

UCZNIE NIENADZOROWANE

ALGORYTM KMEANS

Jednym z zadań uczynienia nienadzorowanego jest tzw. proces analizy skupisk. Proces, w trakcie którego odnajdujemy w nieoznakowanych danych obiekty podobne i 'zamykamy' je w klastry. W trakcie tzw. ostrej klasteryzacji obiekt może należeć tylko do jednego klastra. W klastrze znajdują się obiekty podobne sobie względem przyjętej metryki (np. Euklidesowa). Obiektem może być np.: pacjent w klinice onkologicznej, klient sklepu internetowego a celem procesu uczenia odnalezienie podobnym obiektów. Przykład klasteryzacji dla danych dwuwymiarowych został przedstawiony na rys 1a. Jednym z prostszych algorytmów klasteryzacji jest algorytm kMeans, jego pseudokod został przedstawiony na rys 1b.

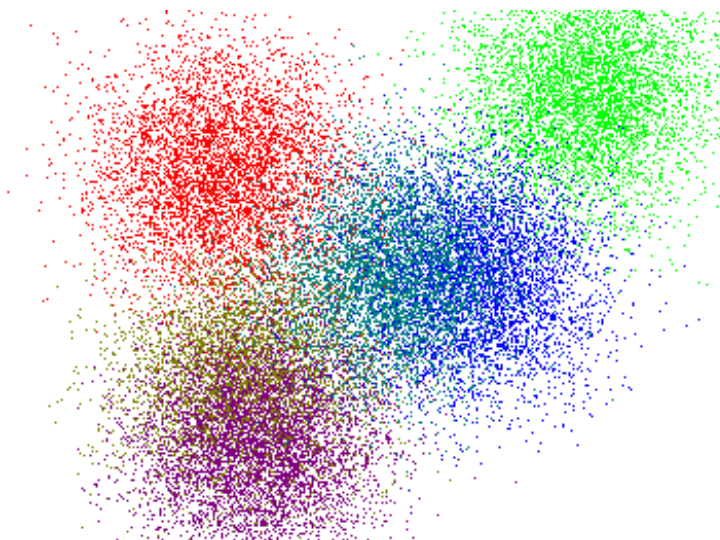
N: number of data objects

K: number of clusters

objects[N]: array of data objects

clusters[K]: array of cluster centers

membership[N]: array of object memberships



```

kmeans_clustering( )
1  while  $\delta/N > \text{threshold}$ 
2     $\delta \leftarrow 0$ 
3    for  $i \leftarrow 0$  to  $N-1$ 
4      for  $j \leftarrow 0$  to  $K-1$ 
5         $\text{distance} \leftarrow | \text{objects}[i] - \text{clusters}[j] |$ 
6        if  $\text{distance} < d_{\min}$ 
7           $d_{\min} \leftarrow \text{distance}$ 
8           $n \leftarrow j$ 
9        if  $\text{membership}[i] \neq n$ 
10          $\delta \leftarrow \delta + 1$ 
11          $\text{membership}[i] \leftarrow n$ 
12          $\text{new\_clusters}[n] \leftarrow \text{new\_clusters}[n] + \text{objects}[i]$ 
13          $\text{new\_cluster\_size}[n] \leftarrow \text{new\_cluster\_size}[n] + 1$ 
14      for  $j \leftarrow 0$  to  $K-1$ 
15         $\text{clusters}[j][*] \leftarrow \text{new\_clusters}[j][*] / \text{new\_cluster\_size}[j]$ 
16         $\text{new\_clusters}[j][*] \leftarrow 0$ 
17         $\text{new\_cluster\_size}[j] \leftarrow 0$ 

```

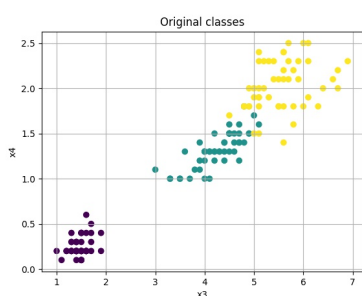
Rys 1 a) Przykład klasteryzacji dla danych dwuwymiarowych; b) pseudokod algorytmu kMeans

OCENA KLASTERYZACJI

Jakość klasteryzacji można ocenić z wykorzystaniem tzw. indeksów klasteryzacji. W zależności czy posiadamy informacje o prawdziwej przynależności obiektów (w praktyce rzadko) czy nie możemy skorzystać z takich wskaźników –wymieniając tylko kilka z wielu- jak homogeniczności, przyległości Randa¹, indeks Dunna² czy współczynnik Silhouette³.

ZADANIA DO WYKONANIA

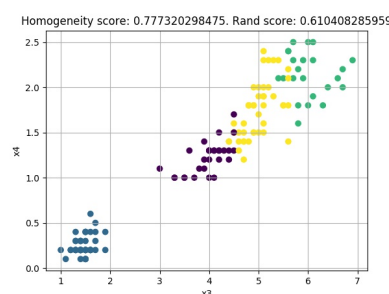
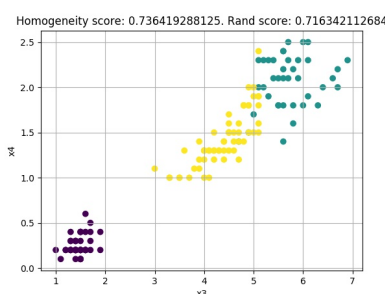
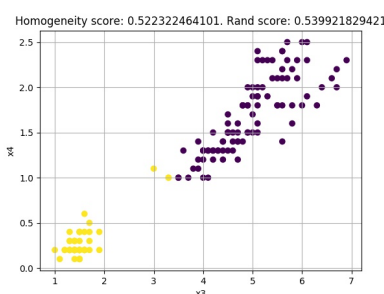
Korzystając z biblioteki `sklearn` (<http://scikit-learn.org>) przetestuj działanie algorytmu `kMeans` na danych "Iris". Sprawdź wyniki dla $k=2$, $k=3$, $k=4$ oraz oceń jakość uzyskanych klastrów korzystając z wskaźnika homogeniczności oraz przyległości Randa.



W celu realizacji zadania należy skorzystać z modułów:

- `sklearn.cluster.KMeans`
- `sklearn.metrics.homogeneity_score`
- `sklearn.metrics.adjusted_rand_score`

Rys 2. Dane z zaznaczonymi prawdziwymi klasami



Rys 3. Przykładowe rozwiązanie dla $k=2$, $k=3$, $k=4$ oraz wskaźnik homogeniczności i Randa

¹ <http://faculty.washington.edu/kayee/pca/supp.pdf>

² <http://www.sciencedirect.com/science/article/pii/S0031320303002838>

³ <http://www.sciencedirect.com/science/article/pii/0377042787901257?via%3Dihub>