauważyć podział, który mógłby być niewidoczny dla badacza bez dostępu do komputera (który sprawia, że czas znalezienia podziału jest znacznie skrócony). Gdy nastąpił już podział na grupy można opracować wiele prostych modeli dla każdej z osobna, zamiast jednego, bardzo złożonego modelu.



pandas

$$y_{it} = \beta' x_{it} + \mu_i + \epsilon_{it}$$

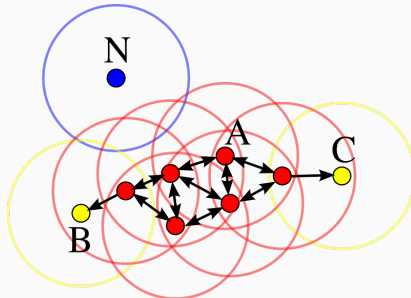


METODY KLASTERYZACJI - PRZYKŁADY

Method name	Parameters	Scalability	Usecase	Geometry (metric used)
K-Means	number of clusters	Very large <code>n_samples</code> , medium <code>n_clusters</code> with MiniBatch code	General-purpose, even cluster size, flat geometry, not too many clusters	Distances between points
Affinity propagation	damping, sample preference	Not scalable with <code>n_samples</code>	Many clusters, uneven cluster size, non-flat geometry	Graph distance (e.g. nearest-neighbor graph)
Mean-shift	bandwidth	Not scalable with <code>n_samples</code>	Many clusters, uneven cluster size, non-flat geometry	Distances between points
Spectral clustering	number of clusters	Medium <code>n_samples</code> , small <code>n_clusters</code>	Few clusters, even cluster size, non-flat geometry	Graph distance (e.g. nearest-neighbor graph)
Ward hierarchical clustering	number of clusters	Large <code>n_samples</code> and <code>n_clusters</code>	Many clusters, possibly connectivity constraints	Distances between points
Agglomerative clustering	number of clusters, linkage type, distance	Large <code>n_samples</code> and <code>n_clusters</code>	Many clusters, possibly connectivity constraints, non Euclidean distances	Any pairwise distance
DBSCAN	neighborhood size	Very large <code>n_samples</code> , medium <code>n_clusters</code>	Non-flat geometry, uneven cluster sizes	Distances between nearest points
Gaussian mixtures	many	Not scalable	Flat geometry, good for density estimation	Mahalanobis distances to centers
Birch	branching factor, threshold, optional global clusterer.	Large <code>n_clusters</code> and <code>n_samples</code>	Large dataset, outlier removal, data reduction.	Euclidean distance between points

Rys. 2: Porównanie algorytmów grupowania w scikit-learn.

- **Density-Based Spatial Clustering of Applications with Noise** - Martin Ester, Hans-Peter Kriegel, Jörg Sander i Xiaowei Xu (1996).
- DBSCAN jest jedną z najprostszych, najpowszechniejszych oraz najszybszych metod klasyfikacji. Polega ona na szukaniu klastrów jako obszarów o zwiększonej gęstości, oddzielonych obszarami o mniejszej gęstości.



Zalety:

- Jako parametru początkowego nie wymaga liczby spodziewanych klastrów w danych.
- Wymaga tylko dwóch parametrów początkowych:
 - Minimalna liczba punktów w sąsiedztwie
 - Maksymalny promień sąsiedztwa
- Potrafi wyszukać klastry nawet o bardzo skomplikowanym kształcie.

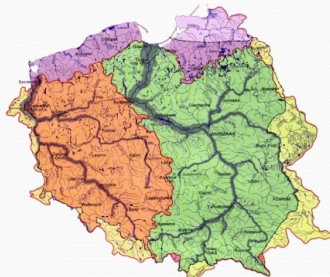
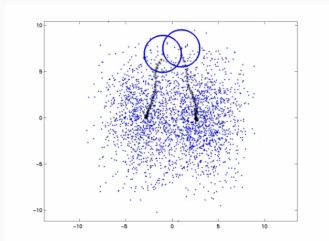
Wady:

- Nieprzydatny dla wielowymiarowych danych (pomiar odległości między punktami).
- Nie działa jeśli różnice w gęstości są zbyt duże.

Algorytm poszukuje klastrów, których centra znajdują się w atraktorach danego układu. Atraktorem w tym wypadku nazywamy punkt, do którego docieramy, podążając od każdego punktu w kierunku, w którym gradient jest największy. Ten krok nazywamy *Średnim przesunięciem* - *mean shift*.

Dla danych dyskretnych musi być stosowany odpowiedni estymator gradientu.

Dla wzrokowców - jest to analogiczne do zlewisk rzek



Zalety:

- Można zastosować go do dowolnej liczby wymiarów.
- Kształt klastra może być dowolny.

Wady:

- Czuły na wybór parametrów początkowych.
- Przy dużej liczbie próbek algorytm znacznie zwalnia.

Metoda affinity propagation szuka klastrów poprzez wymianę informacji między punktami. Algorytm określa podobieństwo między parami punktów (poprzez to, jaką informację wymieniły), jednocześnie je modyfikując (zmienia się odległość tych punktów w przestrzeni wielowymiarowej). Po wielu krokach, jeśli punkty są odpowiednio podobne, czyli odpowiednio bliskie, stają się jednym klastrem. Osobne klastry nie są podobne do siebie i się "odpychają".

Cechy punktów są opisywane za pomocą macierzy – wartości liczbowych.

Zalety:

- Jako parametr początkowy nie jest potrzebne zadanie liczby klastrów
- Może być stosowany w dowolnej liczbie wymiarów
- Stabilny.

Wady:

- Wrażliwy na szумы.
- Trudny wybór parametrów.
- Powolny - złożoność $O(N^2)$.

Dane wykorzystane przez nas do klasteryzacji zebrane zostały poprzez 27 pytańową ankietę, zawierającą odpowiedzi w postaci liczb całkowitych.

<https://goo.gl/forms/aDAEtAKZTxXPf5j32>

Dodatkowo wykorzystaliśmy podobne dane, które zebrane zostały w Portugalii.

<https://www.kaggle.com/uciml/student-alcohol-consumption>

W naszej „analizie” zwróciliśmy uwagę na współczynniki:

- Homogeneity - h - każdy klaster zawiera tylko dane z jednego zestawu
- Completeness - c - wszystkie dane z jednego zestawu znajdują się w tej samej klasie
- V-measure - v - średnia harmoniczna h oraz c

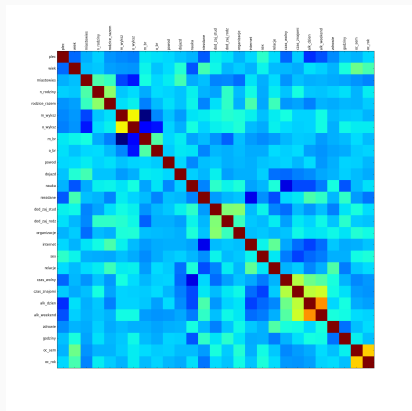
Warunkowa entropia klas przy chwilowym (w danym kroku algorytmu) przydzieleniu do klastrów:

$$H(C|K) = - \sum_{c=1}^{|C|} \sum_{k=1}^{|K|} \frac{n_{c,k}}{n} \cdot \log \left(\frac{n_{c,k}}{n_k} \right) \quad (1)$$

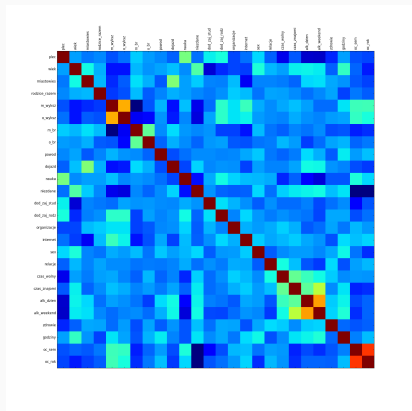
Entropia klas:

$$H(C) = - \sum_{c=1}^{|C|} \frac{n_c}{n} \cdot \log \left(\frac{n_c}{n} \right) \quad (2)$$

$$h = 1 - \frac{H(C|K)}{H(C)} \quad c = 1 - \frac{H(K|C)}{H(K)} \quad v = 2 \cdot \frac{h \cdot c}{h + c} \quad (3)$$

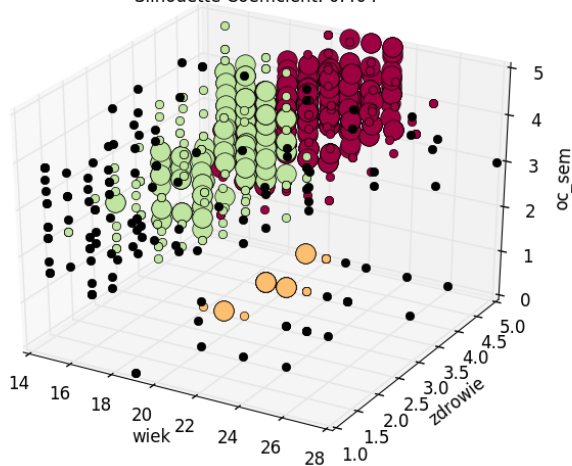


Rys. 3: Wykres korelacji dla naszych danych

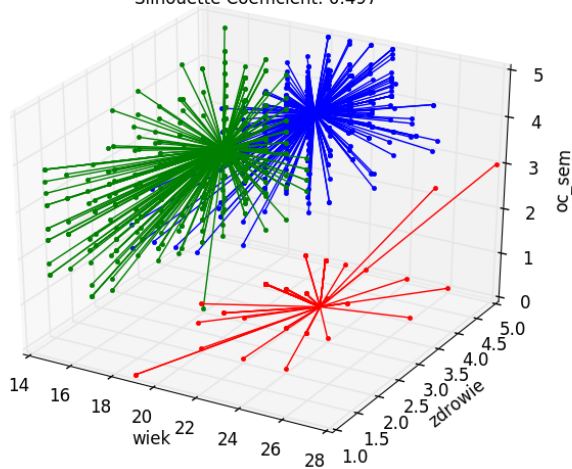


Rys. 4: Wykres korelacji dla danych z Portugalii

Clusters: 3
Homogeneity: 0.793 Completeness: 0.479 V-measure: 0.597
Adjusted Rand Index: 0.644 Adjusted Mutual Information: 0.478
Silhouette Coefficient: 0.404

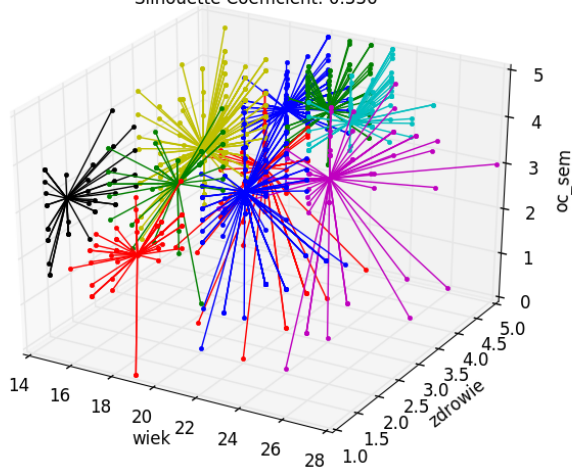


Clusters: 3
Homogeneity: 0.878 Completeness: 0.669 V-measure: 0.759
Adjusted Rand Index: 0.798 Adjusted Mutual Information: 0.668
Silhouette Coefficient: 0.497



AFFINITY PROPAGATION

Clusters: 10
Homogeneity: 0.771 Completeness: 0.239 V-measure: 0.365
Adjusted Rand Index: 0.191 Adjusted Mutual Information: 0.237
Silhouette Coefficient: 0.336



Dziękujemy za uwagę!