

Eksploracja danych

7. Grupowanie

Piotr Szwed
Katedra Informatyki Stosowanej AGH
2021

Grupowanie

Grupa obiektów (ang. cluster) to zbiór obiektów wyodrębniony z danych

- Obiekty wewnątrz grupy powinny cechować się podobieństwem
- Obiekty różnych grup powinny być niepodobne do siebie

Grupowanie, klasteryzacja, analiza skupień (ang. clustering, cluster analysis) ma na celu **odkrywanie wiedzy**: znalezienie interesujących wzorców wewnątrz danych, np. często występujących kombinacji atrybutów lub grup podobnych obiektów.

Wynikiem grupowania jest przydział etykiet (numerów grup) do obiektów.

Porównanie

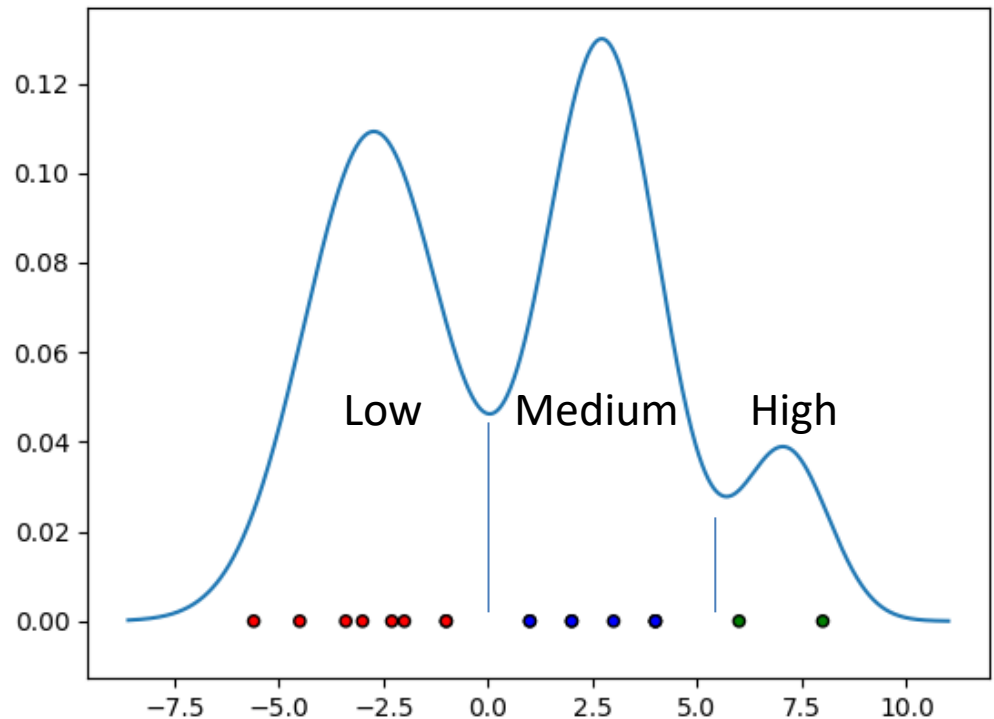
Zagadnienie	Dane uczące	Cel	Metoda oceny
Regresja	$D = \{(x_i, y_i)\}$ $x_i \in X$ $y_i \in R$	Znalezienie funkcji $f: X \rightarrow R$ odwzorowującej całą przestrzeń X	Zbiór testowy Sprawdzan krzyżowy (krosvalidacja)
Klasyfikacja	$D = \{(x_i, y_i)\}$ $x_i \in X$ $C = \{c_1, \dots, c_k\}$ – etykiety klas $y_i \in C$	Znalezienie funkcji $f: X \rightarrow C$ odwzorowującej całą przestrzeń X	Zbiór testowy Sprawdzan krzyżowy (krosvalidacja)
Grupowanie	$D = \{x_i\}$ $x_i \in X$	Znalezienie funkcji $f: D \rightarrow L$ która przypisuje etykiety ze zbioru L instancjom w D	Wartość założonej funkcji celu lub wskaźnika jakości, na podstawie dostępnych etykiet $Y = y_i$, użyteczność w dalszym przetwarzaniu

Uczenie nienadzorowane

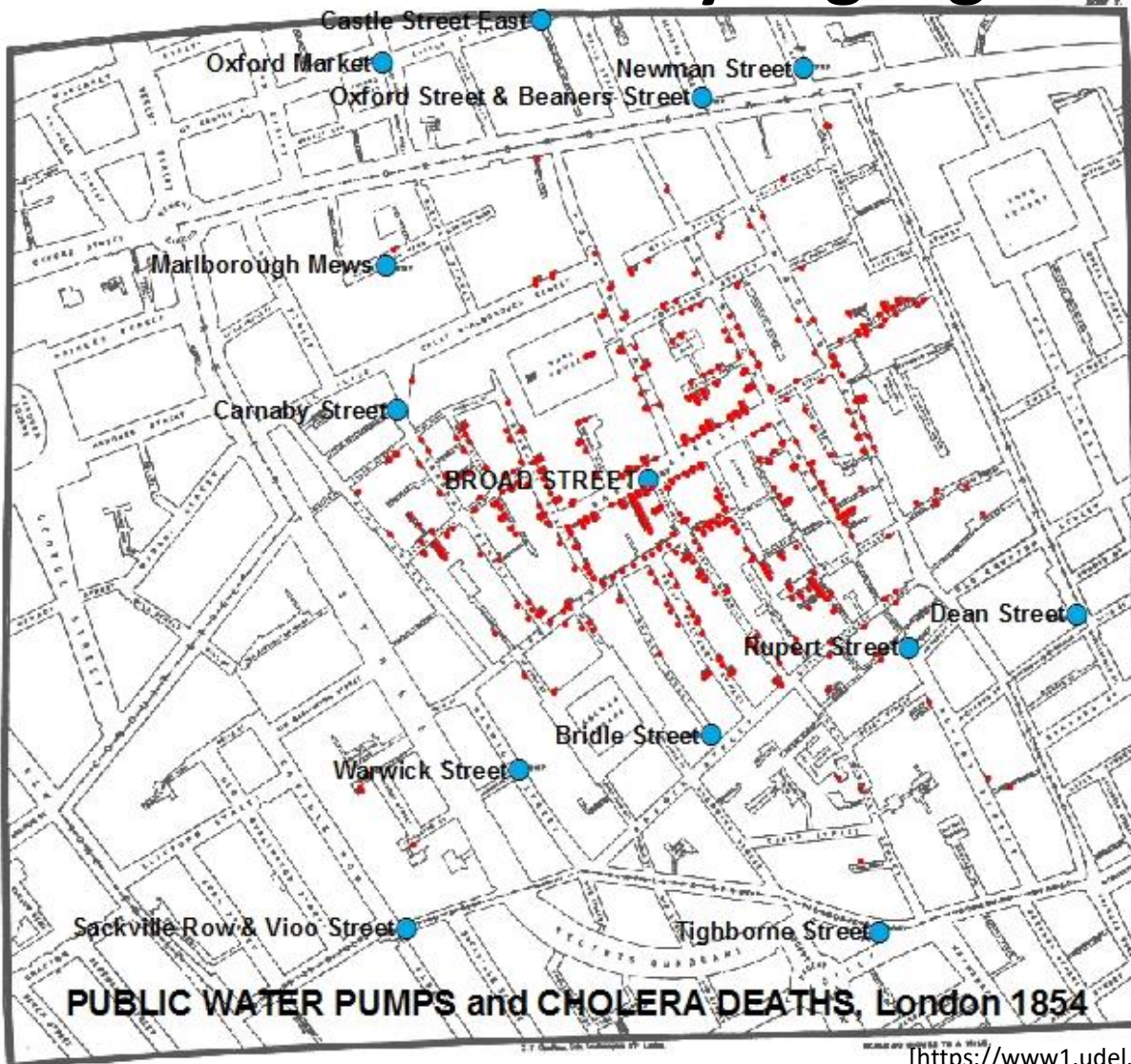
- Jeżeli w procesie uczenia i budowy modeli wykorzystuje się znane wartości wyjściowe y_i dla obserwacji (x_i, y_i) , porównuje się je z wyjściami modelu i koryguje jego parametry – jest to **uczenie nadzorowane**.
- **Klasyfikacja** i **regresja** to przykłady zagadnień uczenia nadzorowanego
- Jeżeli wartości wyjściowe y_i nie są znane/wykorzystywane, podczas budowy modelu lub jego działania jest to **uczenie nienadzorowane**.
- Podczas **grupowania** nie są znane wartości wyjściowe y_i (ciągłe lub etykiety klas). Jest to więc przykład uczenia nienadzorowanego.
- Czy kNN (metodę klasyfikacji) należy zaliczyć do grupy algorytmów uczenia nadzorowanego?

Zastosowania

- Wykorzystywane jako samodzielne narzędzie analizy danych np. do znalezienia interesujących wzorców i rozkładów statystycznych
- Jako wstępny etap przetwarzania danych:
 - Grupowanie danych numerycznych, ustalanie przedziałów i zamiana na nominalne
 - Identyfikacja skupisk w celu automatycznego nadania obserwacjom etykiet klas



Przykłady i zastosowania: analiza danych geograficznych



Skupiska cholery
podczas epidemii w
1854 roku



[<http://tasteforhealth.com/1854-broad-street-cholera-outbreak/>]

[<https://www1.udel.edu/johnmack/frec682/cholera/cholera2.html>]

Przykłady i zastosowania: wyszukiwanie dokumentów tekstowych

The screenshot shows the eTools Web Search interface. At the top, there is a search bar with the query 'trump' and a 'Search' button. Below the search bar, there are navigation links for 'About', 'Search feeds', 'Download', 'Carrot Search', and 'Contact'. The main content area is divided into two sections: a visualization on the left and a list of search results on the right.

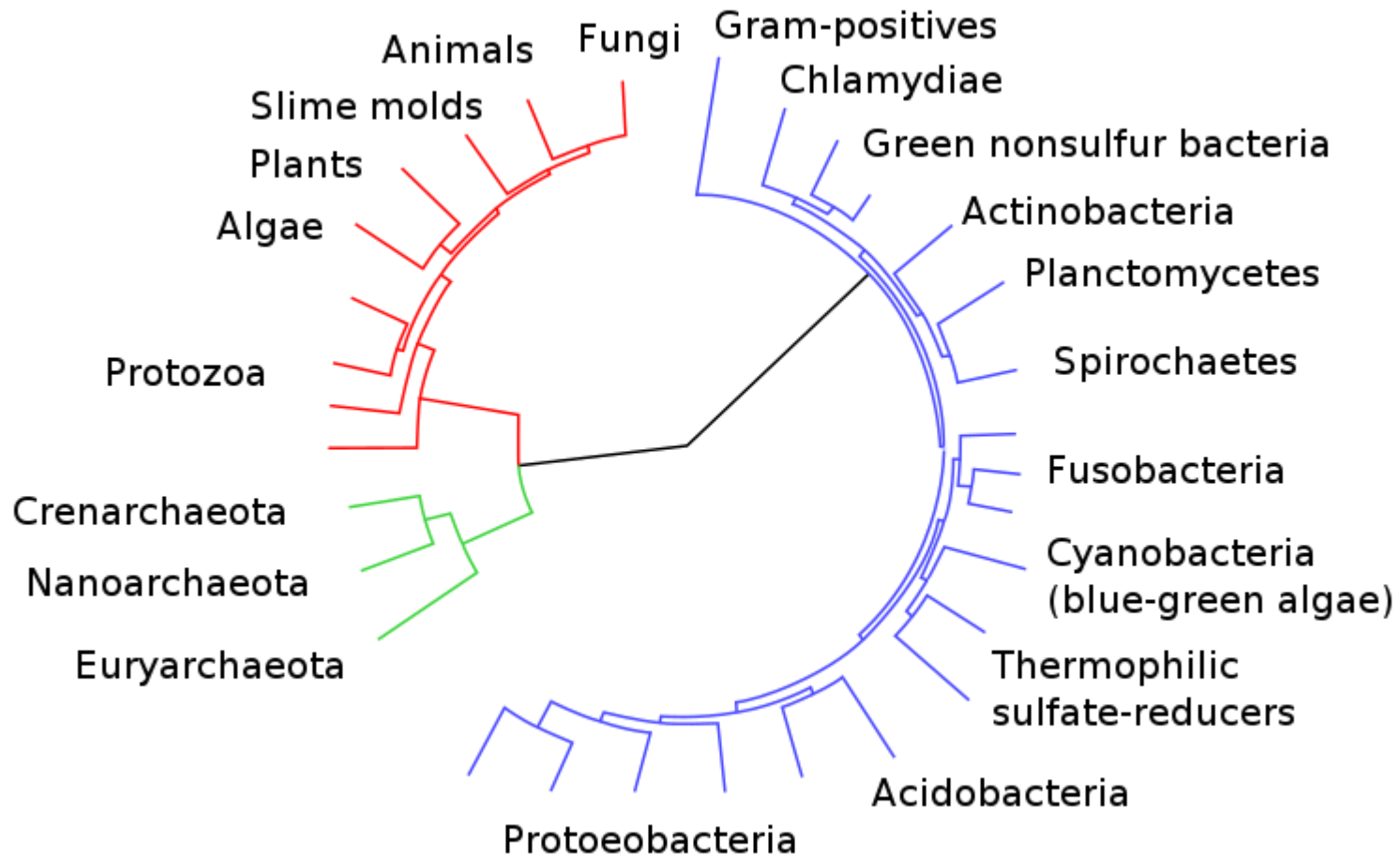
The visualization is a FoamTree cluster, which is a hierarchical tree structure where nodes are represented by colored polygons. The central node is 'President of the United States'. Other nodes include 'Deal', 'Tariffs', 'NOTE Pinyininyin Tone Marking 2016 11', 'Fought', 'Latest News', 'Atlantic', 'Russia', 'Spokesman', 'Campaign', 'White House', 'Family', 'Daily', 'Pinyin Trump', 'America Great', 'Ivana', 'New York', 'Trump Says', 'John', 'China', 'Topic Galleries', 'Other Topics', 'Melania Trump', 'Stormy Daniels Lawyer', 'Case', '45th President of the United States', 'Leader', and 'FoamTree'.

The search results on the right are listed as follows:

- 1 [Donald Trump - Wikipedia](#)
Donald John **Trump** (born June 14, 1946) is the 45th and current President of the United States, in office since January 20, 2017. Before entering politics, he was a businessman and television personality. **Trump** was born and raised in the New York ...
https://en.wikipedia.org/wiki/Donald_Trump [Ask, Google, Wikipedia]
- 5 [Donald John Trump | TheHill](#)
President of the United States, 2017 - PresentThe **Trump** Organization, Founder/Chair/President/Chief Executive Officer, 1975 - PresentThe Apprentice, Executive Producer, 2003 - 2015 ...
<http://thehill.com/people/donald-trump> [Bing, Yahoo]
- 13 [Donald Trump sexism tracker: Every offensive comment in one place](#)
Fat. Pig. Dog. Slob. Disgusting animal. These are just some of the names that Donald **Trump** has called women. The President of the United States has been ...
<https://www.telegraph.co.uk/women/politics/donald-trump-sexism-tracker-every-offensive-comment-in-one-place/> [Bing, Yahoo]
- 19 [Trump](#)
Trump most commonly refers to: Donald **Trump** (born 1946), 45th President of the United States, businessman, and television personality **Trump** (card games) ...

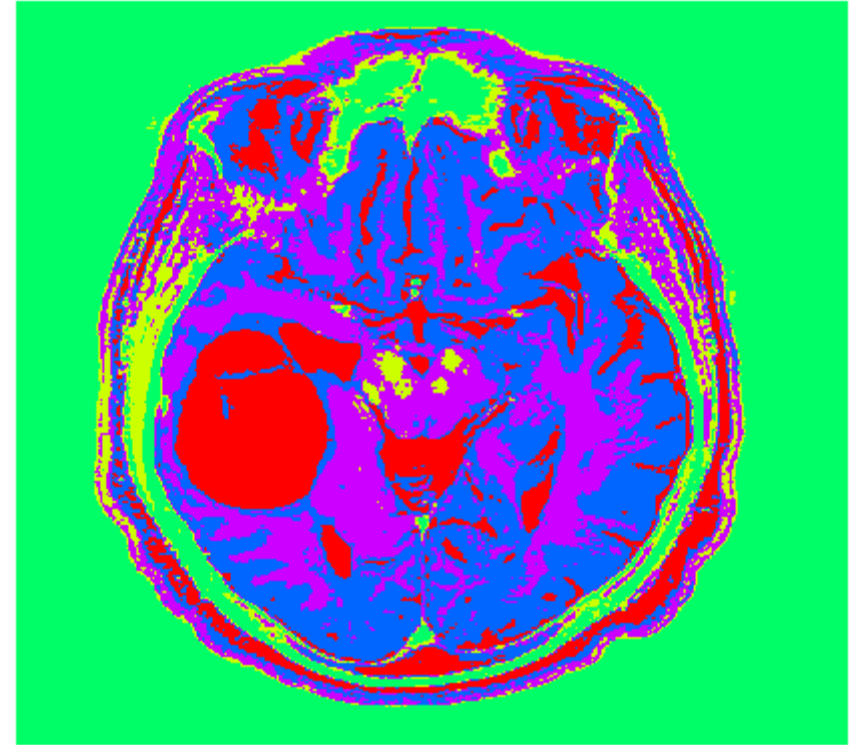
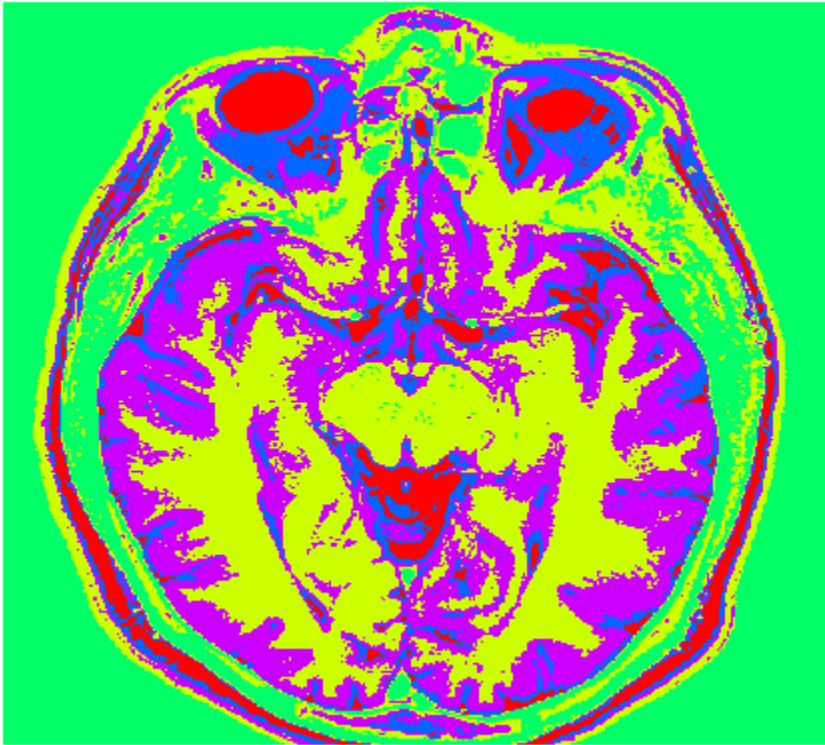
At the bottom of the page, there is a footer with the following information: Query: trump -- Source: eTools Web Search (100 results, 0 ms) -- Clusterer: Lingo v3.15.2-SNAPSHOT | build 5 | 2017-09-11 12:39 © 2002-2018 Stanislaw Osinski, Dawid Weiss

Przykłady i zastosowania: tworzenie taksonomii



[[https://en.wikipedia.org/wiki/Taxonomy_\(biology\)](https://en.wikipedia.org/wiki/Taxonomy_(biology))]

Przykłady i zastosowania: segmentacja obrazów



[<http://mpmendespt.github.io/MRI-image-segmentation.html>]

Zastosowania

- Marketing – poszukiwanie grup podobnych klientów w celu stworzenia dedykowanych programów marketingowych
- Ubezpieczenia – identyfikacja grup klientów o podobnych profilach ryzyka
- Służba zdrowia – analiza kosztów dla różnych wariantów terapii
- Planowanie transportu - analiza ruchu pasażerów
- Rynek nieruchomości – grupowanie ofert o podobnych parametrach

Metody grupowania

- Hierarchiczne
- Rozłączne
- Oparte na modelach statystycznych
- Wykorzystujące informacje o gęstości instancji
- Rozmyte

Grupowanie hierarchiczne

Grupowanie hierarchiczne (aglomeratywne)

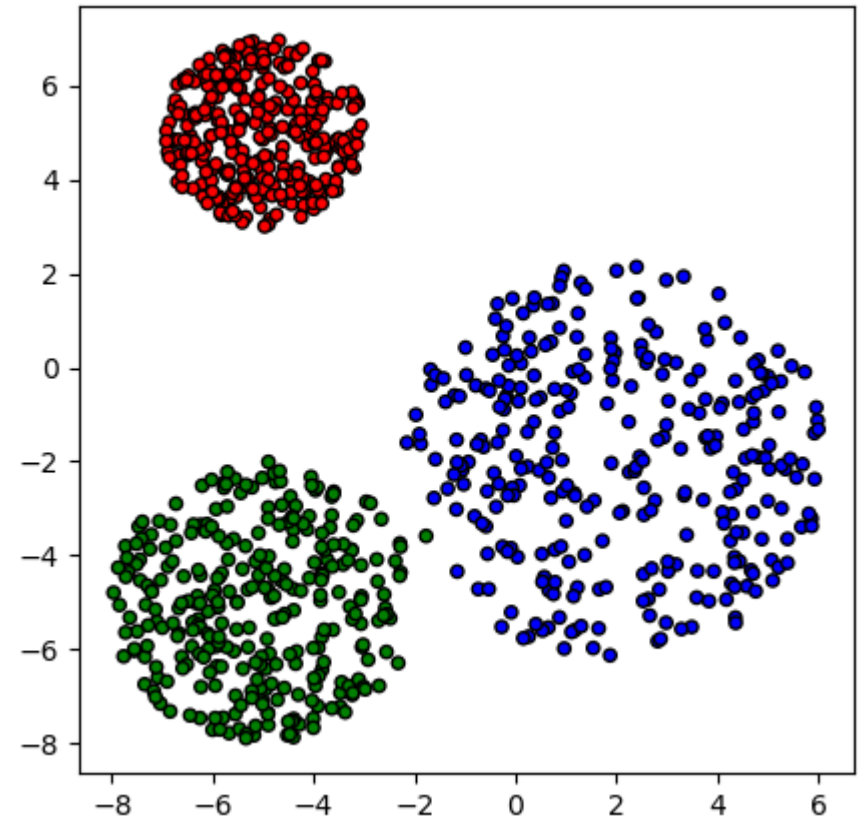
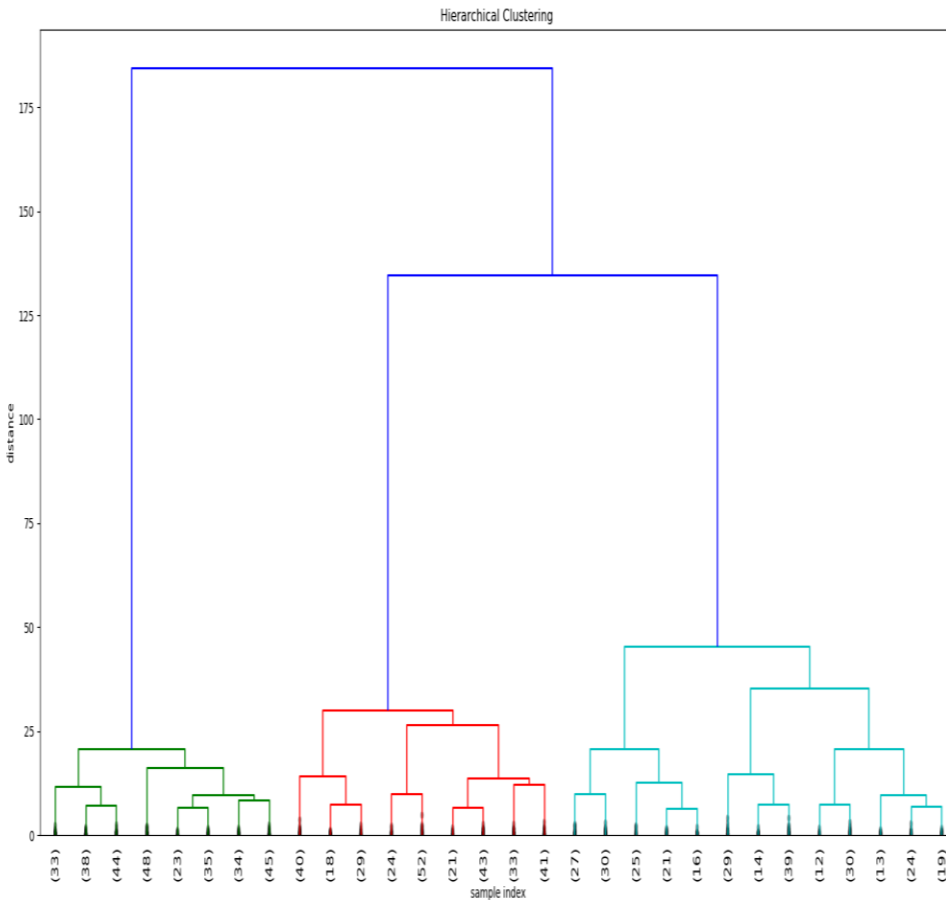
Algorytm

Dane:

- $D = \{x_1, \dots, x_m\} \subset X$
 - $d(A, B)$, gdzie $A, B \subset X$, odległość zbiorów A i B
1. Utwórz tyle grup, ile jest obserwacji $\mathbf{G} \leftarrow \{C_1, \dots, C_m\}$, gdzie $C_i = \{x_i\}$
 2. Dopóki $|\mathbf{G}| > 1$:
 3. Znajdź takie dwie grupy C_r i C_s należące do \mathbf{G} , że odległość między nimi jest najmniejsza
$$(C_r, C_s) = \arg \min_{\substack{(A,B) \in \mathbf{G} \times \mathbf{G} \\ A \neq B}} d(A, B)$$
 4. Połącz grupy: $\mathbf{G} \leftarrow \mathbf{G} \setminus \{C_r, C_s\} \cup \{C_{rs}\}$, gdzie $C_{rs} = C_r \cup C_s$

Informacja o połączeniach grup jest zapisywana: grupa C_{rs} jest nadrzędną względem C_r i C_s

Przykład



Rysunek po prawej stronie pokazuje przydział etykiet do obserwacji na trzecim poziomie dendrogramu.

Odległość zbiorów obserwacji

Może zostać zdefiniowana na wiele sposobów:

- Single link: odległość między najbliższymi punktami

$$d(A, B) = \min_{\substack{x_i \in A \\ x_j \in B}} \text{dist}(x_i, x_j)$$

- Complete link: odległość między najdalszymi punktami

$$d(A, B) = \max_{\substack{x_i \in A \\ x_j \in B}} \text{dist}(x_i, x_j)$$

- Średnia odległość

$$d(A, B) = \sum_{\substack{x_i \in A \\ x_j \in B}} \frac{\text{dist}(x_i, x_j)}{|A| \cdot |B|}$$

Odległość zbiorów obserwacji

- Odległość pomiędzy środkami grup (centroidami)

$$d(A, B) = \text{dist}(\bar{x}_A, \bar{x}_B)$$

$$\text{gdzie } \bar{x}_A = \frac{\sum_{x \in A} x}{|A|}, \bar{x}_B = \frac{\sum_{x \in B} x}{|B|}$$

- Metoda ważona. Zakładając, że grupa A została uformowana przez połączenie dwóch grup A_1 i A_2

$$d(A, B) = \frac{d(A_1, B) + d(A_2, B)}{2}$$

Metoda Warda

- Metoda Warda wybierając grupy do złączenia posługuje się kryterium najmniejszej sumy kwadratów błędów (alternatywnie wariancji -- jak podzielimy RSS przez n).

$$RSS(A) = \sum_{x \in A} (x - \bar{x}_A)^T \cdot (x - \bar{x}_A)$$

- Poszukując grup do połączenia wybierane są te dwie, dla których wzrost kwadratu błędu (wariancji) będzie najmniejszy $(C_i, C_j) = \arg \min_{C_i, C_j \in G} RSS(C_i \cup C_j) - RSS(C_i) - RSS(C_j)$

- Przyrost kwadratu błędu można także obliczać jako:

$$\begin{aligned} \Delta(C_i, C_j) &= RSS(C_i \cup C_j) - RSS(C_i) - RSS(C_j) \\ &= \frac{|C_i| \cdot |C_j|}{|C_i| + |C_j|} \left\| \bar{x}_{C_i} - \bar{x}_{C_j} \right\|^2 \end{aligned}$$

Algorytm k-means k-średnich

Środek grupy - centroid

Założmy, że $x \in R^n$. Punkt środkowy (centroid) grupy $C_i = \{x_{i1}, \dots, x_{ir}\}$ obliczamy jako:

$$\bar{x}(i) = \frac{\sum_{x \in C_i} x}{|C_i|}$$

W idealnym przypadku grupa C_i zawiera się wewnątrz kuli, której środkiem jest centroid $\bar{x}(i)$ oraz kule nie przecinają się.

Miarą, która pozwala określić jakość reprezentacji grupy przez centroid jest suma kwadratów błędów (*residual sum of squares*)

$$RSS(i) = \sum_{x \in C_i} \|x - \bar{x}(i)\|^2$$

Miarą jakości podziału na grupy jest:

$$RSS = \sum_{i=1}^k RSS(i) = \sum_{i=1}^k \sum_{x \in C_i} \|x - \bar{x}(i)\|^2$$

Podział na grupy

- Niech $C = (C_1, \dots, C_k)$ będzie podziałem na grupy zbioru $D = \{x_j: j = 1, m\}$ spełniającym:

$\forall r \neq s: C_r \cap C_s = \emptyset$ oraz

$$\bigcup_{i=1}^k C_i = D$$

- Najlepszy podział na grupy minimalizuje RSS , stąd:

$$(C_1, \dots, C_k)^* = \arg \min_C \sum_{i=1}^k \sum_{x \in C_i} \|x - \bar{x}(i)\|^2$$

Zagadnienie można więc potraktować jako problem optymalizacyjny: znaleźć taki podział na rozłączne grupy, które zminimalizuje funkcję celu RSS .

Problem NP zupełny

Inna postać definicji RSS

Wprowadźmy m zmiennych Z_j o wartościach ze zbioru $\{1, \dots, k\}$ opisujących przydział obserwacji do grupy:

$$Z_j = i \Leftrightarrow x_j \in C_i$$

$$RSS = \sum_{i=1}^k \sum_{j=1}^m \mathbf{1}(Z_j = i) \|x_j - \bar{x}(i)\|^2$$

Wartości m zmiennych Z można wybrać na m^k sposobów (i w wielomianowym czasie obliczyć RSS)

- Funkcja $\mathbf{1}(expr)$ przyjmuje wartość 1, gdy $expr$ jest prawdziwe; 0 w przeciwnym przypadku.
- Funkcja celu ma wiele minimów lokalnych, znalezienie rozwiązania dokładnego jest trudne. W zamian używa się algorytmu zachłannego k-means (k-średnich)

Algorytm k-means

Dane

- $D = \{x_i \in R^n\}, i = 1, \dots, m$ – zbiór obserwacji
- k – zadana liczba grup
- $d(x_1, x_2) = \|x_1 - x_2\|^2$ - odległość euklidesowska

Algorytm

1. Wybierz początkowe położenia k centroidów $\bar{x}(1), \dots, \bar{x}(k)$
2. Przypisz obserwacje do najbliższych centroidów:

for $j = 1 \dots m$:

$$Z_j = \arg \min_{i \in \{1, \dots, k\}} d(x_j, \bar{x}(i))$$

3. Uaktualnij położenie centroidów

for $i = 1 \dots k$:

$$\bar{x}(i) = \frac{\sum_{j=1}^m 1(Z_j = i) \cdot x_j}{\sum_{j=1}^m 1(Z_j = i)}$$

Suma wektorów
należących do
grupy

Liczba elementów
w grupie

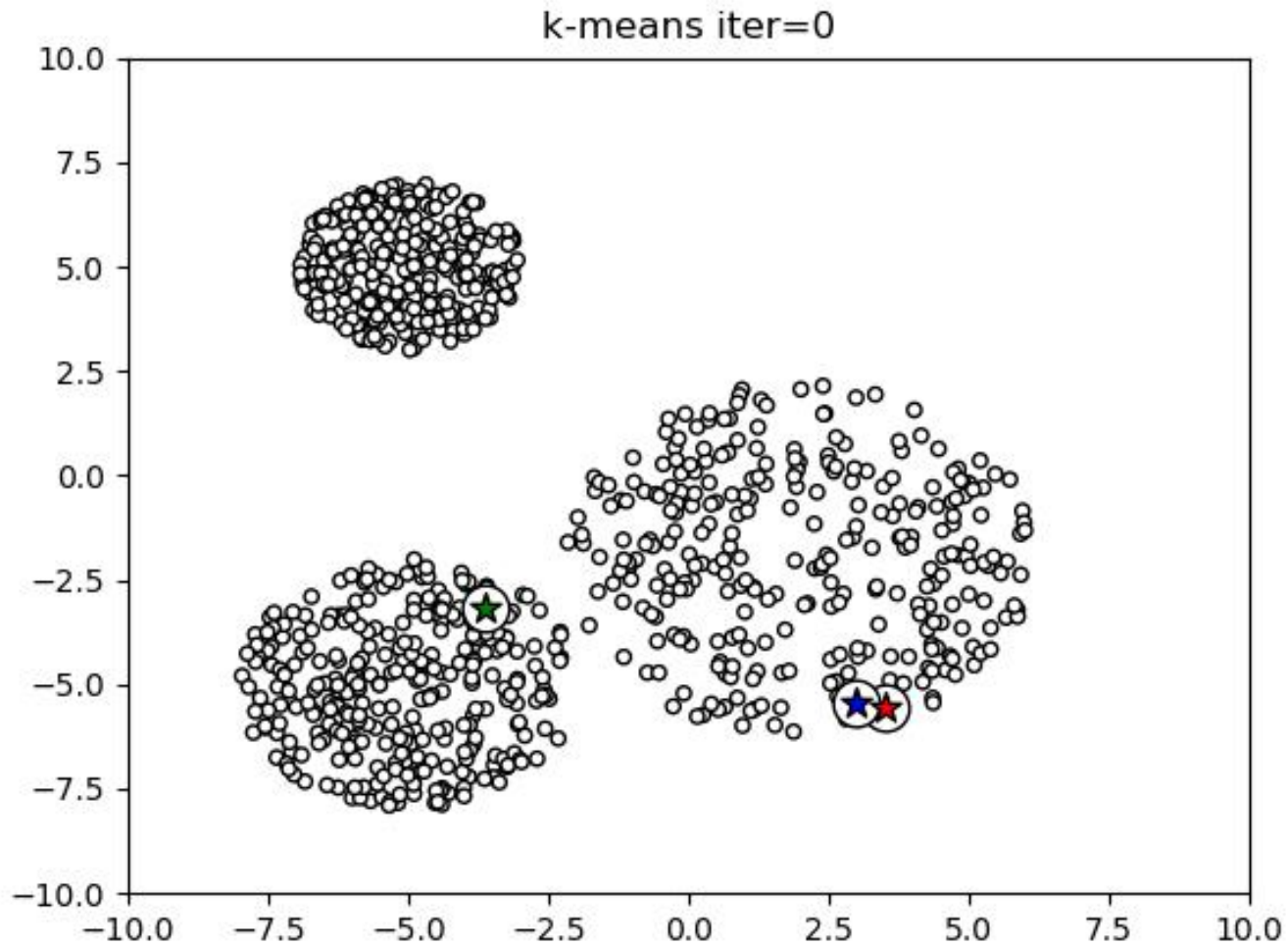
4. Jeżeli nie jest spełnione kryterium stopu – wróć do 2

Algorytm k-means – kryterium stopu

- Brak zmiany przydziału obserwacji do grup (brak zmiany wartości zmiennych Z_j)
- Osiągnięcie maksymalnej liczby iteracji
- Brak poprawy funkcji celu RSS lub względna zmiana mniejsza od zadanego progu:

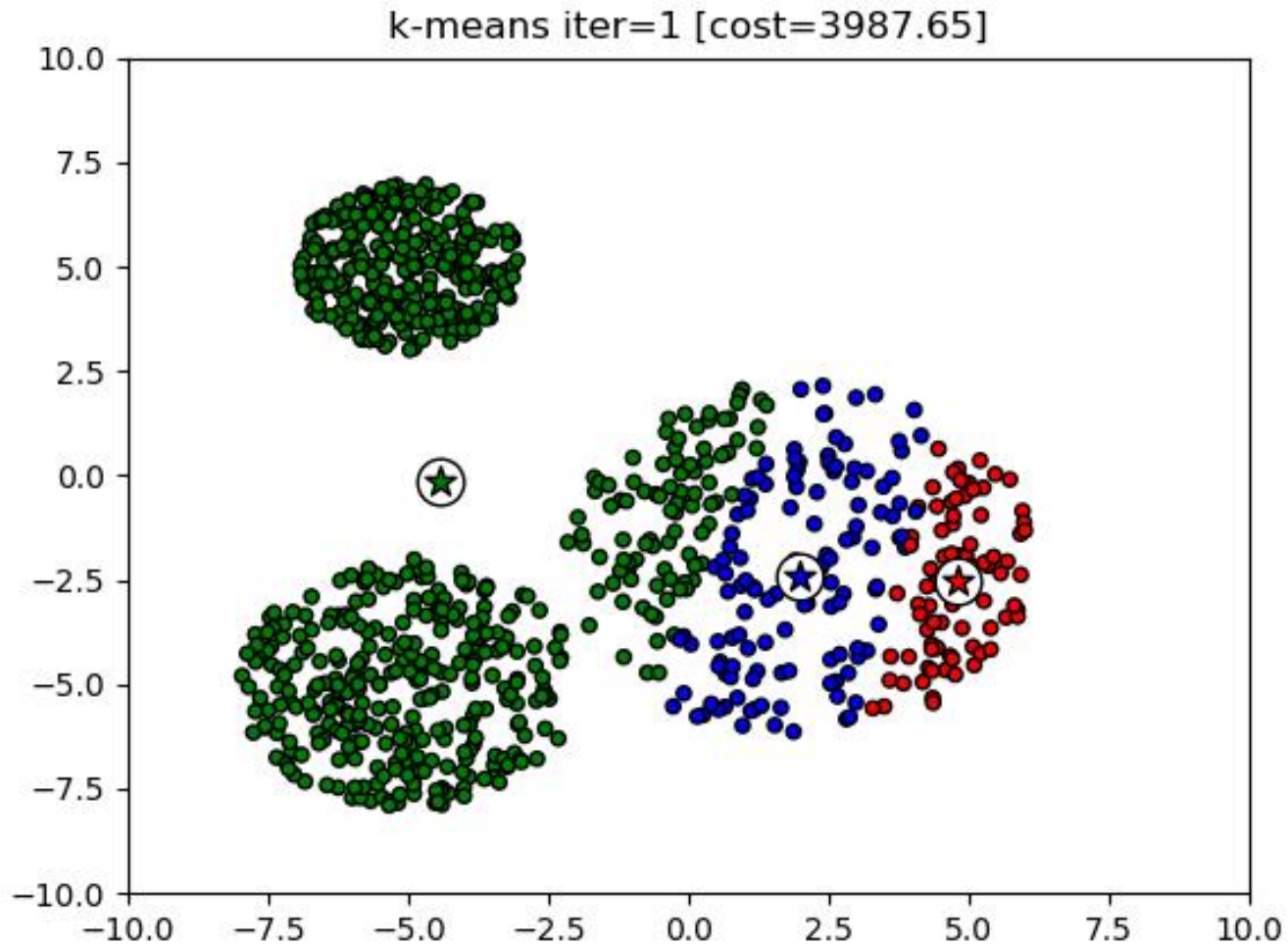
$$\frac{RSS^{(t)} - RSS^{(t-1)}}{RSS^{(t-1)}} < \epsilon$$

Przebieg algorytmu - inicjalizacja



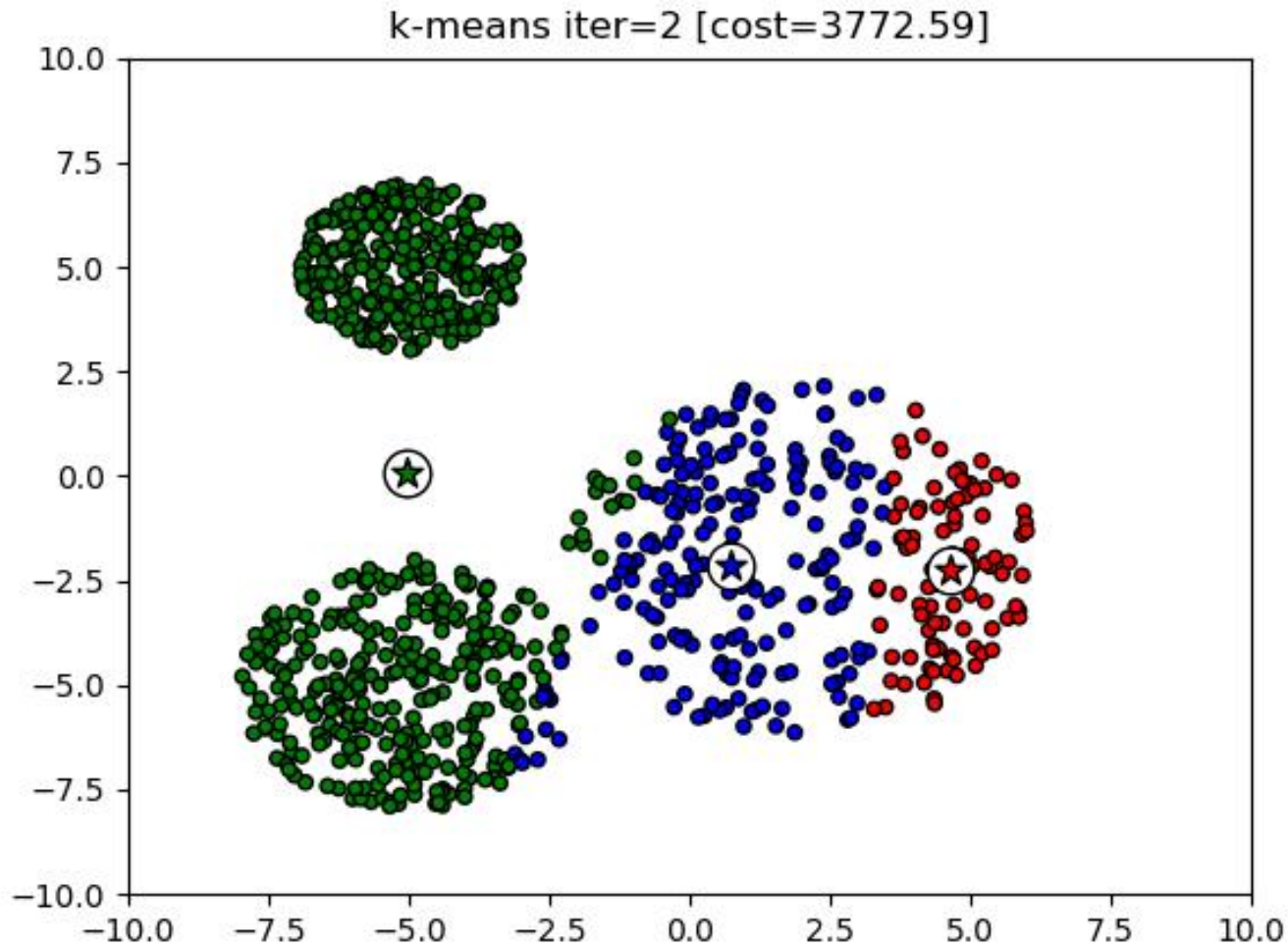
Wylosowano 3 spośród punktów jako środki skupisk

Przebieg algorytmu – po pierwszej iteracji



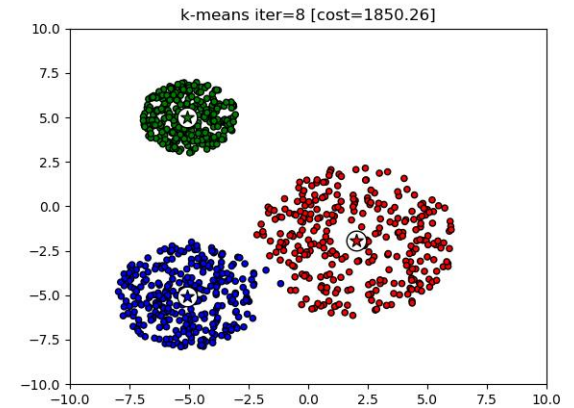
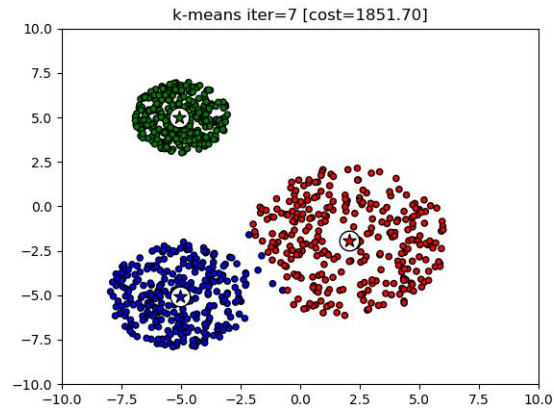
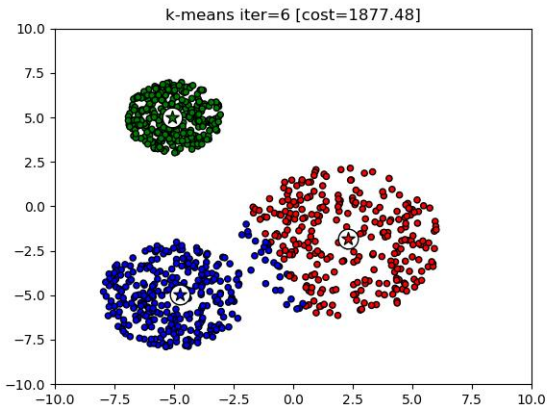
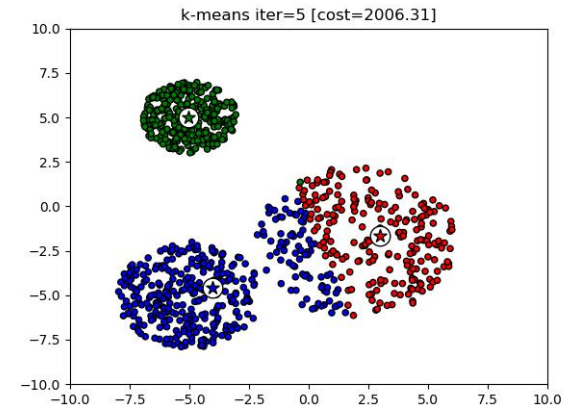
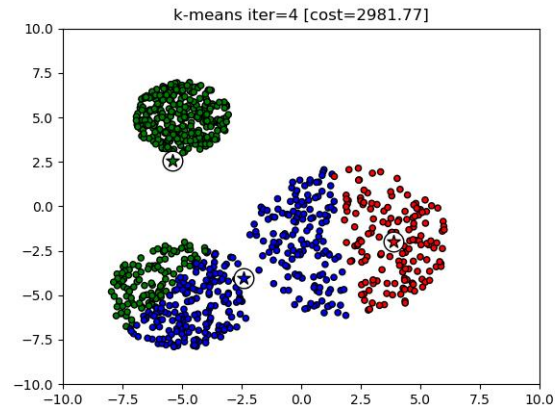
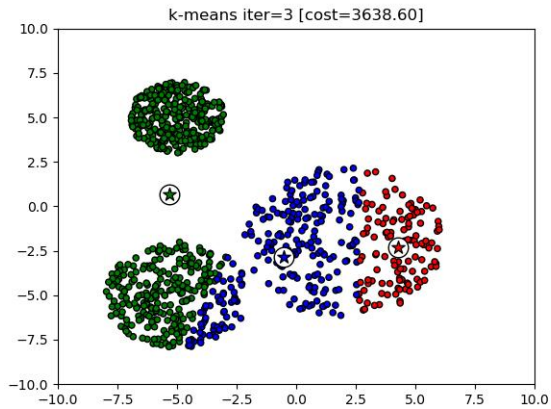
Przesunięcie środków skupisk i uaktualnienie przydziału obserwacji

Przebieg algorytmu – po 2 iteracji

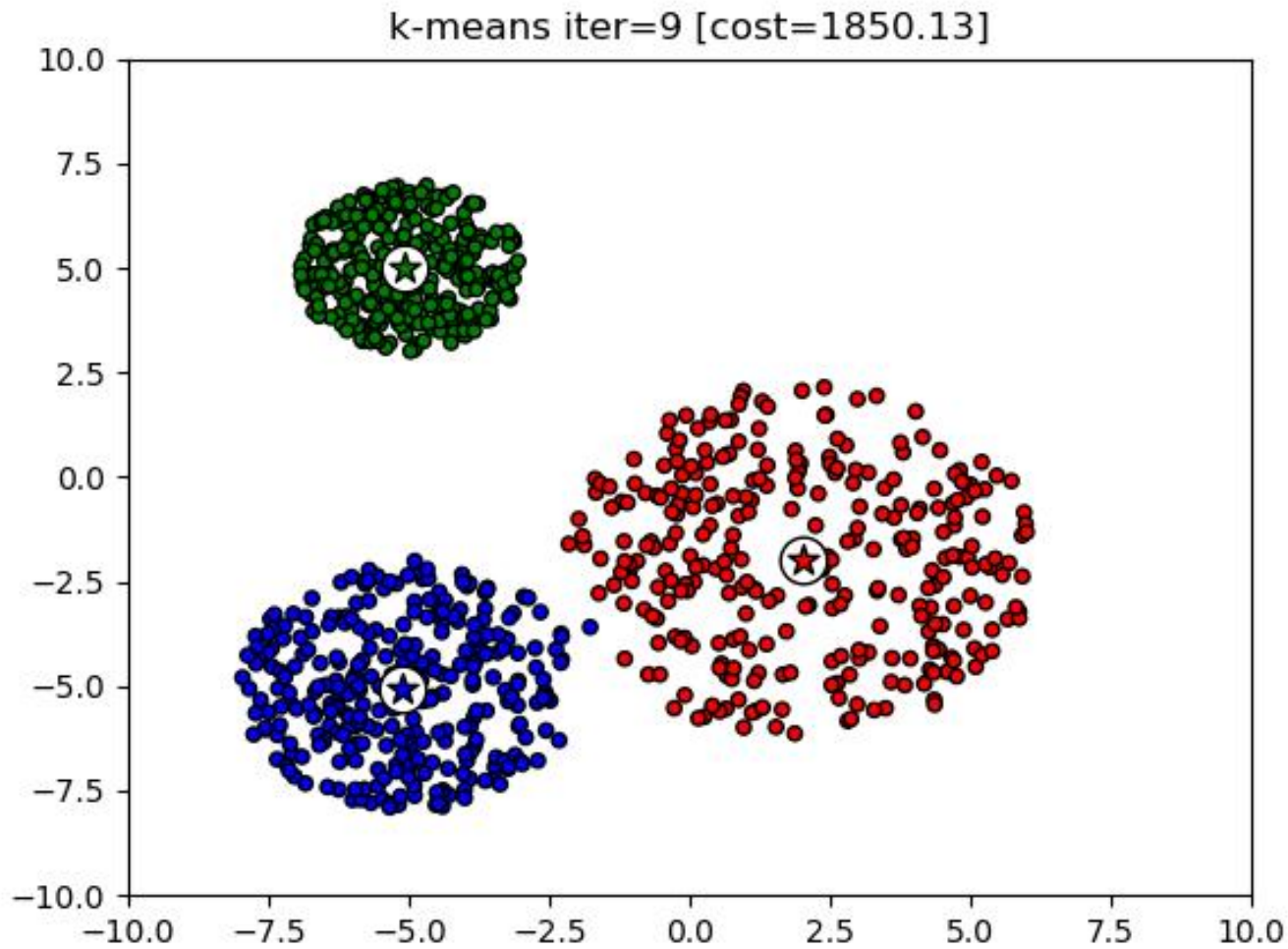


Algorytm jest kontynuowany dopóki następuje zmiana przydziału (etykiety) obserwacji. Zazwyczaj także spada koszt (minimum lokalne)

Przebieg algorytmu – kolejne iteracje

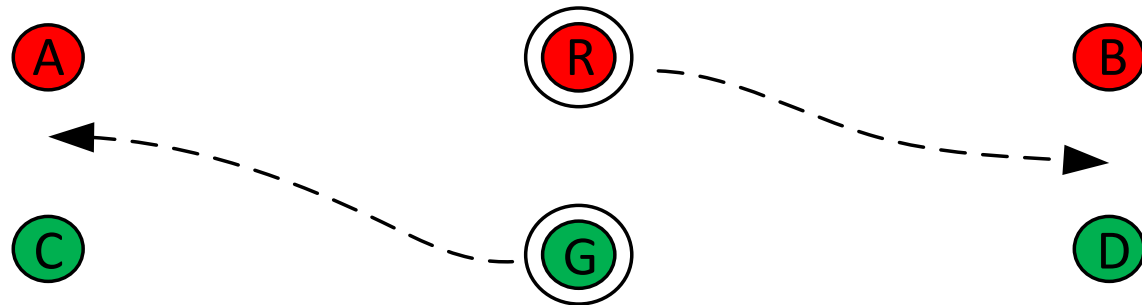


Przebieg algorytmu - koniec



K-means - inicjalizacja

Algorytm jest wrażliwy na wybór początkowych środków skupisk. Niewłaściwy wybór może spowodować utknięcie algorytmu w jednym z minimów lokalnych o bardzo dużej wartości RSS



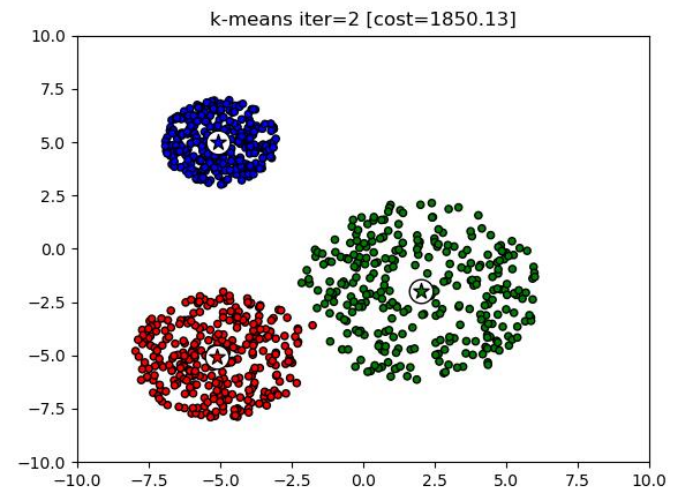
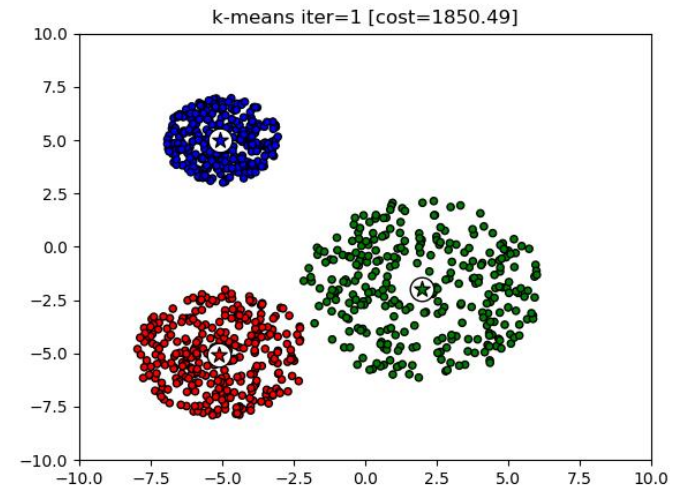
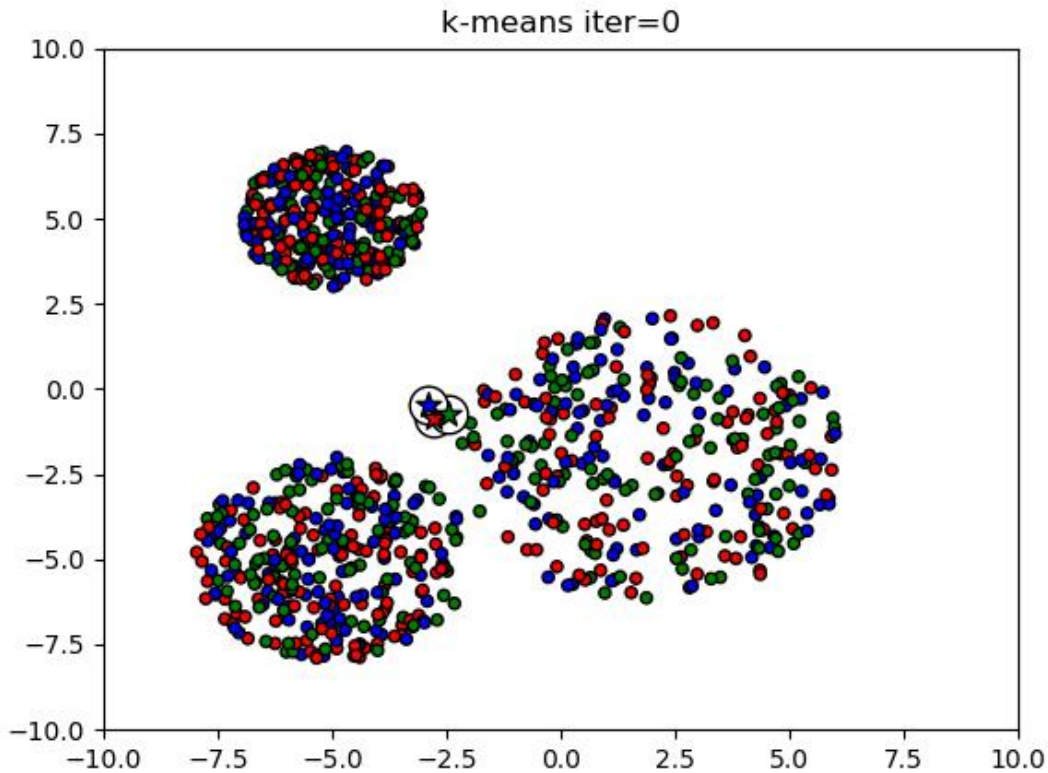
Przykład

- R i G to centroidy będące środkami grup $\{A, B\}$ i $\{C, D\}$
- Optymalny podział na grupy to $\{A, C\}$ i $\{B, D\}$.
- Jednakże zmiana podziału nie jest możliwa, np. punkt A nigdy nie zmieni grupy, ponieważ $|AG| > |AR|$. Podobnie zachowują się pozostałe punkty B, C i D .

Trzy strategie inicjalizacji

- Losowanie obserwacji (pokazane wcześniej)
- Losowy przydział obserwacji do skupisk
- K-means++
Losowy wybór kolejnych punktów z prawdopodobieństwem
wzrastającym wraz z odległością od wybranych środków

Inicjalizacja – losowy przydział obserwacji



Przy losowaniu przydziału do grup zazwyczaj centroidy są położone w pobliżu środka masy całego zbioru danych

Algorytm inicjalizacji k-means++

1. Przypisz $i \leftarrow 1$. Wylosuj środek $\bar{x}(1)$.

2. Przypisz $i \leftarrow i + 1$

3. Dla każdej obserwacji x_j oblicz najmniejszą odległość od ustalonych środków:

$$d_{min}(x_j) \leftarrow \min\{d(x_j, \bar{x}(r)): r = 1, i - 1\}$$

4. Przypisz każdej obserwacji prawdopodobieństwo wyboru

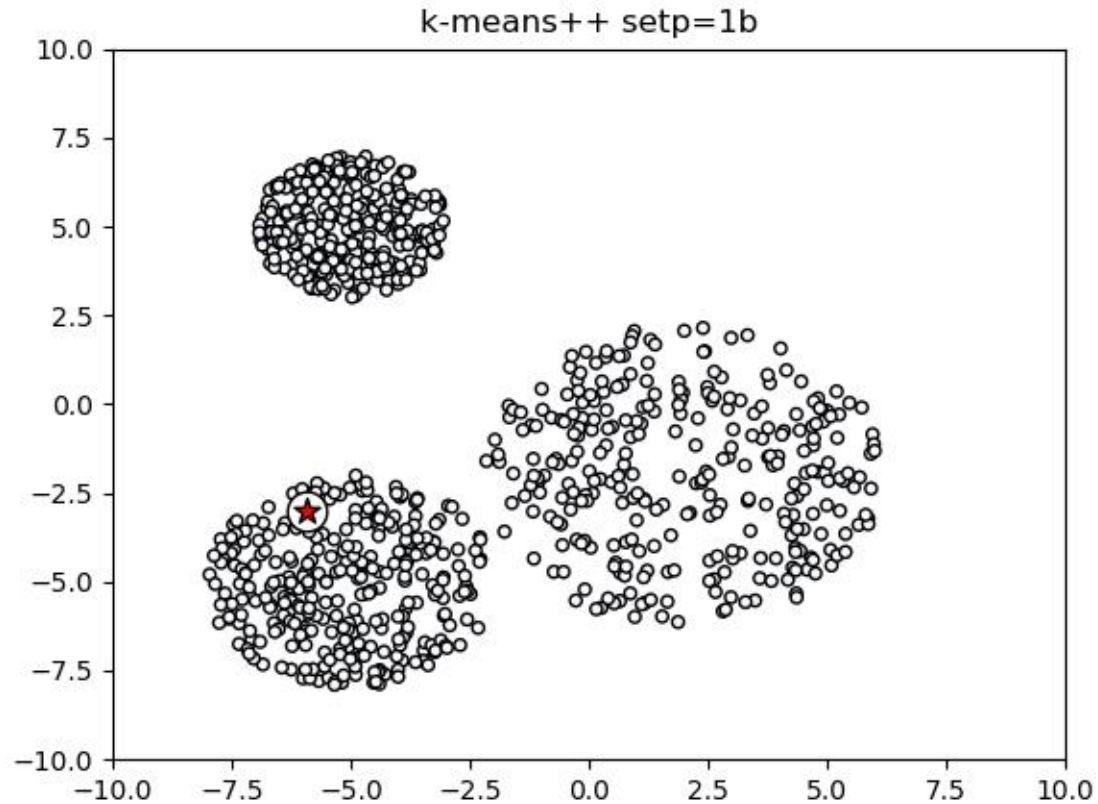
$$P(x_j) \sim d_{min}(x_j)^2$$

5. Wylosuj $\bar{x}(i)$

6. Jeżeli $i = k$, STOP. W przeciwnym przypadku przejdź do 2

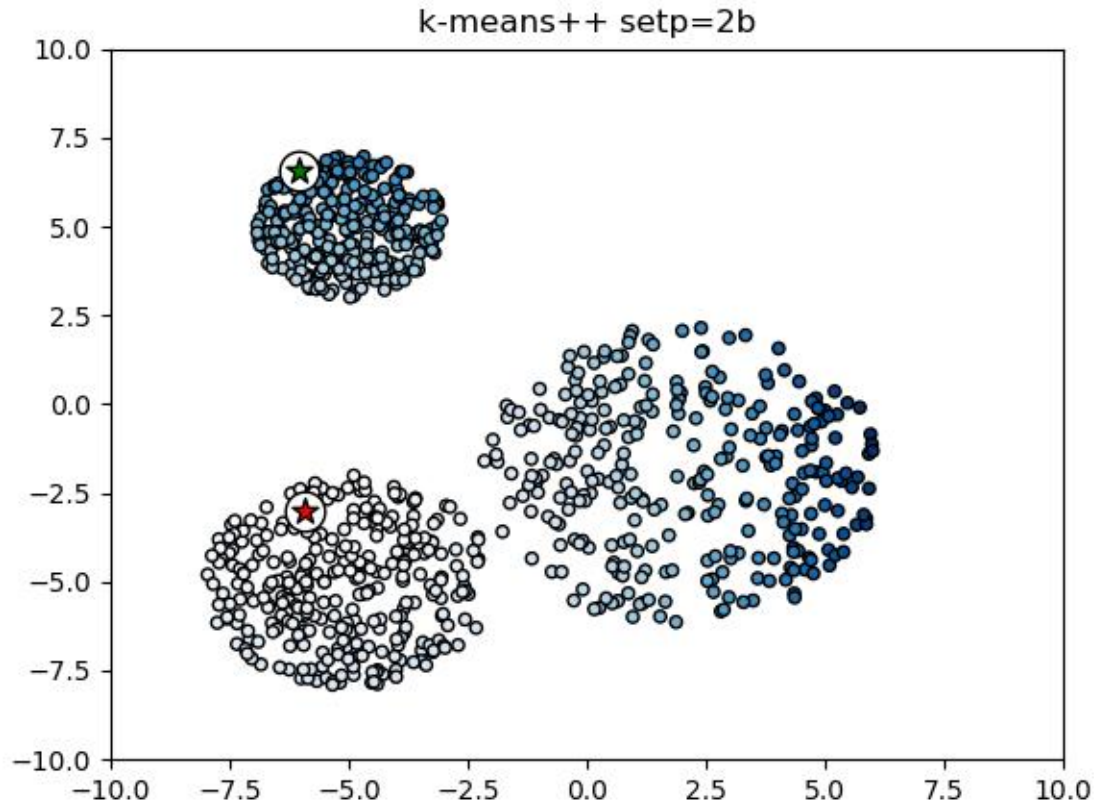
Jest to obecnie najczęściej stosowany algorytm inicjalizacji k-means (z reguły standardowa opcja).

K-means++



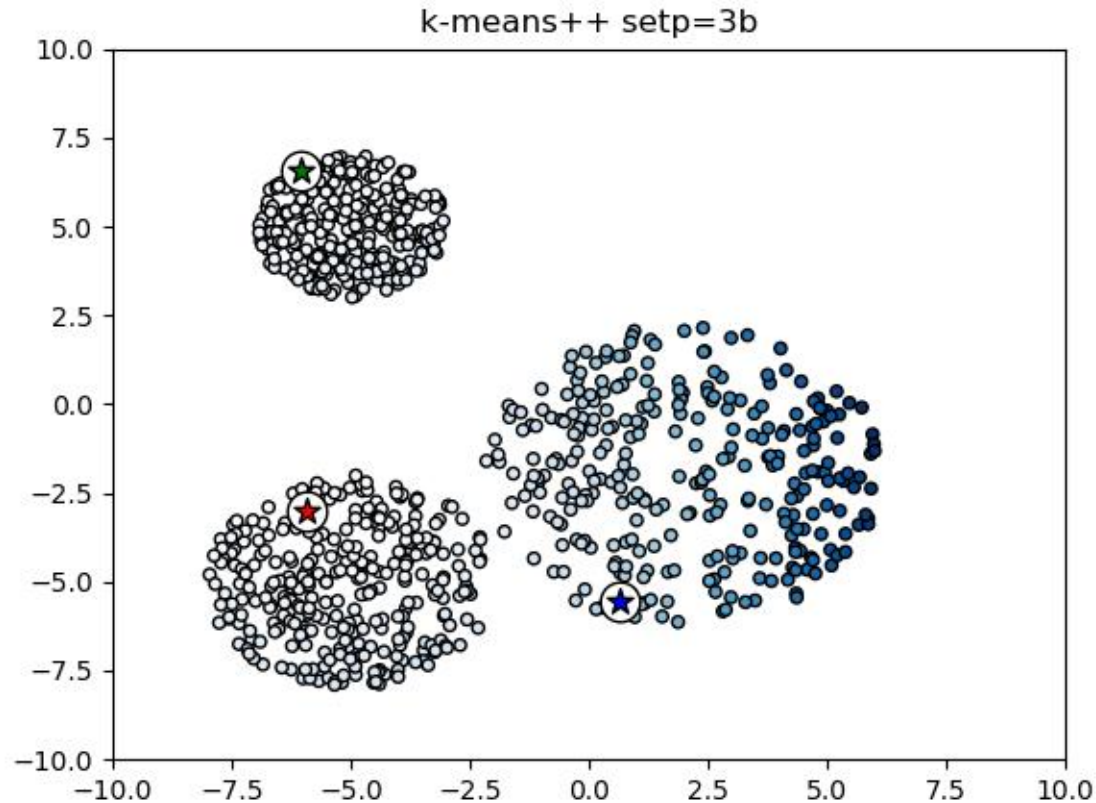
Losowany jest pierwszy punkt centralny...

K-means++



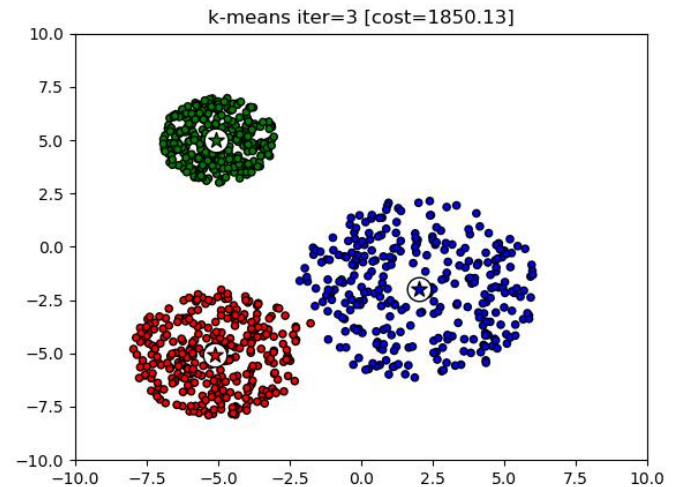
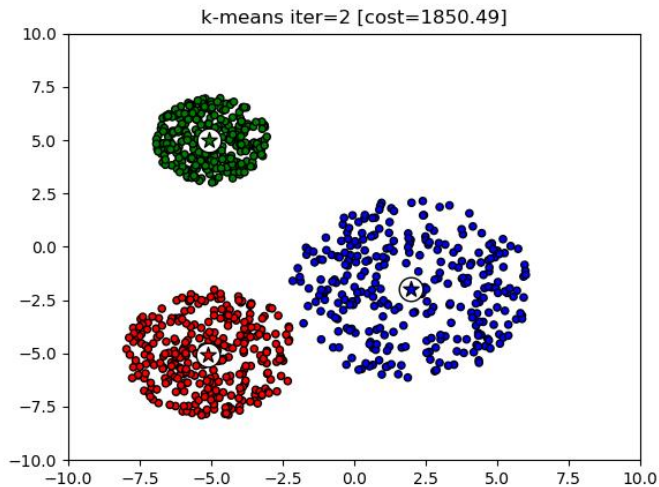
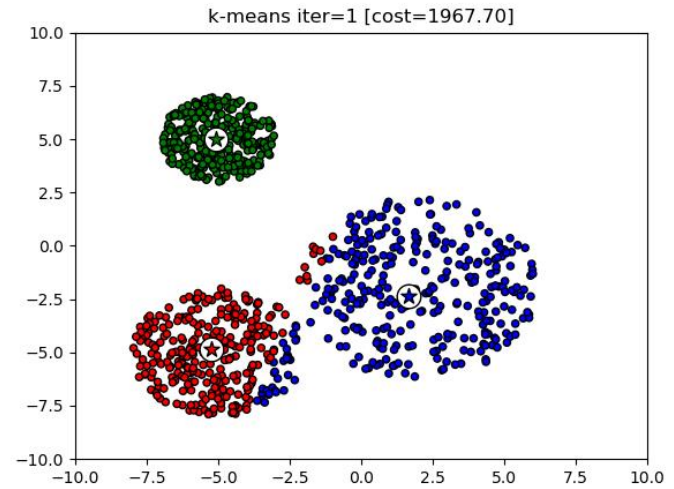
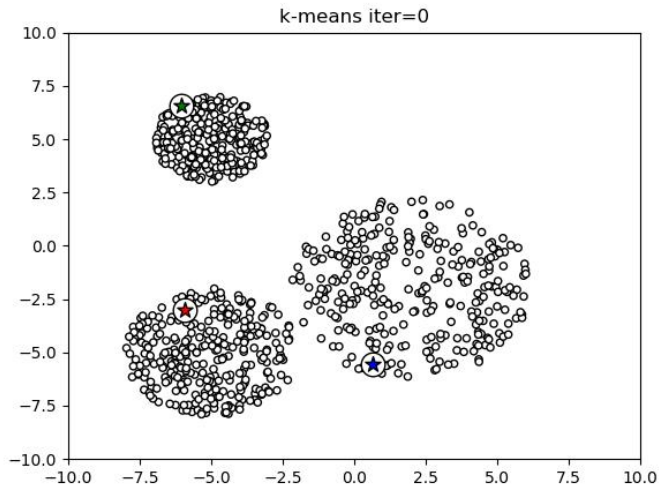
Wylosowany zostaje kolejny środek grupy. Prawdopodobieństwo rośnie wraz z kwadratem odległości od pierwszego punktu centralnego.

K-means++



Losowany jest trzeci punkt. Nie musi być to najdalej położona obserwacja (ale prawdopodobieństwo wylosowania dalej położonych punktów jest większe)

Dalszy przebieg algorytmu k-means



EM

Mieszananina rozkładów Gaussa (ang. Gaussian Mixture Model)

Rozkład Gaussa

- Rozkład Gaussa (normalny) dla $x \in \mathbb{R}$

$$p(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}$$

- μ – wartość oczekiwana (średnia wartość)
- σ^2 - wariancja

- Rozkład Gaussa dla $x \in \mathbb{R}^n$

$$p(x|\mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^n |\Sigma|}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)}$$

- $\mu \in \mathbb{R}^n$ - wartość oczekiwana dla x , $\mu = E(X)$
- Σ – macierz kowariancji $\Sigma = E[(X - \mu)(X - \mu)^T]$
- $|\Sigma|$ - wyznacznik macierzy, musi spełniać $|\Sigma| > 0$

Rozkład Gaussa – przypadek 2D

- $x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \in R^2$
- Środek rozkładu: $\mu = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix} \in R^2$
- Macierz kowariancji $\Sigma = \begin{bmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{21} & \sigma_2^2 \end{bmatrix} = \begin{bmatrix} \sigma_1^2 & r\sigma_1\sigma_2 \\ r\sigma_1\sigma_2 & \sigma_2^2 \end{bmatrix},$

$\sigma_{12} = \sigma_{21}$ - kowariancja X_1 i X_2 (r – to współczynnik Pearsona)

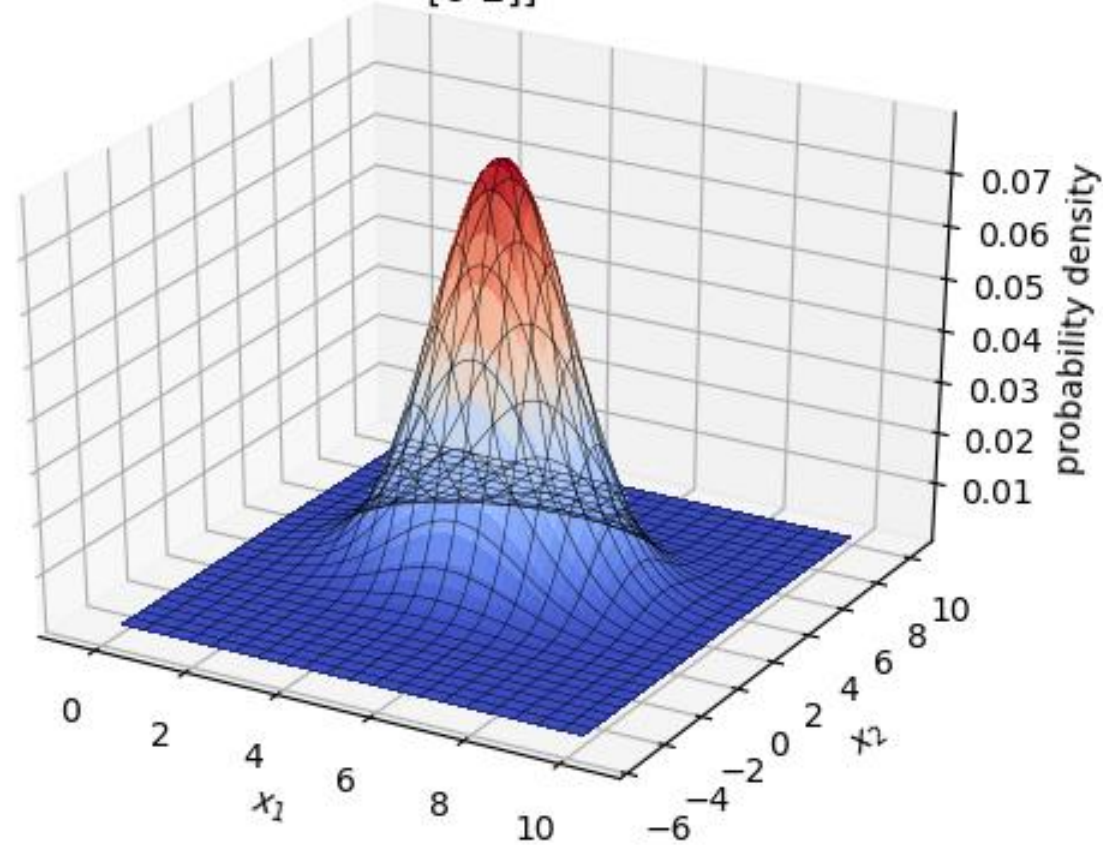
- Punkty, o jednakowej wartości prawdopodobieństwa tworzą elipsę opisaną równaniem:

$$\frac{(x_1 - \mu_1)^2}{\sigma_1^2} - 2r \frac{(x_1 - \mu_1)(x_2 - \mu_2)}{\sigma_1\sigma_2} + \frac{(x_2 - \mu_2)^2}{\sigma_2^2} = \text{const}$$

- Długości osi zależą od wartości σ_1 oraz σ_2
- Kąt obrotu zależy od wartości współczynnika Pearsona r (lub korelacji)

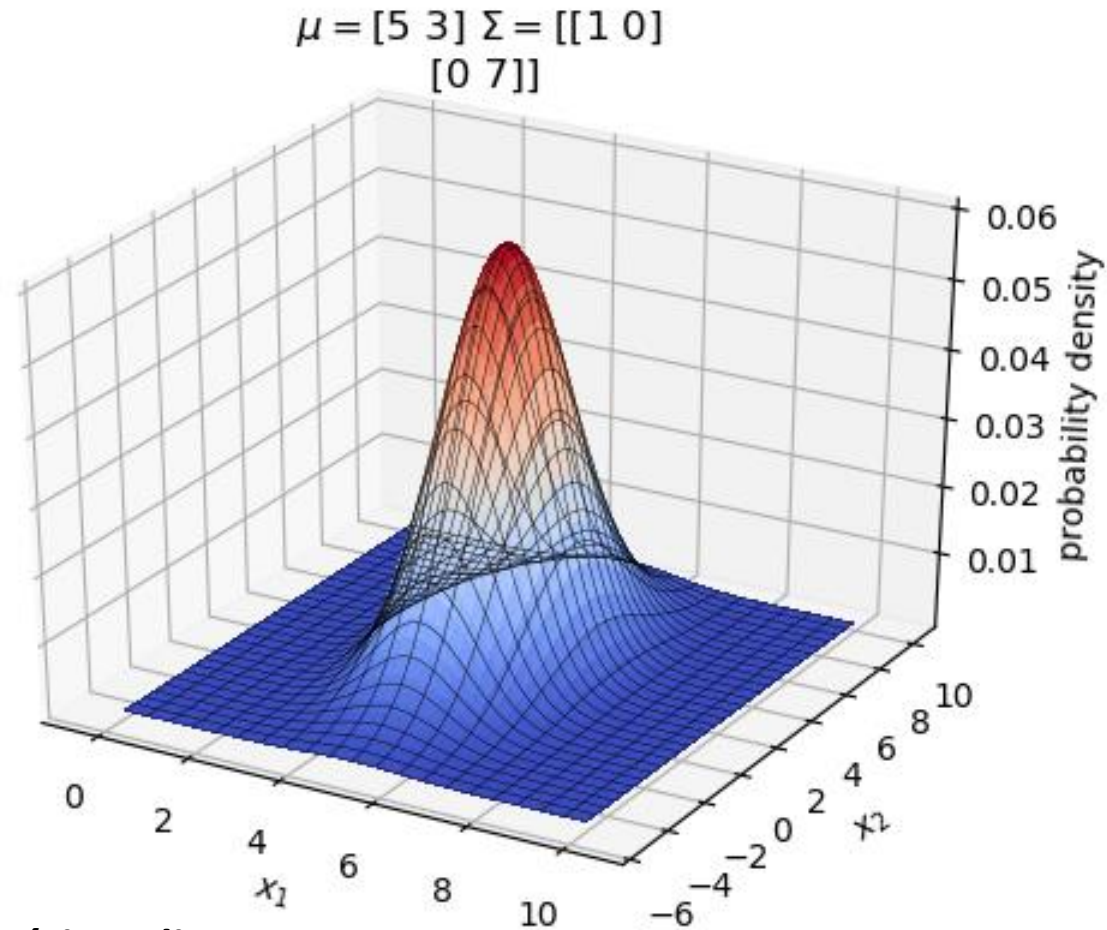
Przykład: $\sigma_1^2 = \sigma_2^2$

$$\mu = [5 \ 3] \quad \Sigma = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$$



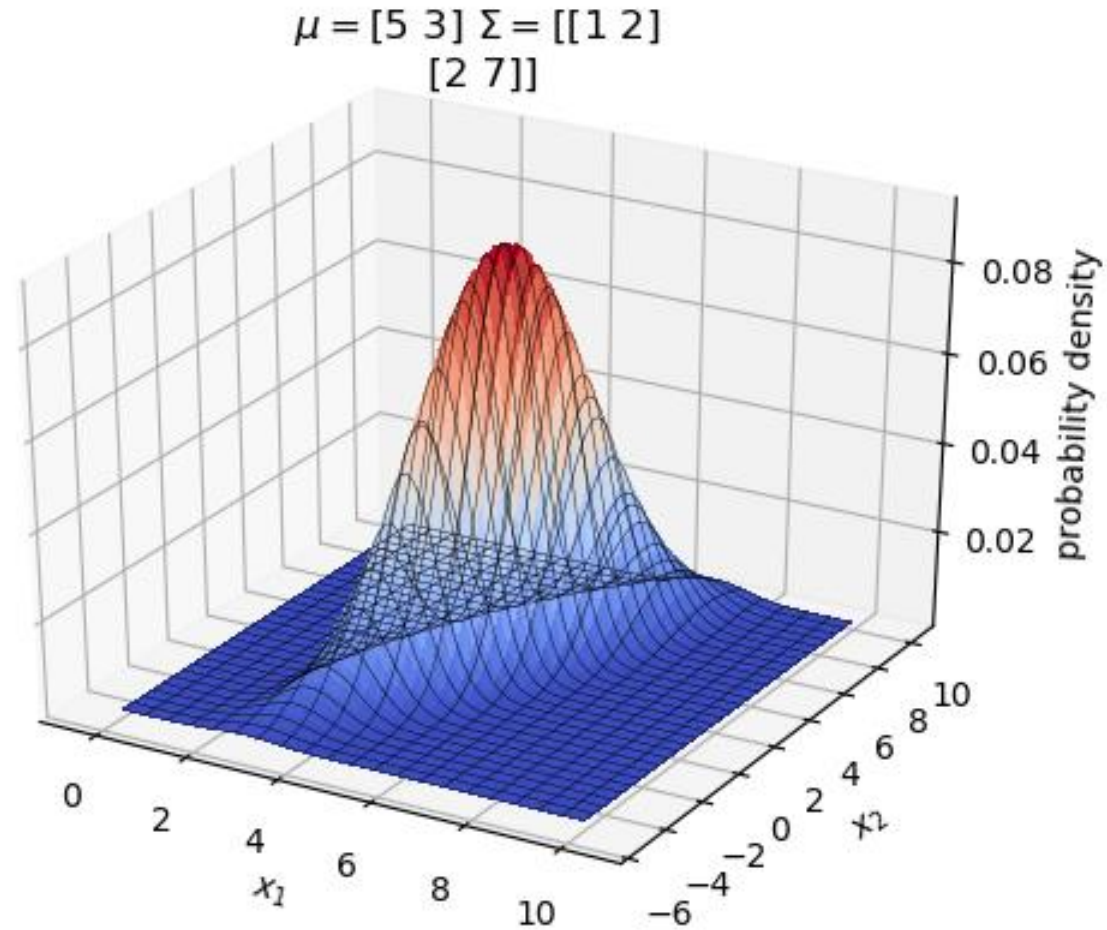
Przekrojami funkcji gęstości są okręgi

Przykład: $\sigma_1^2 \neq \sigma_2^2, r = 0$



Przekrojami funkcji gęstości są elipsy o osiach zgodnych z kierunkami x_1 i x_2

Przykład: $\sigma_1^2 \neq \sigma_2^2, r \neq 0$



Przekrojami funkcji gęstości są obrócone elipsy

Mieszanina k rozkładów Gaussa

- Załóżmy, że model składa się z k klas C_1, \dots, C_k , które mogą generować dane z prawdopodobieństwem opisanym rozkładem Gaussa $p(x|\mu_i, \Sigma_i)$. Modelem klasy jest więc para $\theta_i = (\mu_i, \Sigma_i)$.
- Proces losowania zaobserwowanych danych $x_j \in \mathbb{R}^n$ przebiega następująco:

- Wpierw losowana jest i -ta klasa z prawdopodobieństwem π_i
- Następnie losowany jest punkt x_j zgodnie z rozkładem

$$p(x|\mu_i, \Sigma_i) = p(x|\theta_i)$$

- Rozkład prawdopodobieństwa $f(x|\Theta)$ jest sumą rozkładów przeskalowanych przez współczynniki π_i (prawdopodobieństwa a priori klas) spełniające: $\forall i: \pi_i \geq 0$ oraz $\sum_{i=1}^k \pi_i = 1$

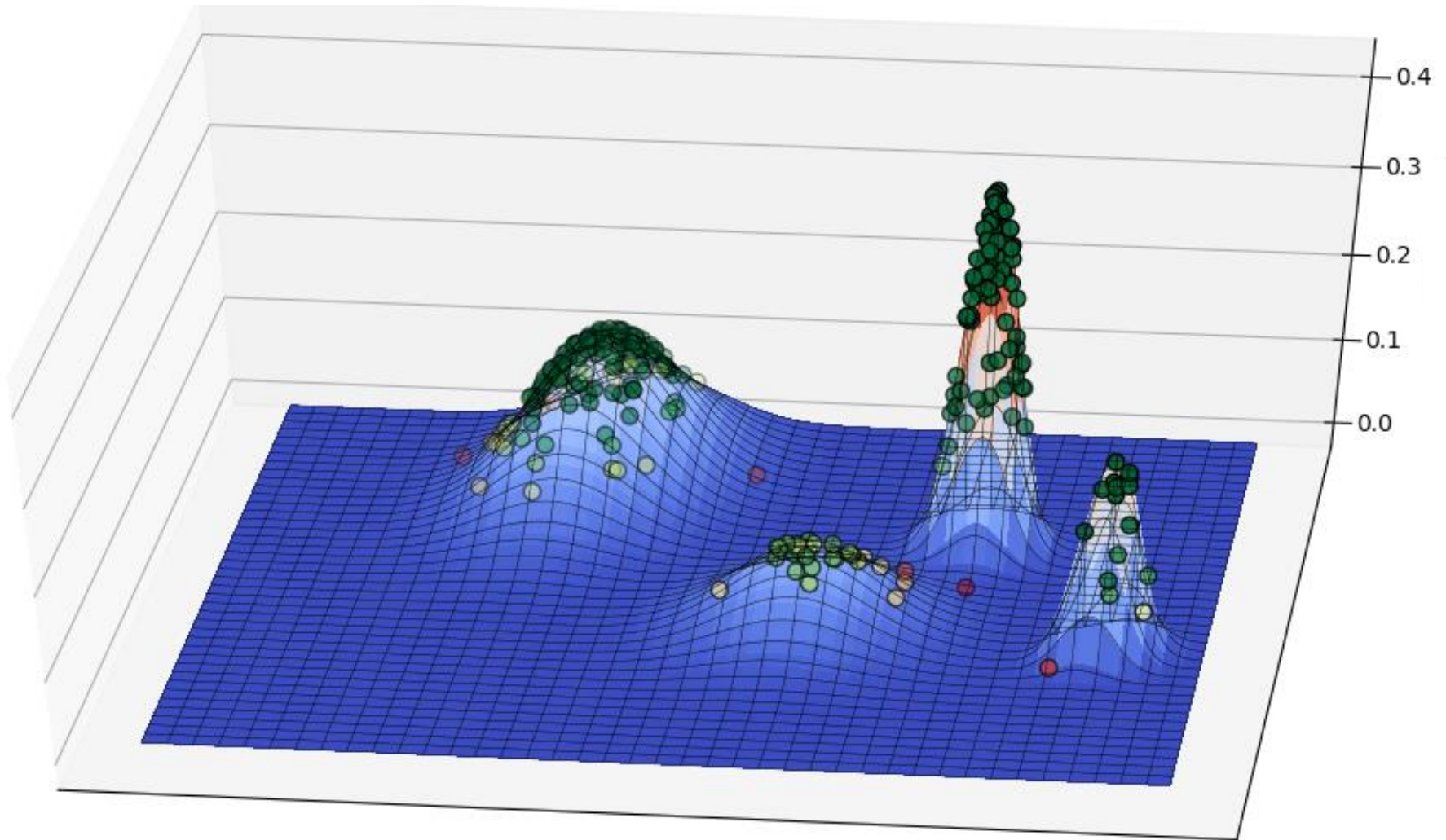
$$p(x|\Theta) = \sum_{i=1}^k \pi_i p(x|\theta_i)$$

- Pełne parametry modelu

$$\Theta = (\pi_1, \dots, \pi_k, \theta_1, \dots, \theta_k)$$

obejmują zarówno prawdopodobieństwa π_i , jak i parametry rozkładów Gaussa $\theta_i = (\mu_i, \Sigma_i)$.

Przykład 2D: mieszanina k rozkładów Gaussa



Mieszanina k rozkładów Gaussa (podział przestrzeni)

- Jeżeli znany jest pełny model $\Theta = (\pi_1, \dots, \pi_k, \theta_1, \dots, \theta_k)$, wówczas prawdopodobieństwo przynależności dowolnego punktu $x \in R^n$ do grupy C_i można obliczyć ze wzoru Bayesa:

$$P(C_i|x) = \frac{\pi_i p(x|\theta_i)}{\sum_{j=1}^k \pi_j p(x|\theta_j)}$$

- π_i to prawdopodobieństwo a priori
 - θ_j to model dla C_i
 - mianownik $\sum_{j=1}^k \pi_j p(x|\theta_j)$ to prawdopodobieństwo zaobserwowania x dla modelu Θ
- Na podstawie tego prawdopodobieństwa -- lub wyłącznie korzystając z licznika $\pi_i p(x|\theta_i)$ - można wyznaczyć naturalny podział przestrzeni na grupy

$$C_i = \left\{ x: \pi_i p(x|\theta_i) > \max_{j \neq i} \pi_j p(x|\theta_j) \right\}$$

Mieszanka k rozkładów Gaussa

- Zakładając, że obserwacje ze zbioru danych $D = \{x_j\}$, $j = 1, \dots, m$ są niezależne, prawdopodobieństwo jego zaobserwowania, dla danego zestawu parametrów Θ wynosi:

$$\begin{aligned} P(D|\Theta) &= \prod_{j=1}^m p(x_j|\Theta) = \prod_{j=1}^m \left(\sum_{i=1}^k \pi_i \cdot p(x_j|\theta_i) \right) \\ &= \prod_{j=1}^m \left(\sum_{i=1}^k \pi_i \cdot \frac{1}{\sqrt{(2\pi)^n |\Sigma|}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)} \right) \end{aligned}$$

- Parametry modelu Θ są jednak nieznane, natomiast znane są obserwacje D . Prawdopodobieństwo (ang. likelihood), że znane obserwacje D powstały jako wynik losowania dla badanego modelu Θ ma taką samą wartość

$$L(\Theta|D) = P(D|\Theta)$$

Mieszanina k rozkładów Gaussa

- Po zastosowaniu logarytmu (zamiana iloczynu małych wartości prawdopodobieństwa na sumę ich logarytmów) otrzymujemy log-likelihood

$$LL(\Theta|D) = \ln \prod_{j=1}^m \left(\sum_{i=1}^k \pi_i \cdot \frac{1}{\sqrt{(2\pi)^n |\Sigma|}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)} \right) =$$
$$\sum_{j=1}^m \ln \left(\sum_{i=1}^k \pi_i \cdot \frac{1}{\sqrt{(2\pi)^n |\Sigma|}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)} \right)$$

oraz przepis na optymalną wartość parametrów

$$\Theta^* = \arg \max_{\Theta} LL(\Theta|D)$$

- Niestety logarytm naturalny sumy jest trudny do optymalizacji...

Algorytm EM 1

- W algorytmie EM zakłada się, że model zawiera m (m to liczba danych) ukrytych zmiennych $Z = \{Z_j\}$.
- Zmienne te przyjmują wartości ze zbioru $\{1, \dots, k\}$. Jeżeli $Z_j = i$, to oznacza, że obserwacja x_j została wygenerowana przez i -ty komponent mieszaniny zgodnie z rozkładem

$$\pi_i \cdot p(x_j | \theta_i) = \pi_i \cdot \frac{1}{\sqrt{(2\pi)^n |\Sigma|}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)}$$

(czyli x_j należy też do i -tej grupy).

Algorytm EM 2

- Wówczas zamiast maksymalizować

$$L(\Theta|D) = P(D|\Theta)$$

maksymalizuje się

$$L(\Theta|D, Z) = P(D, Z|\Theta)$$

- Przekształcając:

$$P(D, Z|\Theta) = \frac{P(D, Z, \Theta)}{P(\Theta)} = \frac{P(Z|D, \Theta) \cdot P(D, \Theta)}{P(\Theta)} = P(Z|D, \Theta) \cdot P(D|\Theta)$$

Pozwala to skonstruować iteracyjną procedurę, złożoną z dwóch kroków:

Expectation. Wpierw wyznacza się prawdopodobieństwa $P(Z|D, \Theta^{(t-1)})$ poszczególnych wartości zmiennych Z na podstawie obserwacji ze zbioru D i parametrów z modelu poprzedniej iteracji $\Theta^{(t-1)}$

Maximization. Następnie poszukuje się nowego modelu $\Theta^{(t)}$ maksymalizującego prawdopodobieństwo $P(D|\Theta^{(t)})$ przy ustalonym w fazie E rozkładzie zmiennych Z

Expectation

- Oznaczmy przez w_{ji} prawdopodobieństwo, że $Z_j = i$ (czyli, że obserwacja x_j została wygenerowana przez i -ty komponent)

$$w_{ji} = \frac{\pi_i \cdot p(x_j | \theta_i)}{\sum_{l=1}^k \pi_l \cdot p(x_j | \theta_l)}$$

- Wartość w_{ji} może być traktowana jako waga określająca przynależność obserwacji x_j do i -tej grupy. Opisuje ona niepewność, który komponent wygenerował obserwację x_j dla danych parametrów Θ

Maximization

W_{ji}

0.5	0.2	0.3
0.2	0.4	0.4
...
0.2	0.2	0.6
0.1	0.7	0.2



m wierszy

$$\sum_{i=1}^k w_{ji} = 1$$

k kolumn

Parametry nowego modelu $\Theta^{(t)}$:

- Wagi (prawdopodobieństwa) komponentów:

$$\pi_i = \frac{\sum_{j=1}^m w_{ji}}{m}$$

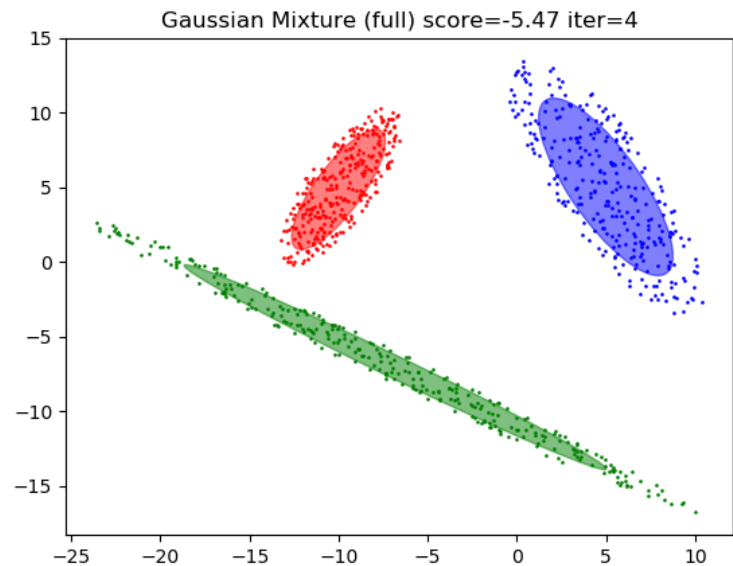
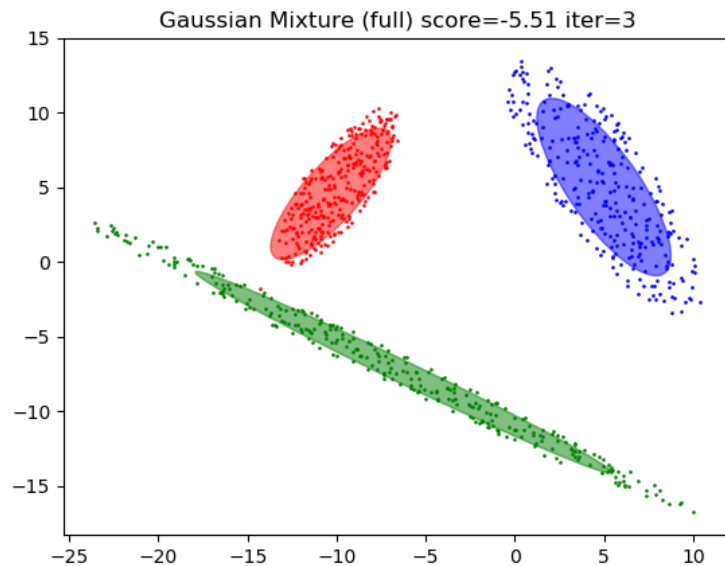
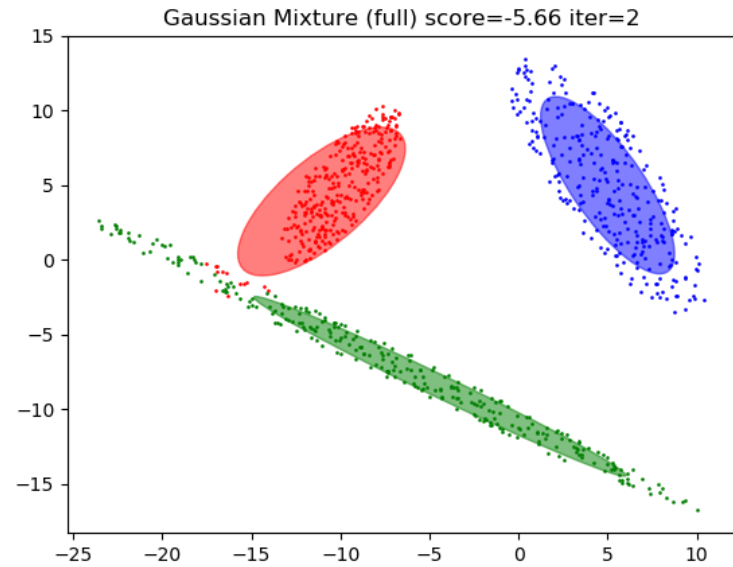
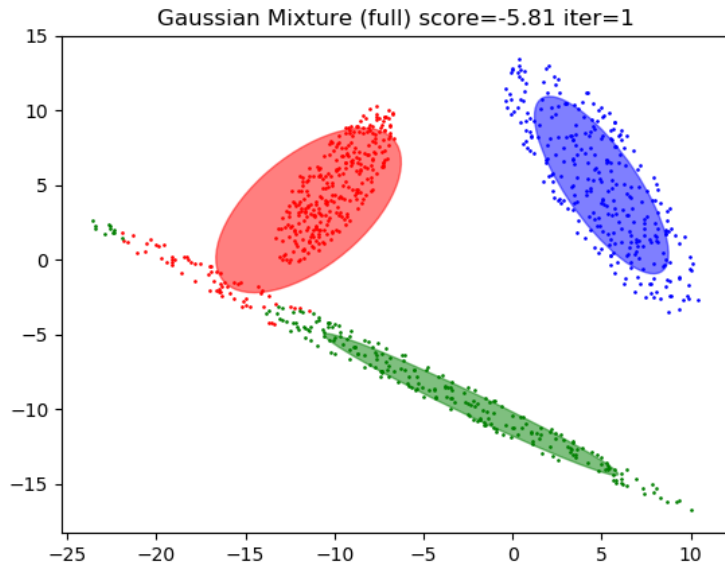
- Środki komponentów

$$\mu_i = \frac{\sum_{j=1}^m w_{ji} \cdot x_j}{\sum_{j=1}^m w_{ji}}$$

- Macierze kowariancji

$$\Sigma_i = \frac{\sum_{j=1}^m w_{ji} \cdot (x_j - \mu_i)(x_j - \mu_i)^T}{\sum_{j=1}^m w_{ji}}$$

Przykładowy przebieg algorytmu



Inicjalizacja i kryterium stopu

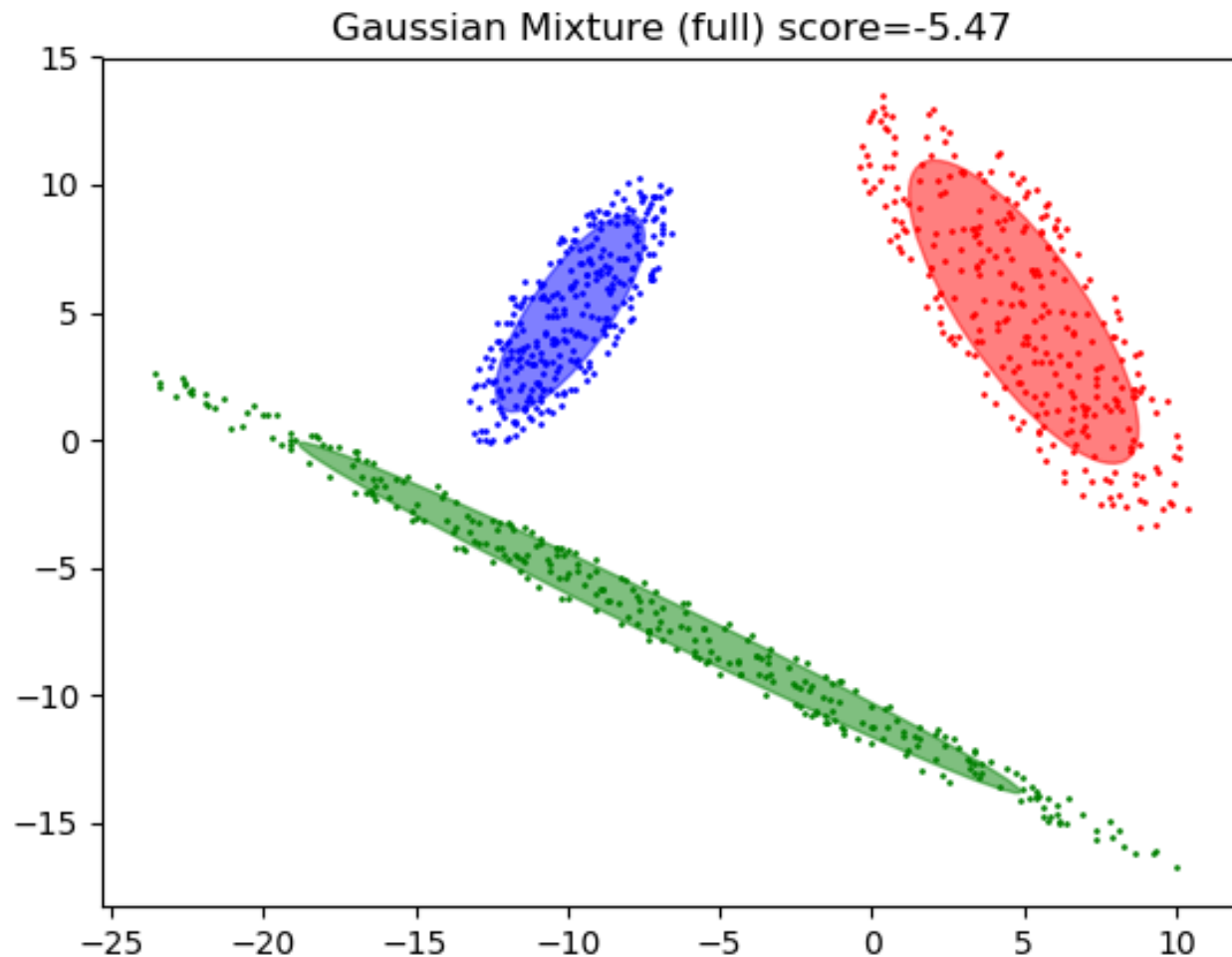
- Inicjalizacja możliwa jest na dwa sposoby
 - Wstępny wybór parametrów Θ i przeprowadzenie kroku E
 - Wstępny wybór wag w_{ji} i przeprowadzenie kroku M
 - W obu przypadkach dane mogą być wybrane losowo (np. przez wylosowanie punktów środkowych i wartości wariancji) lub też przez zastosowanie jakiegoś algorytmu heurystycznego. Typowe rozwiązanie – użycie k-means.
 - Kryteria stopu
 - Brak poprawy funkcji celu
- $$J(\Theta) = \sum_{j=1}^m \ln \left(\sum_{i=1}^k \pi_i \cdot \frac{1}{\sqrt{(2\pi)^n |\Sigma|}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)} \right)$$
- Maksymalna liczba iteracji

Postać modelu

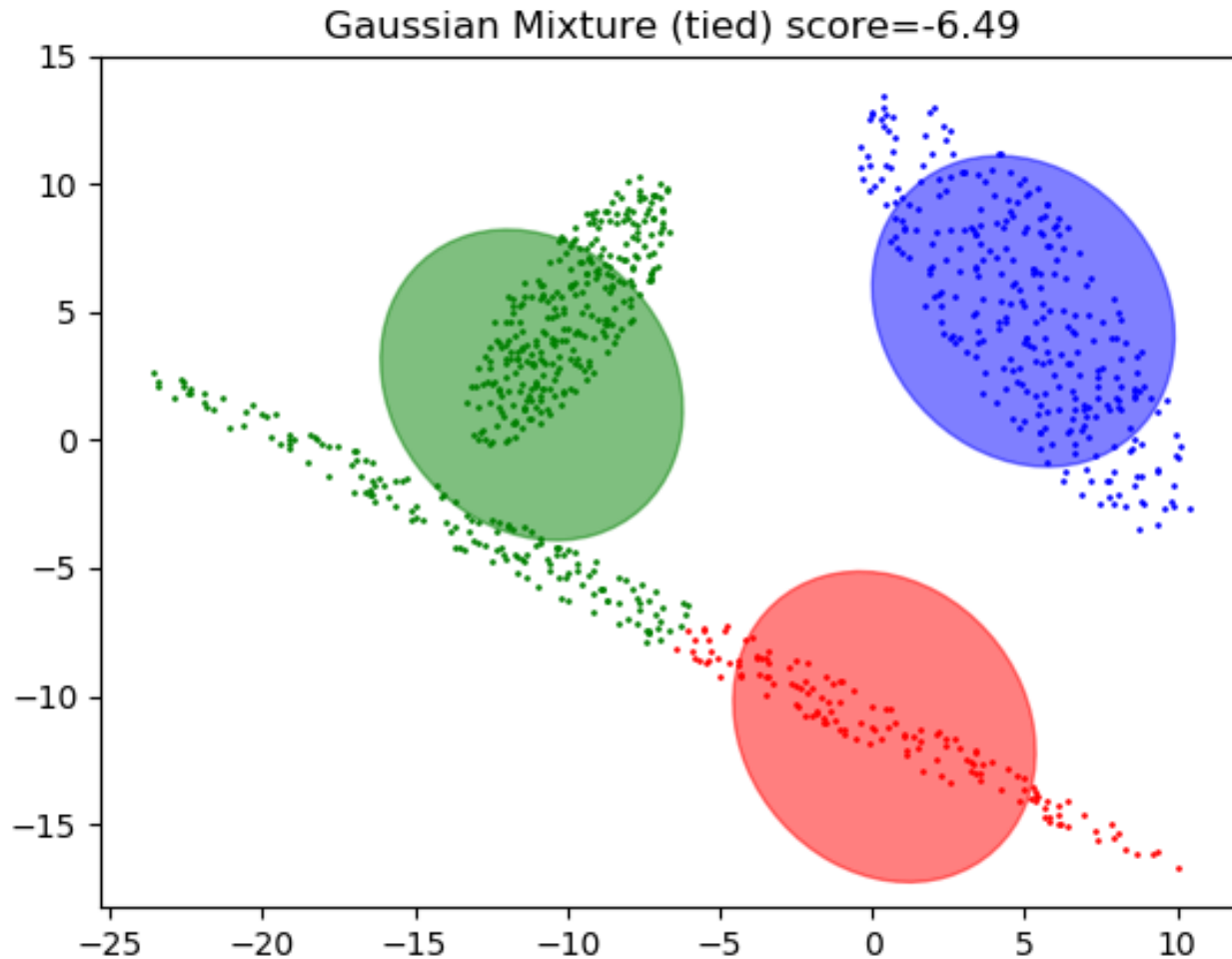
W zależności od implementacji mogą być budowane uproszczone modele.

- Python – wybór pomiędzy
 - full – pełna macierz kowariancji
 - tied – jedna macierz kowariancji dla wszystkich grup
 - diag – diagonalna macierz kowariancji (czyli zawierająca wyłącznie wariacje dla poszczególnych atrybutów)
 - spherical – jednakowa wariacja dla wszystkich wymiarów i grup
- Weka – implementuje wyłącznie model diag

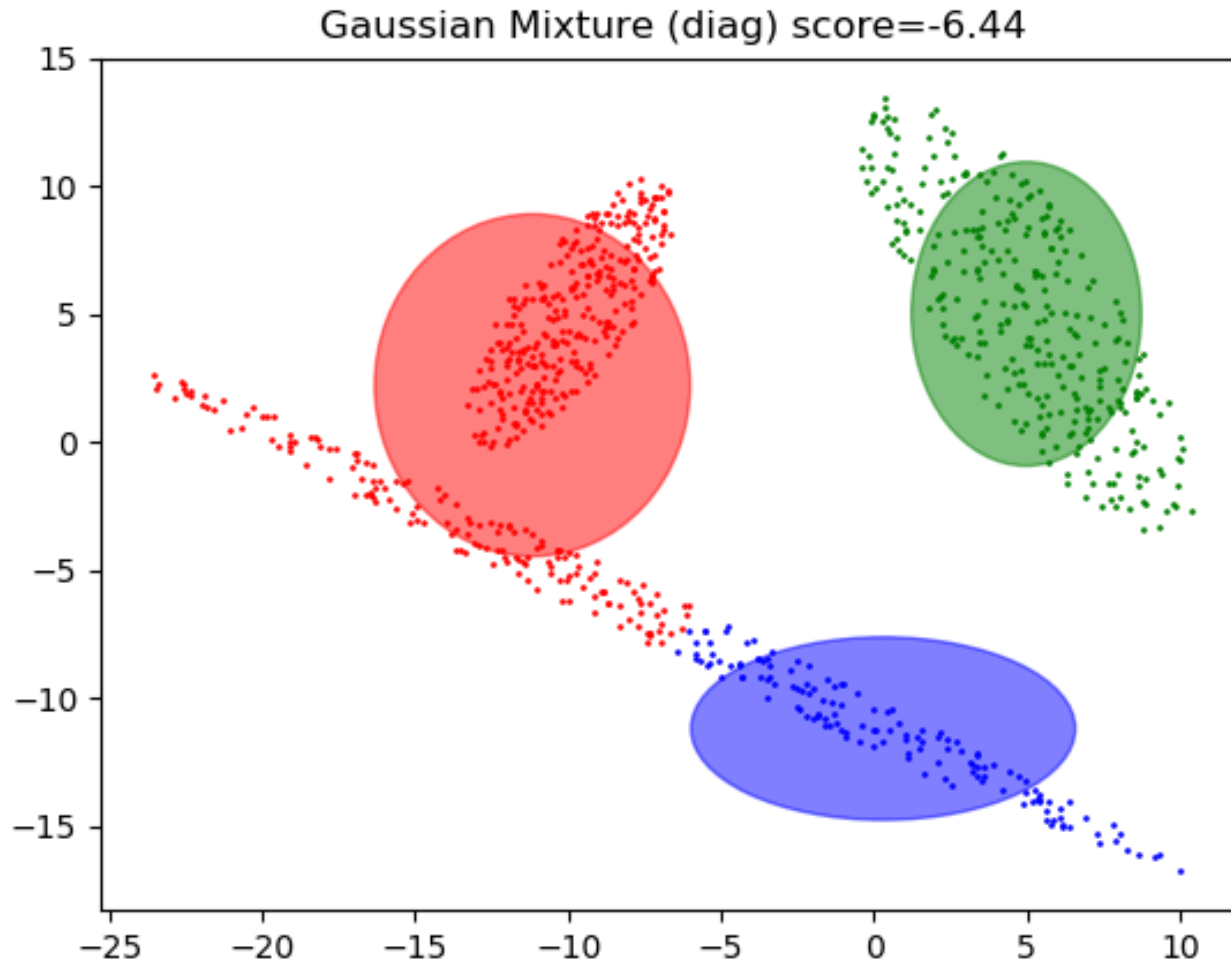
Model full – pełna macierz kowariancji



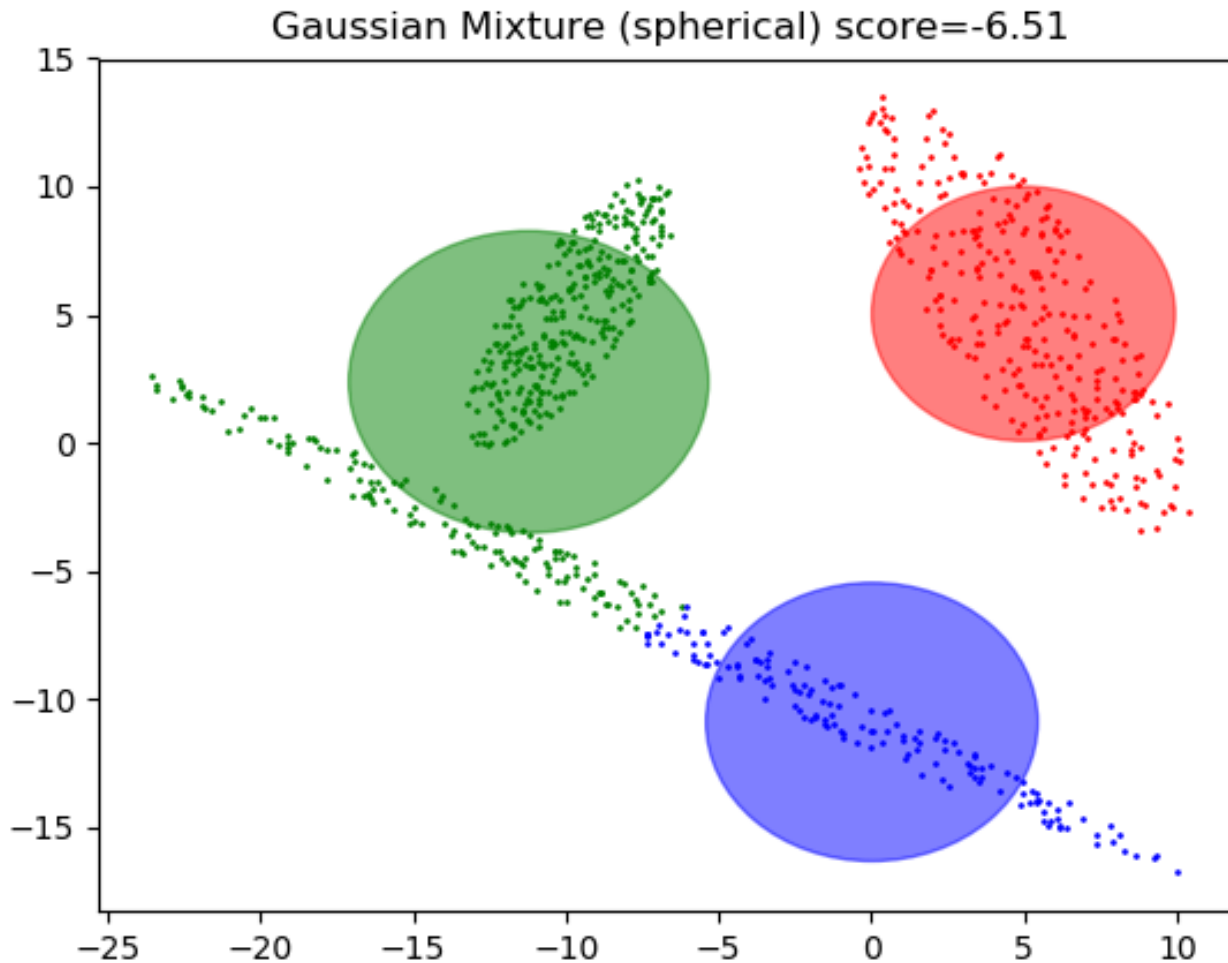
Model tied – jedna macierz kowariancji dla wszystkich grup



Model diag - wyłącznie wariancje dla poszczególnych atrybutów



Model spherical - jednakowa wariancja dla wszystkich wymiarów i grup



K-means vs. EM

- Model spherical mieszanki Gaussa odpowiada założeniom k-means.
- Dodatkowo , w k-means wariancje muszą być małe (tak, aby kule zmieściły się w komórkach Woronoja)
- W rzeczywistości algorytm k-means może być potraktowany jako uproszczona procedura EM:
 - $\Theta = \{\bar{x}(1), \dots, \bar{x}(k)\}$
 - **Expectation** – „twardy” przydział obserwacji do grup przy ustalonych położeniach centroidów z poprzedniej iteracji $\Theta^{(t-1)}$
 - **Maximization** – uaktualnienie parametrów modelu $\Theta^{(t)}$ przy założonym przydziale obserwacji do grup

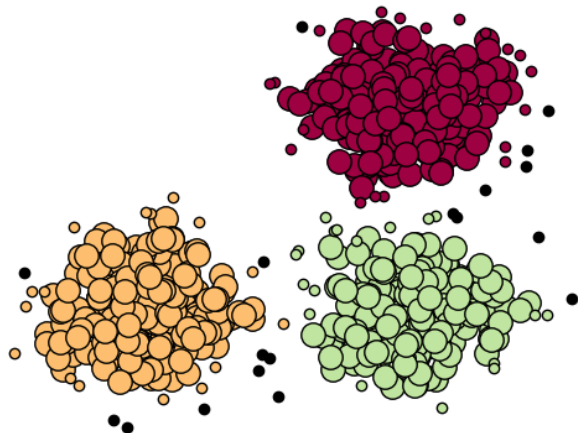
DBSCAN

DBSCAN - założenia

Algorytm DBSCAN (Density-Based Spatial Clustering of Applications with Noise) wykorzystuje informacje o gęstości punktów w przestrzeni do przeprowadzenia grupowania.

Założenia:

- Skupienia mogą zostać zidentyfikowane jako zbiory gęsto (blisko) położonych punktów
- W przestrzeni pomiędzy skupieniami mogą znaleźć się rzadko rozmieszczone punkty, które są uznane za szum.



[http://scikit-learn.org/stable/auto_examples/cluster/plot_dbscan.html#sphx-glr-auto-examples-cluster-plot-dbscan-py]

DBSCAN - definicje

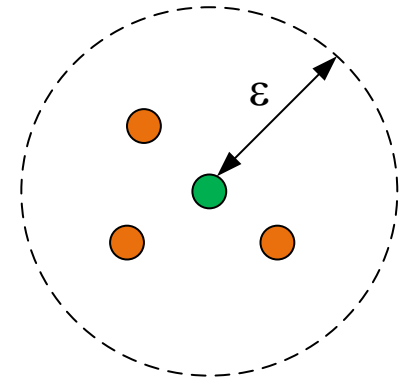
Parametry algorytmu:

- $d(x_1, x_2)$ – odległość pomiędzy obiektami x_1 i x_2
- ε – promień otoczenia obiektu
- $minPts$ – minimalna liczba punktów w otoczeniu

Definicje

- Otoczenie $N(x) = \{z \in D: d(x, z) < \varepsilon\}$
- Punkt jest punktem **wewnętrznym** (ang. core point), jeżeli jego otoczenie zawiera co najmniej $minPts$ punktów
- Punkt z jest **bezpośrednio osiągalny** (ang. directly reachable) z x , jeżeli x jest punktem wewnętrznym i z należy do sąsiedztwa x

$$core(x) \Leftrightarrow |N(x)| \geq minPts$$
$$x \xrightarrow{dr} z \Leftrightarrow core(x) \wedge z \in N(x)$$



DBSCAN -definicje

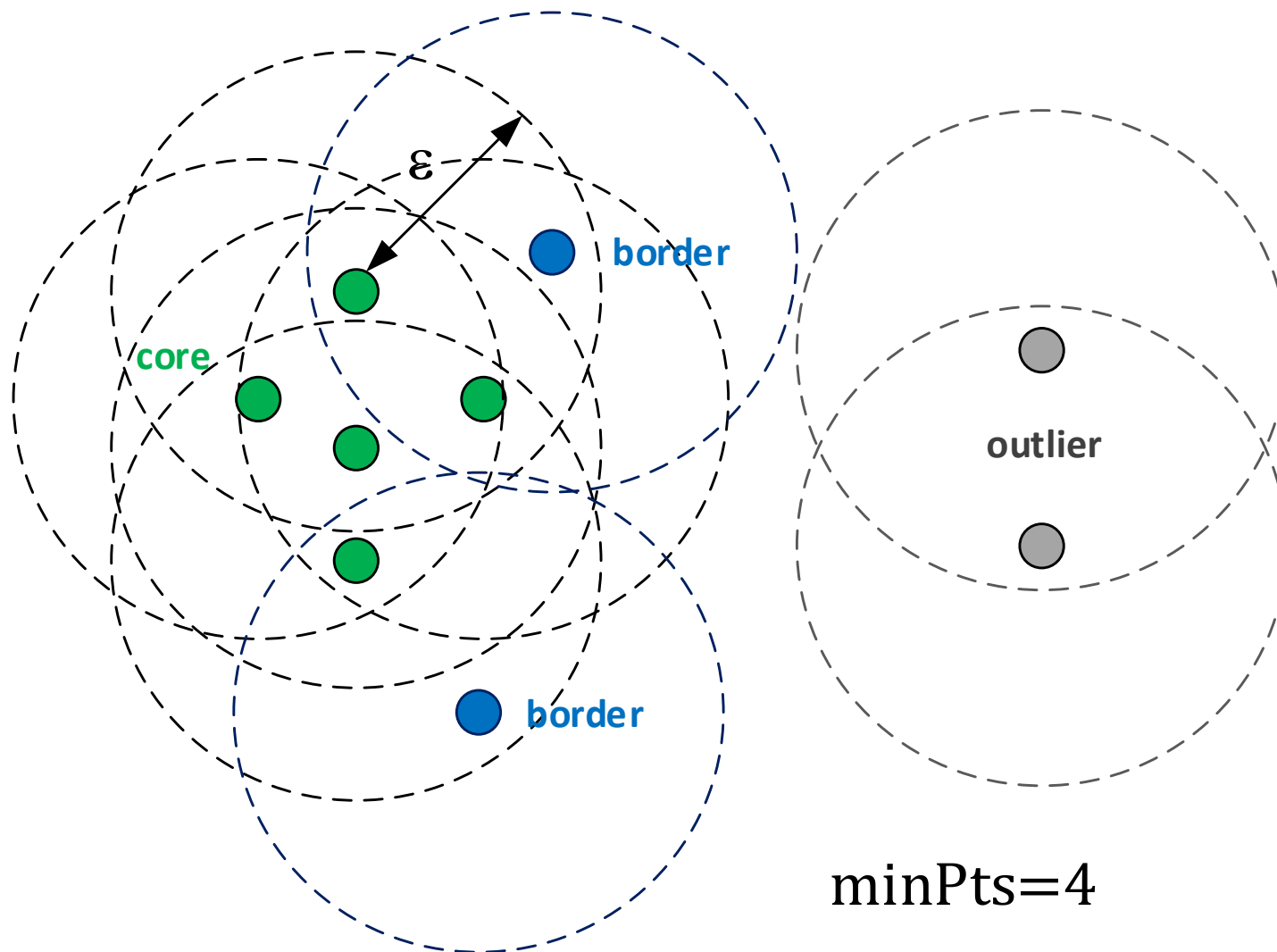
- Punkt z jest **osiągalny** z x , jeżeli istnieje ścieżka bezpośrednio osiągalnych punktów prowadząca od x do z

$$x \xrightarrow{dr} v_1 \xrightarrow{dr} \dots \xrightarrow{dr} v_k \xrightarrow{dr} z$$

czyli wszystkie punkty x, v_1, \dots, v_n muszą być punktami wewnętrznymi (mieć co najmniej $minPts$ sąsiadów) i być położone w odległości nie większej niż ϵ od siebie

- Jeżeli z nie jest punktem wewnętrznym (w jego otoczeniu jest mniej niż $minPts$ punktów), to nazywany jest punktem **brzegowym**.
- Punkty, które nie są osiągalne z żadnego punktu wewnętrznego są punktami poza grupami (outliers)

DBSCAN definicje - ilustracja



DBSCAN - algorytm

ALGORITHM 1: Pseudocode of Original Sequential DBSCAN Algorithm

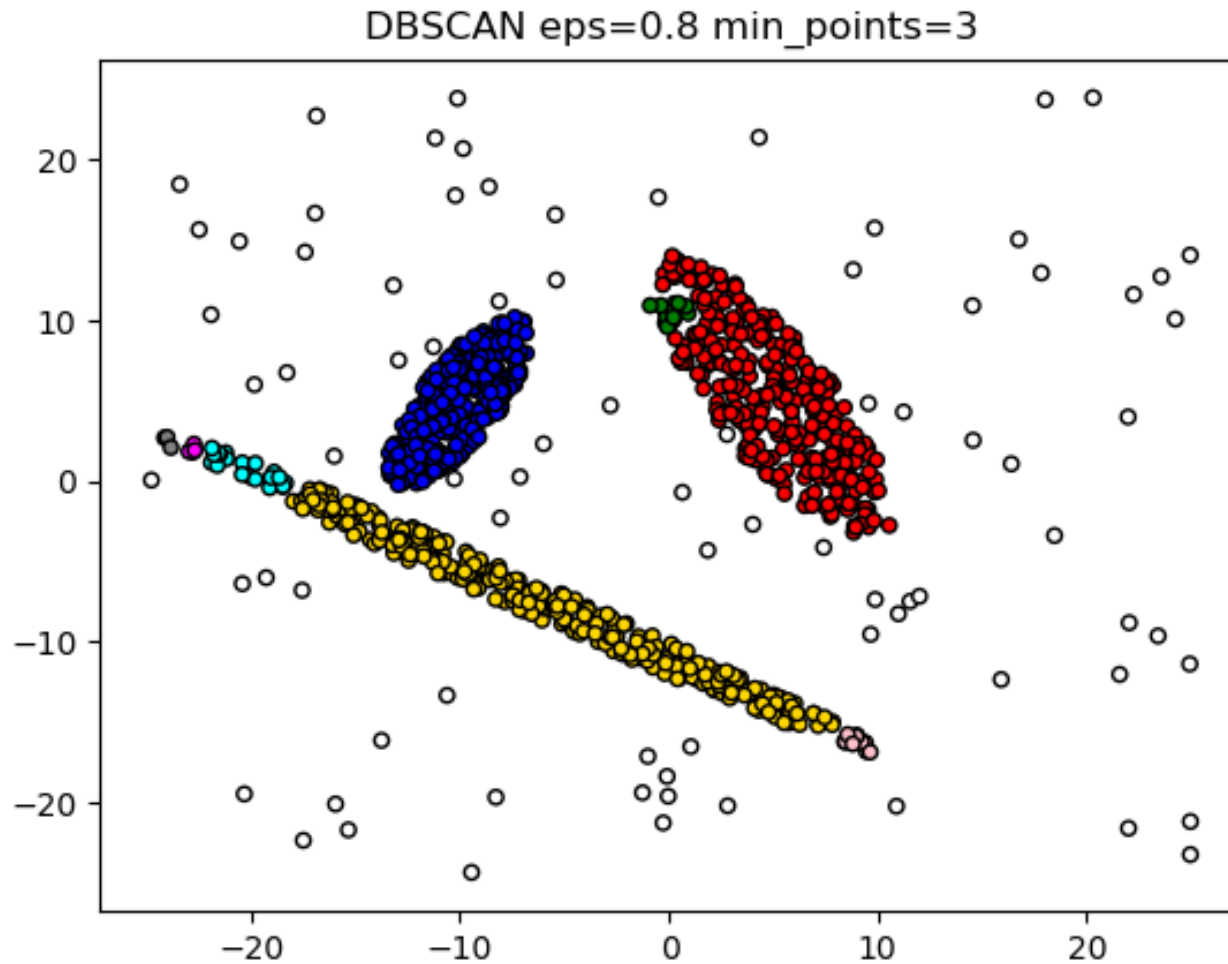
```
Input: DB: Database
Input:  $\epsilon$ : Radius
Input: minPts: Density threshold
Input: dist: Distance function
Data: label: Point labels, initially undefined
1 foreach point p in database DB do // Iterate over every point
2   if label(p) ≠ undefined then continue // Skip processed points
3   Neighbors N  $\leftarrow$  RANGEQUERY(DB, dist, p,  $\epsilon$ ) // Find initial neighbors
4   if  $|N| < \textit{minPts}$  then // Non-core points are noise
5     label(p)  $\leftarrow$  Noise
6     continue
7   c  $\leftarrow$  next cluster label // Start a new cluster
8   label(p)  $\leftarrow$  c
9   Seed set S  $\leftarrow$   $N \setminus \{p\}$  // Expand neighborhood
10  foreach q in S do
11    if label(q) = Noise then label(q)  $\leftarrow$  c
12    if label(q) ≠ undefined then continue
13    Neighbors N  $\leftarrow$  RANGEQUERY(DB, dist, q,  $\epsilon$ )
14    label(q)  $\leftarrow$  c
15    if  $|N| < \textit{minPts}$  then continue // Core-point check
16    S  $\leftarrow$   $S \cup N$ 
```

Niejednoznaczność dla punktów brzegowych

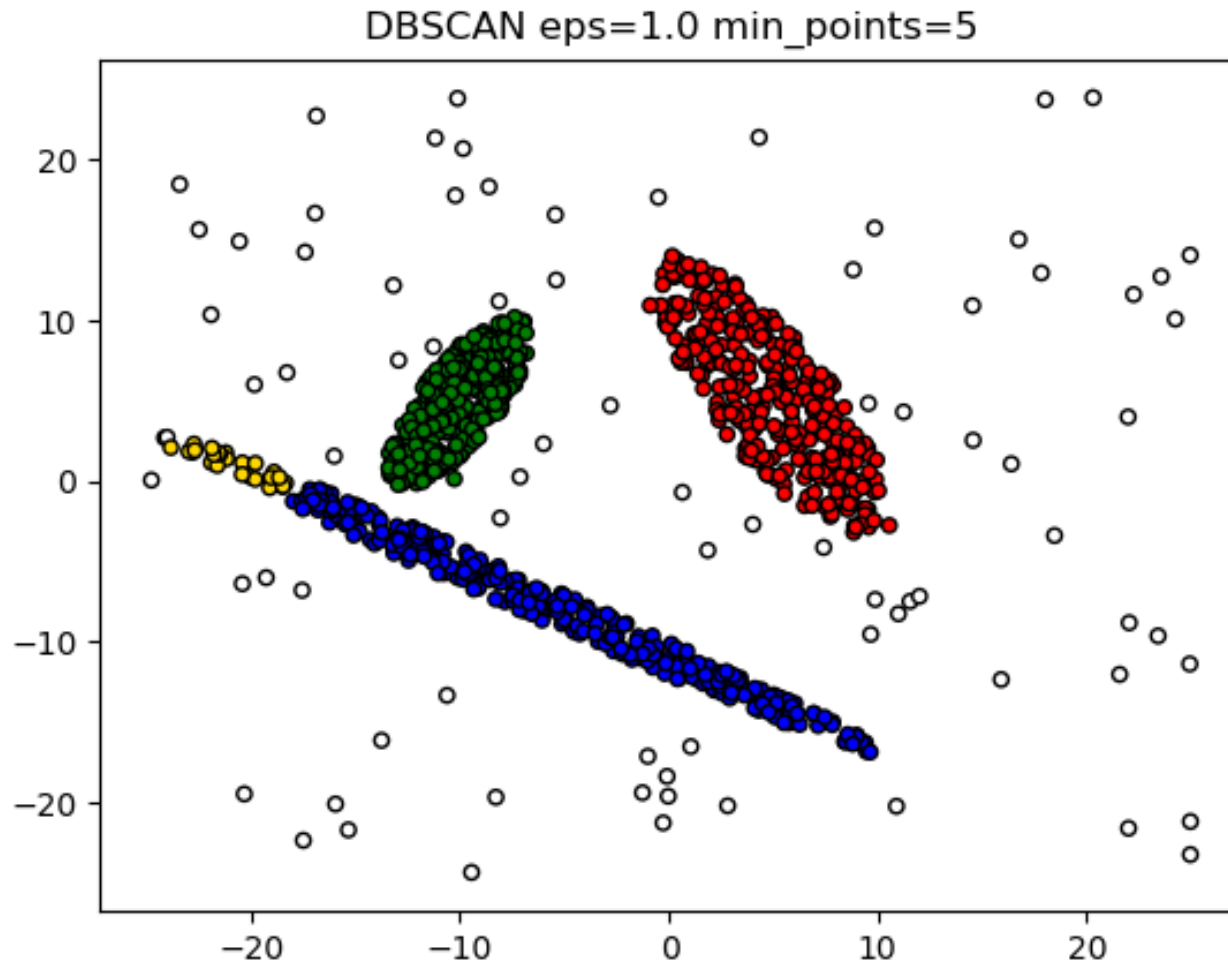
Schubert, Erich; Sander, Jörg; Ester, Martin; [Kriegel, Hans Peter](#); Xu, Xiaowei (July 2017).

"DBSCAN Revisited, Revisited: Why and How You Should (Still) Use DBSCAN". *ACM Trans. Database Syst.* 42 (3): 19:1–19:21.

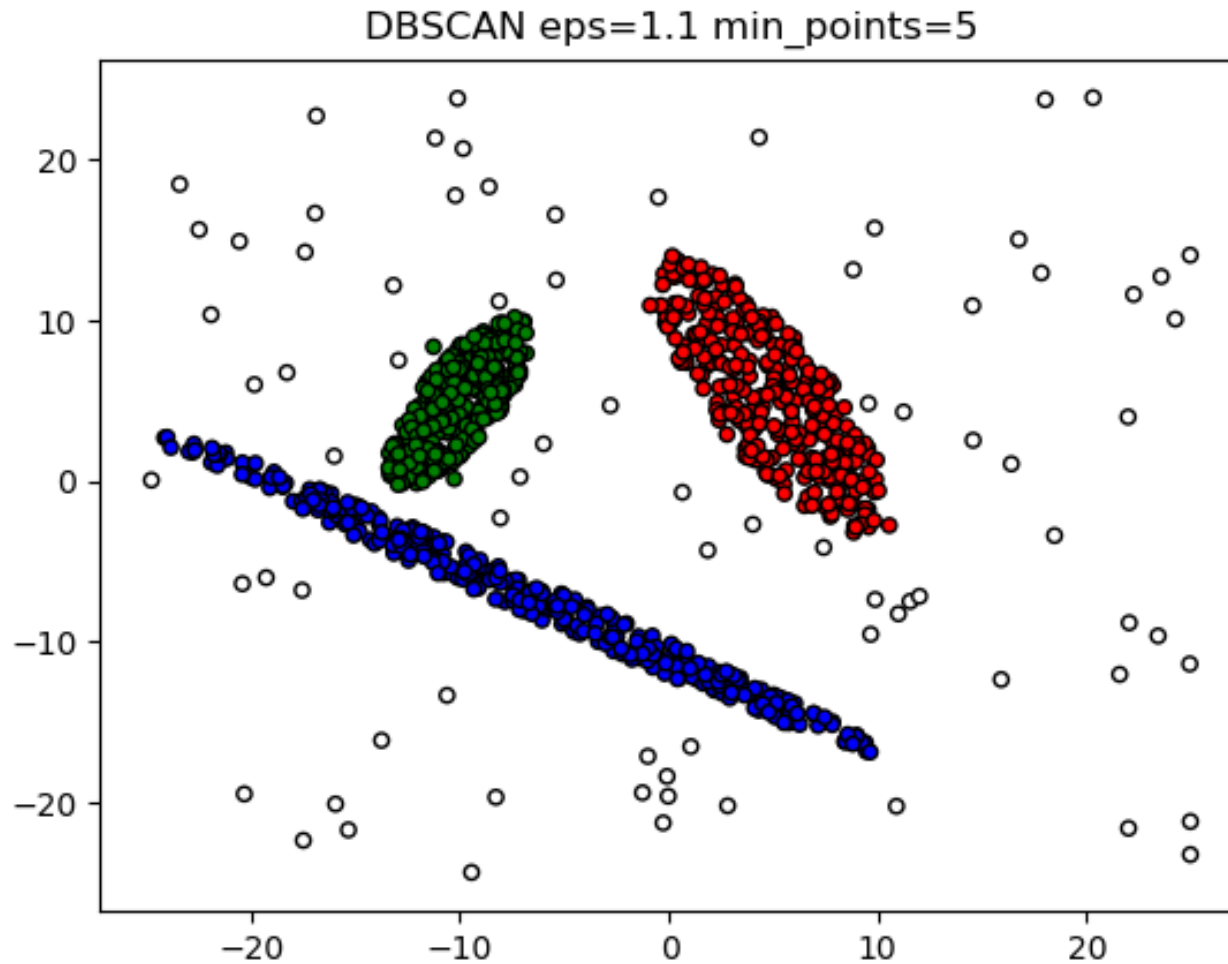
DBSCAN – przykład 1



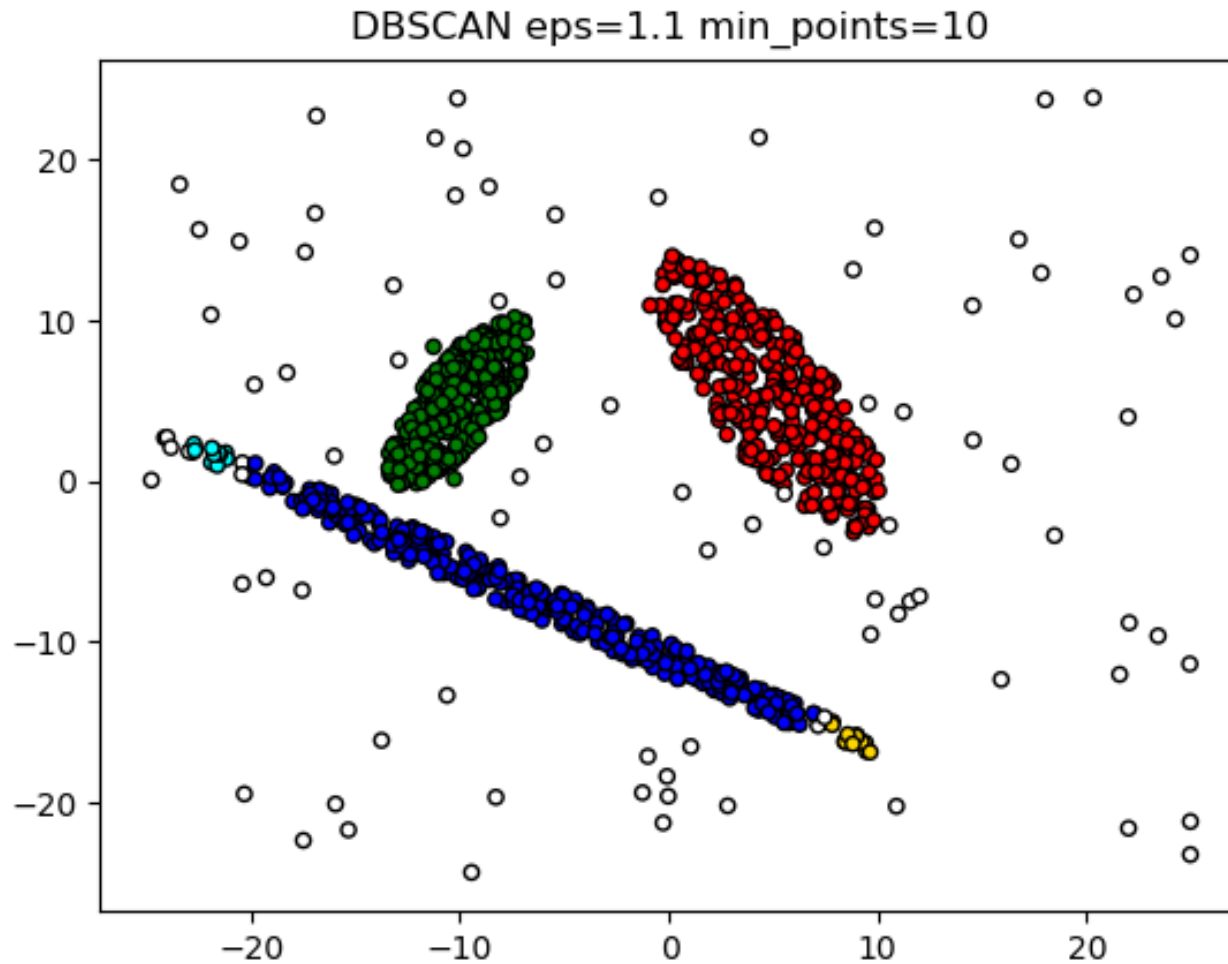
DBSCAN – przykład 2



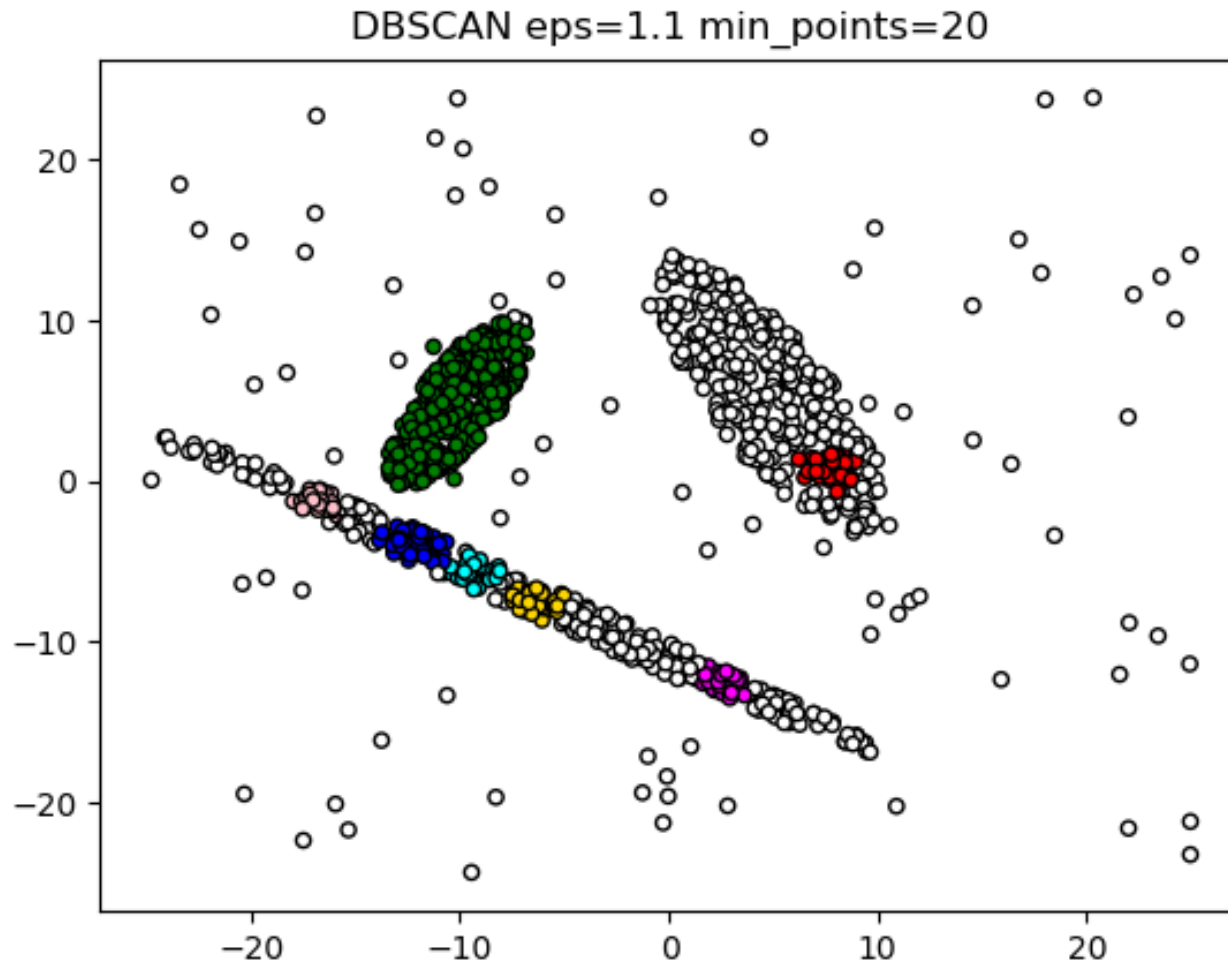
DBSCAN – przykład 3



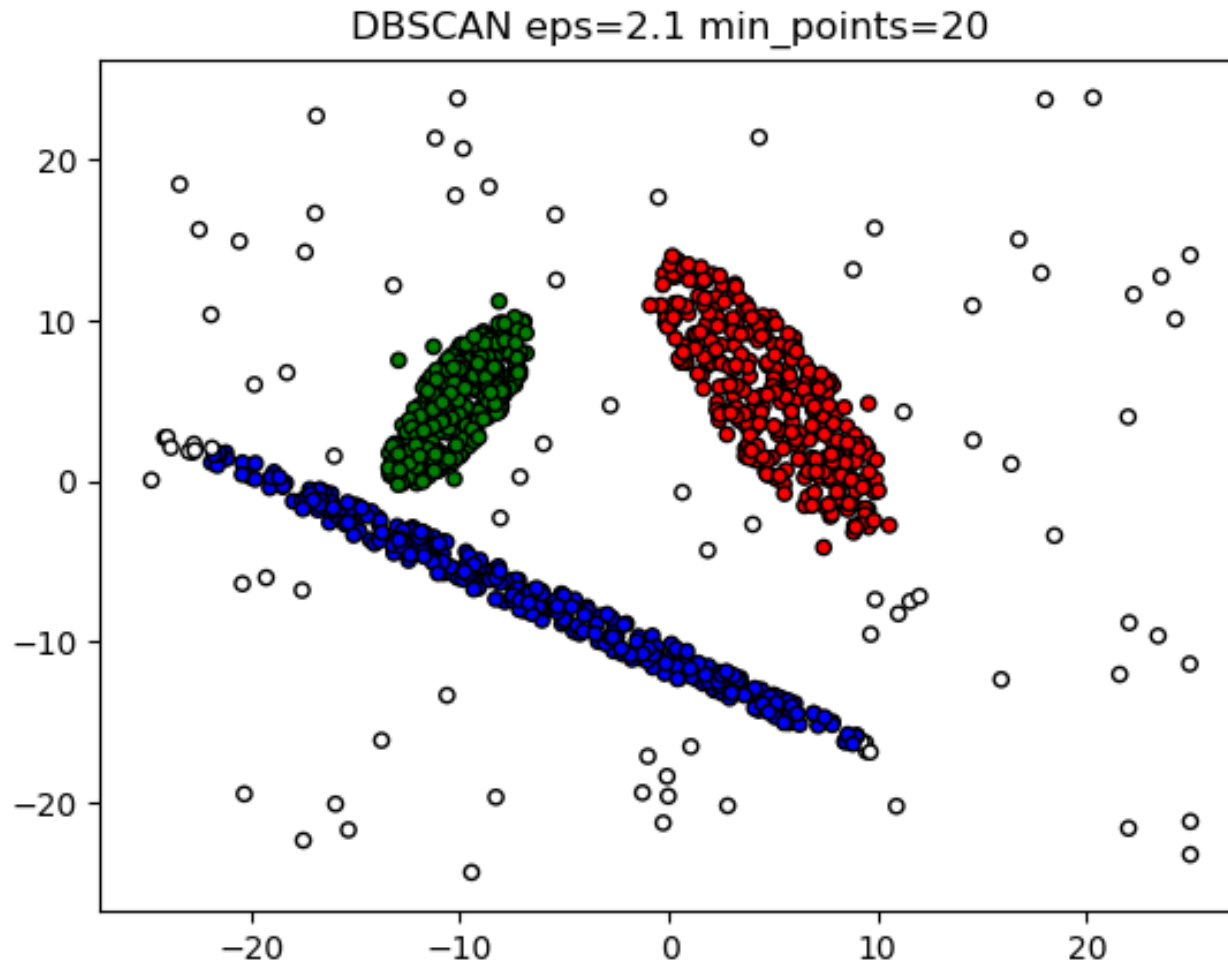
DBSCAN – przykład 4



DBSCAN – przykład 5

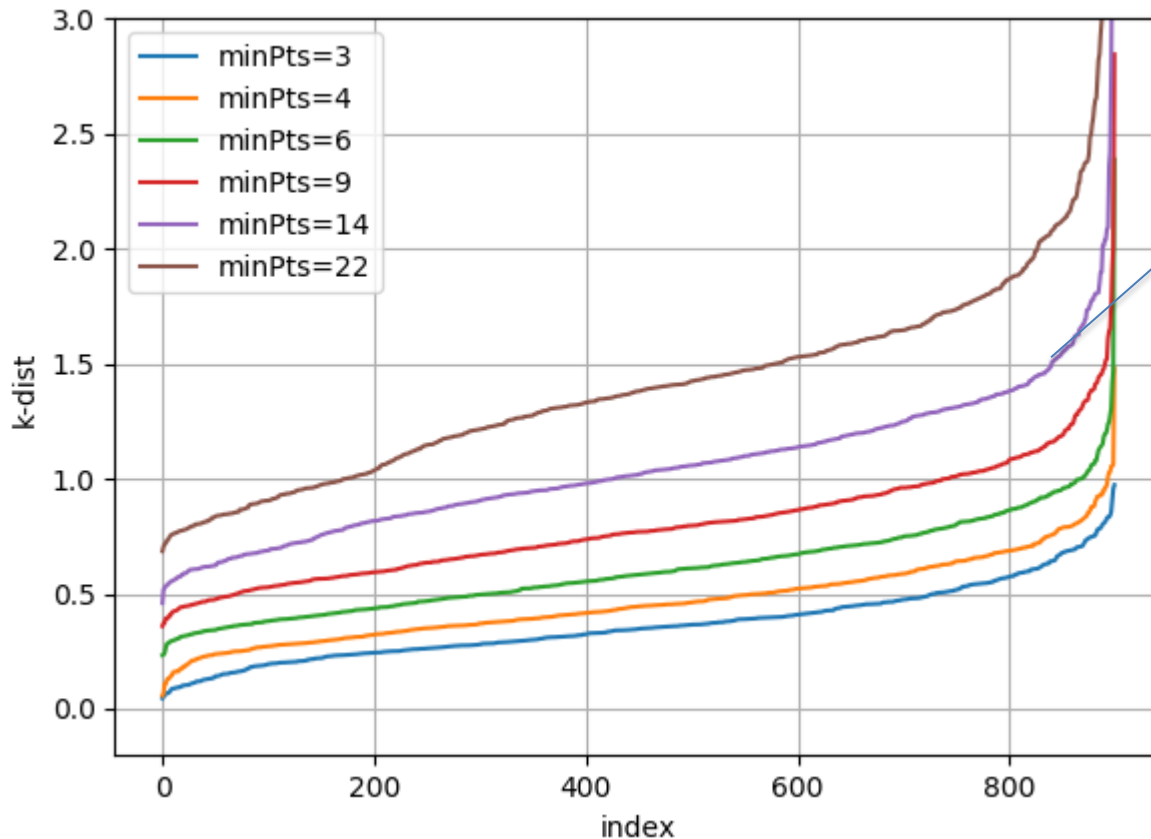


DBSCAN – przykład 6



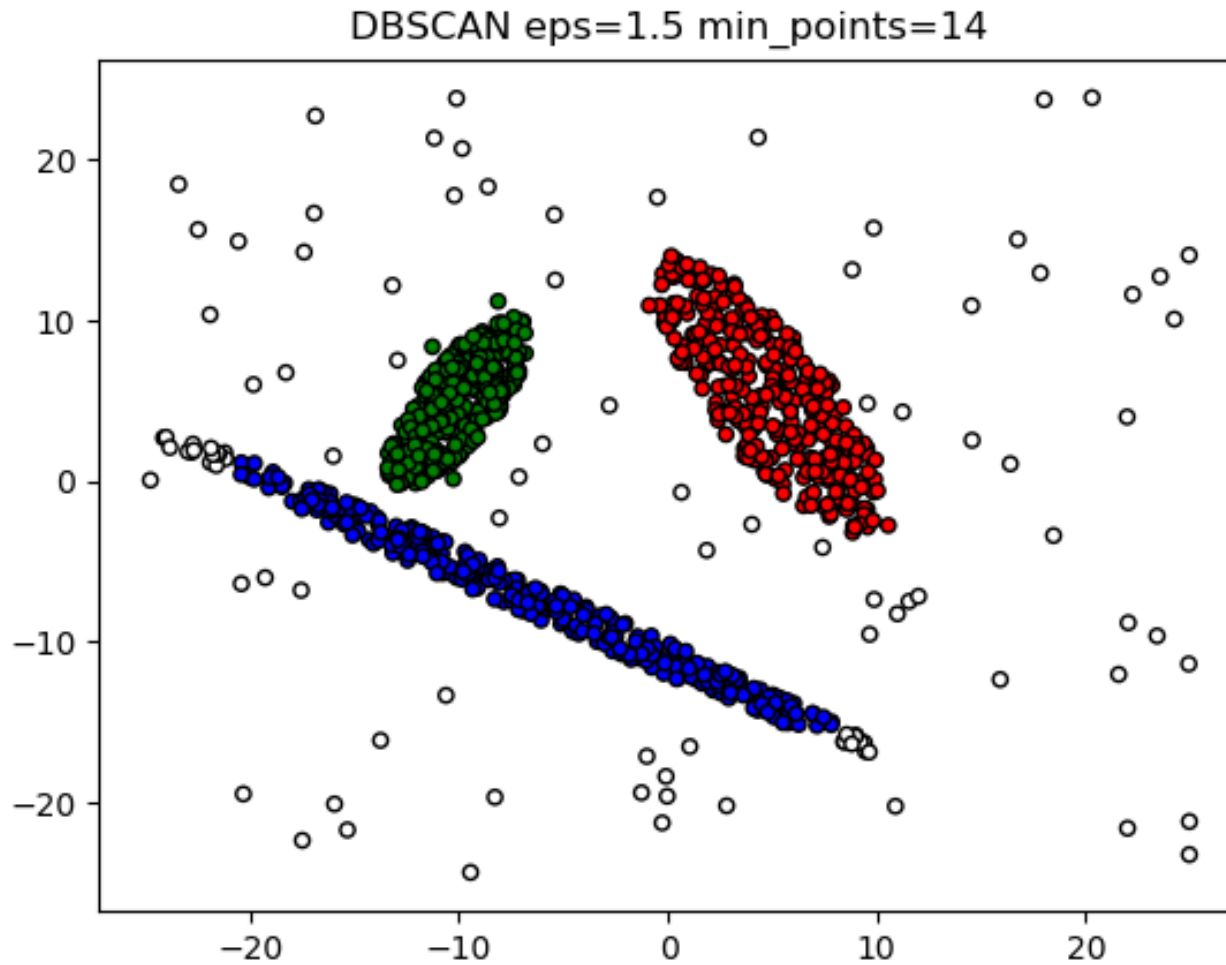
Jak dobrać $minPts$ i ϵ ?

Wykres k-dist pokazuje posortowany zbiór odległości do najbliższego k-tego sąsiada. Należy go sporządzić dla $k = minPts - 1$



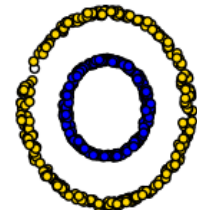
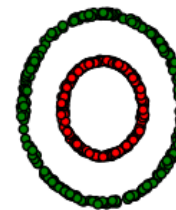
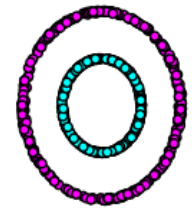
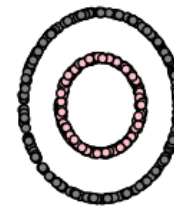
Odczytaj dist dla punktu przegięcia i wybierz go jako ϵ , np. $minPts = 14$, $\epsilon = 1.5$

Jak dobrać minPts i ϵ -wynik ?



DBSCAN - podsumowanie

- Nie wymaga podania liczby grup
- Może znaleźć grupy o dowolnym kształcie
- Może zidentyfikować odstające punkty (outliers)
- Jeżeli punkty są przetwarzane w takiej samej kolejności – zachowuje się deterministycznie
- Może dawać różne przydziały dla punktów brzegowych
- Dobór parametrów nie zawsze jest oczywisty (zwłaszcza dla większych wymiarów)
- Może używać kD-tree lub ball tree do przechowywania danych



Metody oceny grupowania

Podział metod

- Ocena **wewnętrzna** (ang. internal evaluation) – ocena na podstawie danych poddanych grupowaniu z użyciem funkcji oceny. Zwykle powiązanej z algorytmem grupowania.
- Ocena **zewnętrzna** (ang. external evaluation) – ocena na podstawie danych które nie były użyte do grupowania, np. znanych etykiet klas reprezentujących ground truth
- Ocena **ekspercka** (ang. manual evaluation)
- Ocena **pośrednia** (ang. indirect) – ocena użyteczności grupowania

Ocena wewnętrzna - indeks Daviesa-Boudina

- Oznaczmy przez $\bar{d}(i)$ średnią odległość od środka wewnątrz i -tej grupy:

$$\bar{d}(i) = \frac{\sum_{j=1}^m \mathbf{1}(Z_j = i) d(x_j, \bar{x}(i))}{\sum_{j=1}^m \mathbf{1}(Z_j = i)}$$

- Indeks Daviesa-Boudina jest zdefiniowany jako:

$$DB = \frac{1}{k} \sum_{i=1}^k \max_{j \neq i} \left\{ \frac{\bar{d}(i) + \bar{d}(j)}{d(\bar{x}(i), \bar{x}(j))} \right\}$$

- Indeks odzwierciedla założenie, że grupy powinny skupiać obiekty położone blisko siebie (mała wartość licznika), natomiast grupy powinny być oddalone (duża wartość mianownika). Niższa wartość DB oznacza lepsze grupowanie.

Ocena wewnętrzna – indeks Dunna

- Niech $d(i, j)$ oznacza odległość pomiędzy i -tą i j -tą grupą, natomiast $\bar{d}(i)$ odległość wewnątrz grupy. Obie te metryki mogą być zdefiniowane na różne sposoby (patrz hierarchiczne grupowanie dla $d(i, j)$).
- Indeks Dunna jest zdefiniowany jako:

$$D = \frac{\min\{d(i, j): 1 \leq i, j \leq k\}}{\max\{\bar{d}(i): 1 \leq i \leq k\}}$$

Oczekiwane jest duże podobieństwo/maća odległość wewnątrz grupy (mianownik mały) i duża odległość pomiędzy grupami, stąd większe wartości D są lepsze.

Ocena wewnętrzna – współczynnik kształtu

- Współczynnik kształtu (ang. Silhouette Coefficient) dla pojedynczej próbki jest zdefiniowany jako

$$s = \frac{b - a}{\max\{a, b\}}$$

gdzie:

- a – to średnia odległość pomiędzy próbką i innymi punktami z tej samej grupy
- b – to średnia odległość pomiędzy próbką i punktami z **najbliżej położonej** grupy
- Współczynnik przyjmuje wartości pomiędzy -1 (dla złego grupowania) a 1 (dla idealnego grupowania)
- Współczynnik jest wyższy, dla gęstych i wypukłych grup (np. otrzymanych z k-means, natomiast może być zawodny dla niewypukłych grup, które może zwrócić np. DBSCAN.

Ocena wewnętrzna – indeks Calińskiego Harabasza

Oznaczmy przez $\bar{x}(i)$ środek i -tej grupy, a przez c punkt środkowy całego zbioru danych.

Macierze W_k i B_k definiują kowariancję wewnątrz grup (W_k) oraz pomiędzy grupami (B_k)

- $W = \sum_{i=1}^k \sum_{x \in C_i} (x - \bar{x}(i)) (x - \bar{x}(i))^T$
- $B = \sum_{i=1}^k (\bar{x}(i) - c) (\bar{x}(i) - c)^T$

Wskaźnikiem określającym stopień rozproszenia (dyspersji) dla macierzy kowariancji Σ jest jej ślad $tr(\Sigma)$, czyli suma elementów na przekątnych (wariancji).

Indeks CH jest zdefiniowany jako:

$$CH(k) = \frac{tr(B)}{tr(W)} \frac{m - k}{k - 1}$$

Indeks rośnie dla gęstych (mała wartość $tr(W)$) i dobrze odseparowanych (duża wartość $tr(B)$) grup.

Ocena zewnętrzna

- W ocenie zewnętrznej używa się etykiet klas (np. wprowadzonych manualnie dla części danych) lub też informacji o rozkładach użytych do generacji syntetycznych zbiorów danych.
- Zdefiniowane wskaźniki mierzą, podobieństwa pomiędzy strukturą klas i grup (stopień zgodności etykiet klas i grup).

Ocena zewnętrzna – indeks Randa

Indeks Randa (William Rand) pozwala ocenić zgodność dwóch podziałów zbioru na rozłączne podzbiory

Niech $O = \{o_1, \dots, o_n\}$ będzie zbiorem obiektów, $X = \{X_1, \dots, X_r\}$ oraz $Y = \{Y_1, \dots, Y_r\}$ dwoma podziałami.

Rozważmy parę różnych elementów z O : (o_i, o_j) . Mamy cztery możliwe przypadki:

- A. o_i, o_j należą do jednego ze zbiorów w X oraz do jednego ze zbiorów w Y
- B. o_i, o_j należą do dwóch różnych zbiorów w X oraz różnych zbiorów w Y
- C. o_i, o_j należą do tego samego zbioru w X i różnych zbiorów w Y
- D. o_i, o_j należą do różnych zbiorów w X i tego samego zbioru w Y

Oznaczmy przez a, b, c, d liczby par spełniające warunki A, B, C i D.

$$RI = \frac{a + b}{a + b + c + d} = \frac{a + b}{\binom{n}{2}}$$

RI zwraca wartości pomiędzy 0 i 1 i może być interpretowany jako prawdopodobieństwo, że dowolnie wybrana para jest w analogiczny sposób sklasyfikowana w obu grupowaniach.

Skorygowany indeks Randa (ang. Adjusted Rand Index) wprowadza poprawkę uwzględniającą prawdopodobieństwo dwa algorytmy grupowania zachowując się losowo równocześnie rozdzielią parę lub dołączą do jednej grupy (rozkład hipergeometryczny).

$$ARI = \frac{RI - E(RI)}{\max(RI) - E(RI)}$$

[<https://link.springer.com/content/pdf/10.1007%2F978-3-319-08075-5.pdf>]

[https://en.wikipedia.org/wiki/Rand_index]

Entropia (przypomnienie)

- Niech X będzie zmienną losową przyjmującą k dyskretnych wartości, każdą z nich z prawdopodobieństwem $P(X = x_k)$
- Entropia zdefiniowana jest jako:

$$H(X) = - \sum_{i=1}^k P(X = x_k) \log_2 P(X = x_k)$$

- Entropia (mierzona w bitach) podaje minimalną liczbę bitów niezbędną do zakodowania losowo wybranej wartości X .
- Jeśli $P = 0$, przyjmuje się, że $0 \cdot \log_2 0 = 0$

Symbol	Częstość
A	0.17
B	0.33
C	0.50

- $H = -(0.17 \cdot \log_2(0.17) + 0.33 \cdot \log_2(0.33) + 0.5 \cdot \log_2(0.5)) = 1.46$
- Proponowane kodowanie (średnio 1.55 bitu/symbol):
C: 0
B: 10
A: 11

Entropia warunkowa (przypomnienie)

Entropia warunkowa obliczana jest dla dwóch zmiennych Y i X

$$H(Y|X) = - \sum_{i=1}^k \sum_{j=1}^r P(Y = y_i \wedge X = x_j) \log_2 \frac{P(Y = y_i \wedge X = x_j)}{P(X = x_j)}$$

przekształcając:

$$H(Y|X) = - \sum_{i=1}^k \sum_{j=1}^r P(Y = y_i|X = x_j) P(X = x_j) \log_2 P(Y = y_i|X = x_j)$$

$$H(Y|X) = - \sum_{j=1}^r P(X = x_j) \sum_{i=1}^k P(Y = y_i|X = x_j) \log_2 P(Y = y_i|X = x_j)$$

- Wydzielany jest podzbiór $X \times Y$, gdzie $X = x_j$
- Obliczane jest entropia w tym podzbiorze
- Obliczana jest suma z wagami $P(X = x_j)$

Ocena zewnętrzna V-measure

Założmy, że mamy podział zbioru obserwacji na klasy (ground truth) $C = C_1, \dots, C_n$ oraz zbiór grup $K = K_1, \dots, K_m$ będących efektem działania algorytmu.

Podziałowi K stawiane są dwa wymagania:

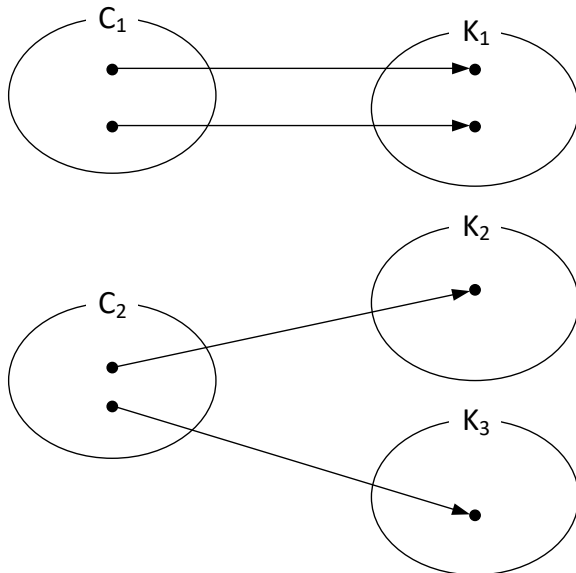
- Jednorodności (ang. homogeneity) – grupowanie powinno przydzielić do jednej grupy **wyłącznie** elementy jednej klasy (nie powinny być pomieszane).
- Zupełności (ang. completeness) – **wszystkie** elementy jednej klasy powinny zostać przydzielone do **jednej grupy**

Jednorodność

- Formalnie jednorodność jest zdefiniowana jako

$$h = 1 - \frac{H(C|K)}{H(C)}$$

- W idealnym przypadku $H(C|K) = 0$, ponieważ dla każdej klasy C_i i grupy K_j prawdopodobieństwo $P(C_i|K_j)$ wynosi 1 lub 0. Wtedy $h = 1$. Czynniki $H(C)$ jest użyty do skalowania. Entropia $H(C)$ zależy od rozkładu liczby obserwacji w klasach.



Przykład

Grupy zawierają jednorodne obserwacje (pochodzące z tej samej klasy).

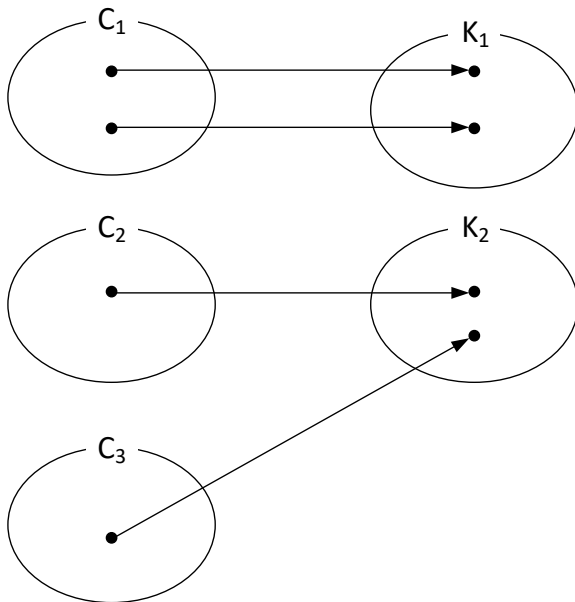
Np. $P(C_1|K_1) = 1$, $P(C_2|K_3) = 1$
Wówczas $h = 1$

Zupełność

- Zupełność jest zdefiniowana jako

$$c = 1 - \frac{H(K|C)}{H(K)}$$

- Jest to miara zdefiniowana symetrycznie, odpowiadająca wymaganiu, aby wszystkie elementy danej klasy trafiły do tej samej grupy



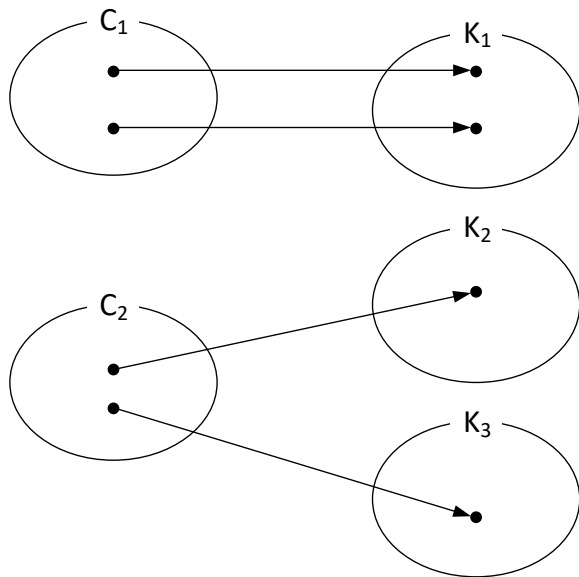
Przykład

Wszystkie elementy klas trafiły do tych samych grup.

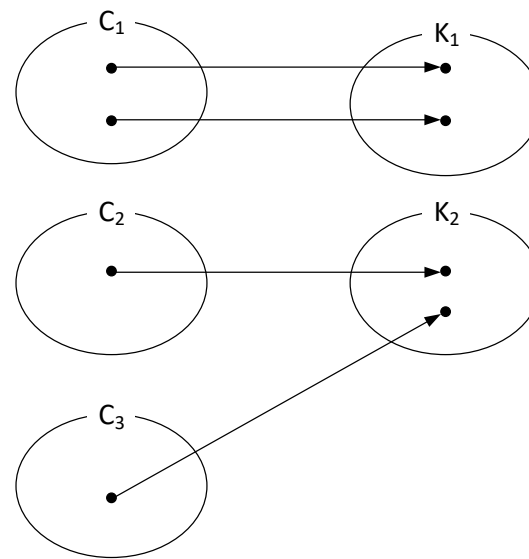
Np. $P(K_1|C_1) = 1$, $P(K_2|C_2) = 1$

Wówczas $c = 1$

Ocena zewnętrzna V-measure



homogeneity $h = 1$
completeness $c = 2/3$



homogeneity $h = 2/3$
completeness $c = 1$

- Miara V-measure jest średnią harmoniczną h i c :

$$v = 2 \frac{h \cdot c}{h + c}$$

[Andrew Rosenberg and Julia Hirschberg: V-Measure: A conditional entropy-based external cluster evaluation measure
<http://www.aclweb.org/anthology/D07-1043>]

Ocena zewnętrzna – wzajemna informacja

- Dla dwóch zmiennych losowych X i Y wskaźnik wzajemnej informacji (ang. mutual information) mierzy w bitach ile na podstawie jednej zmiennej można uzyskać informacji o drugiej. Formalnie zdefiniowany jest jako

$$I(Y, Z) = \sum_{y \in Y} \sum_{z \in Z} p(y, z) \log_2 \left(\frac{p(y, z)}{p(y)p(z)} \right)$$

- Niech $Y = \{Y_1, \dots, Y_i, \dots, Y_r\}$ oraz $Z = \{Z_1, \dots, Z_j, \dots, Z_r\}$ będą rozłącznymi podziałami zbioru n obiektów. Wówczas

- $p(i) = \frac{|Y_i|}{n}$ będzie prawdopodobieństwem, że losowo wybrany obiekt znajdzie się w Y_i ,
- $p(j) = \frac{|Z_j|}{n}$ że znajdzie się w Z_j
- $p(i, j) = \frac{|Y_i \cap Z_j|}{n}$ – że losowo wybrany obiekt znajdzie się w Y_i oraz Z_j

- Indeks wzajemnej informacji MI dla dwóch podziałów na grupy jest zdefiniowany jako:

$$MI(Y, Z) = \sum_{i=1}^{|Y|} \sum_{j=1}^{|Z|} p(i, j) \log_2 \left(\frac{p(i, j)}{p(i)p(j)} \right)$$

- Indeks przyjmuje wartości z przedziału $[0, 1]$.
- Wraz ze wzrostem liczby (małych) grup indeks przyjmuje większe wartości, stąd korekta uwzględniająca przypadkową zgodność przydziału (adjusted mutual information)

$$AMI = \frac{MI - E(MI)}{\max(MI) - E(MI)}$$