

# Eksploracja danych

## KLASTERYZACJA I SEGMENTACJA cz. 2

***Wojciech Waloszek***

*wowal@eti.pg.gda.pl*

***Teresa Zawadzka***

*tegra@eti.pg.gda.pl*

*Katedra Inżynierii Oprogramowania  
Wydział Elektroniki, Telekomunikacji i Informatyki  
Politechnika Gdańska*



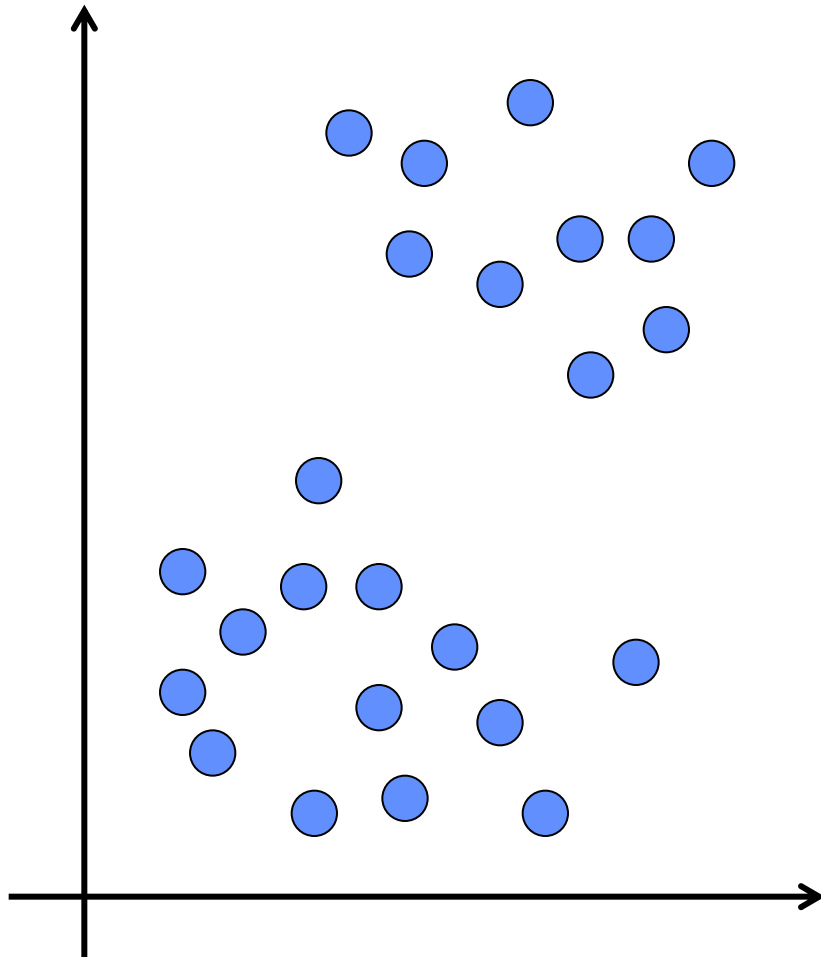
# Inne ważne metody grupowania

- K-means:
  - Iteracyjne grupowanie na zasadzie odnajdywania centrów podzbiorów przykładów,
- EM:
  - Próba odnalezienia optymalnego globalnie podziału poprzez poszukiwanie parametrów rozkładów Gaussa.

# K-means

- Poszukuje „centrów” dla zadanej liczby klastrów,
- Pierwszy dobór centrów odbywa się losowo, kolejne iteracje algorytmu przesuwają „centra”, aby lepiej odzwierciedlały one podział przestrzeni przykładów,
- Kryterium stopu jest zazwyczaj pewien próg odległości, o którą zostały przesunięte „centra” w kolejnej iteracji.

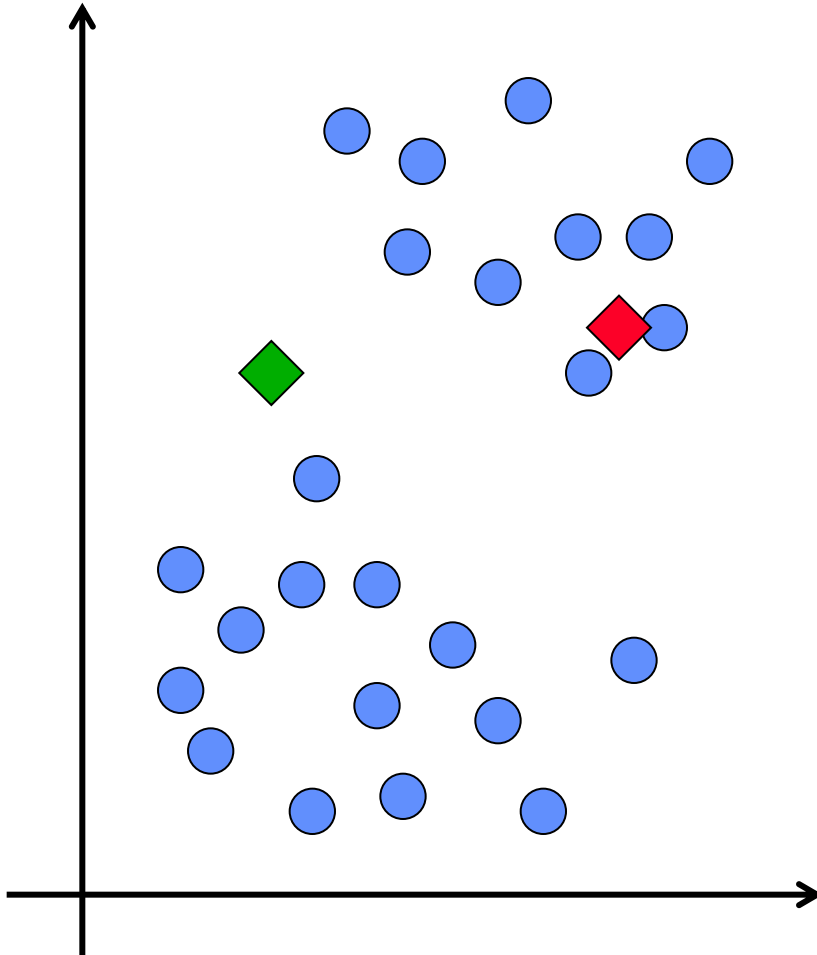
# K-means – przykład działania



- Przestrzeń przykładów przedstawiamy sobie jako  $n$ -wymiarową przestrzeń euklidesową, w której przykłady reprezentowane są jako punkty

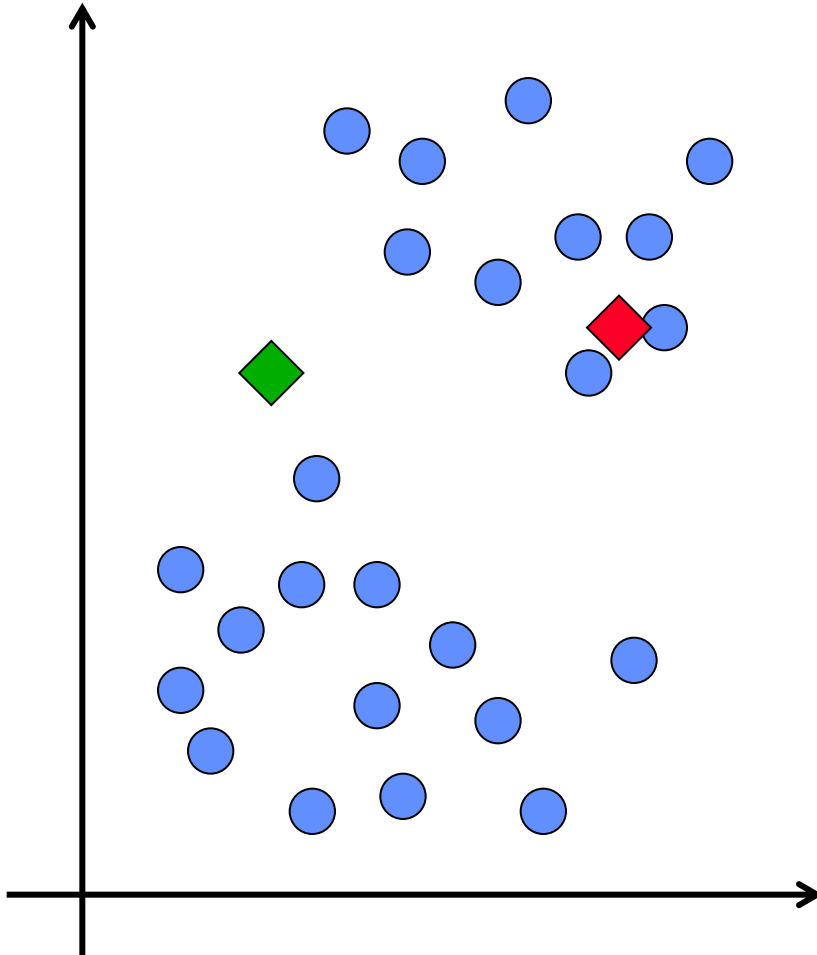
# K-means – przykład działania (2)

- Na początek losujemy położenie  $k$  centrów (tutaj przyjmujemy  $k = 2$ )



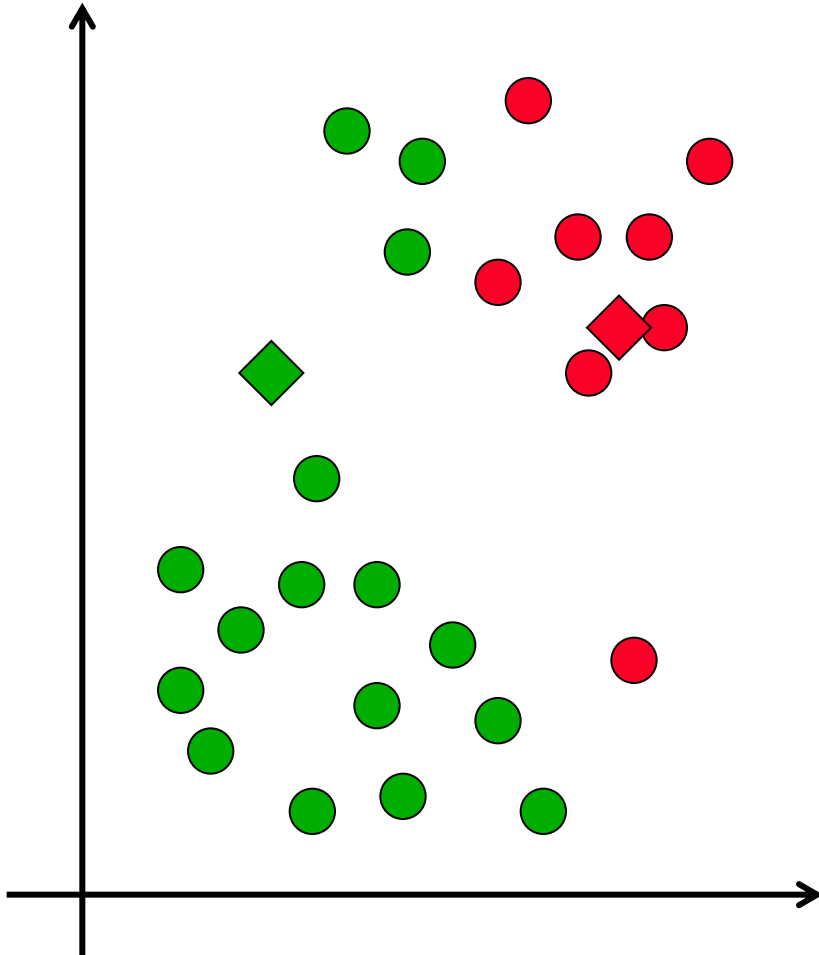
# K-means – przykład działania (3)

- Oznaczamy, które przykłady leżą najbliżej którego „centrum”

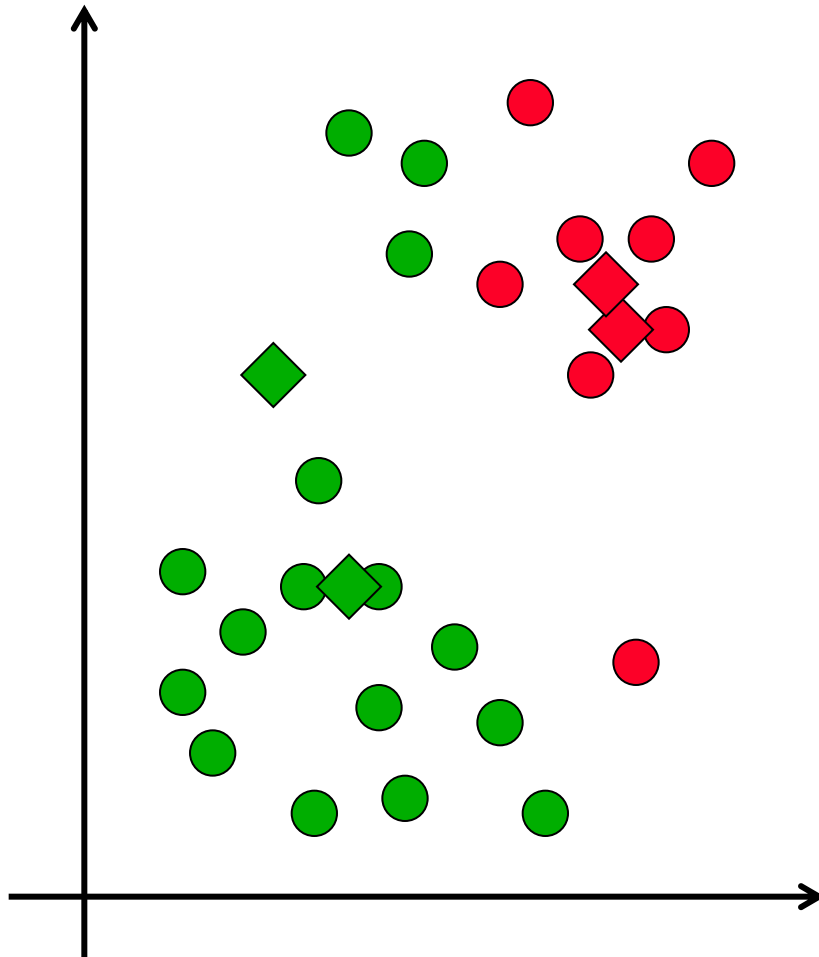


# K-means – przykład działania (3)

- Oznaczamy, które przykłady leżą najbliżej którego „centrum”



# K-means – przykład działania (4)

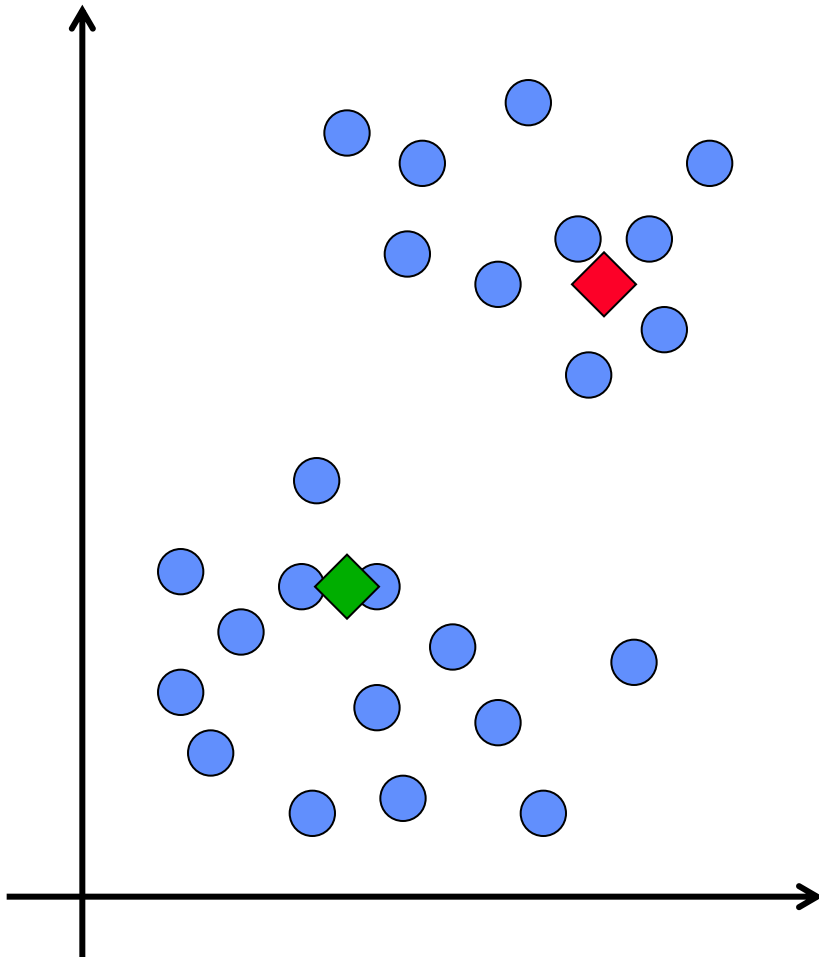


- Przenosimy dotychczasowe „centra” do „środków” przypisanych im zbiorów przykładów



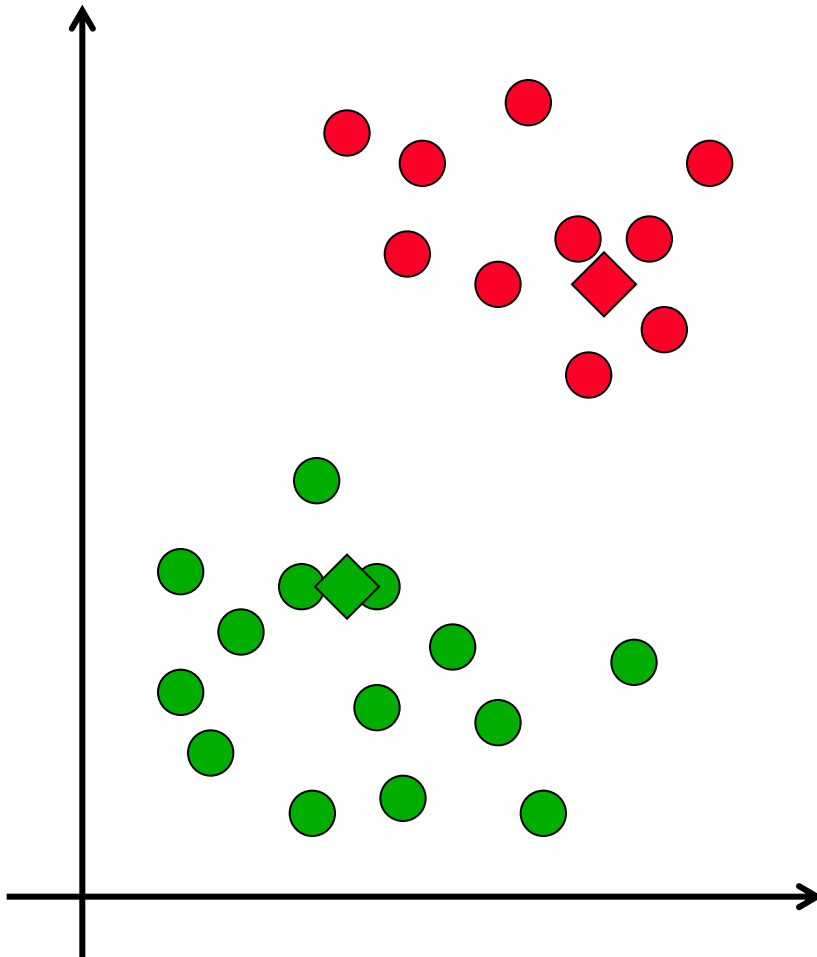
# K-means – przykład działania (5)

- Ponownie przypisujemy przykłady „centrom”



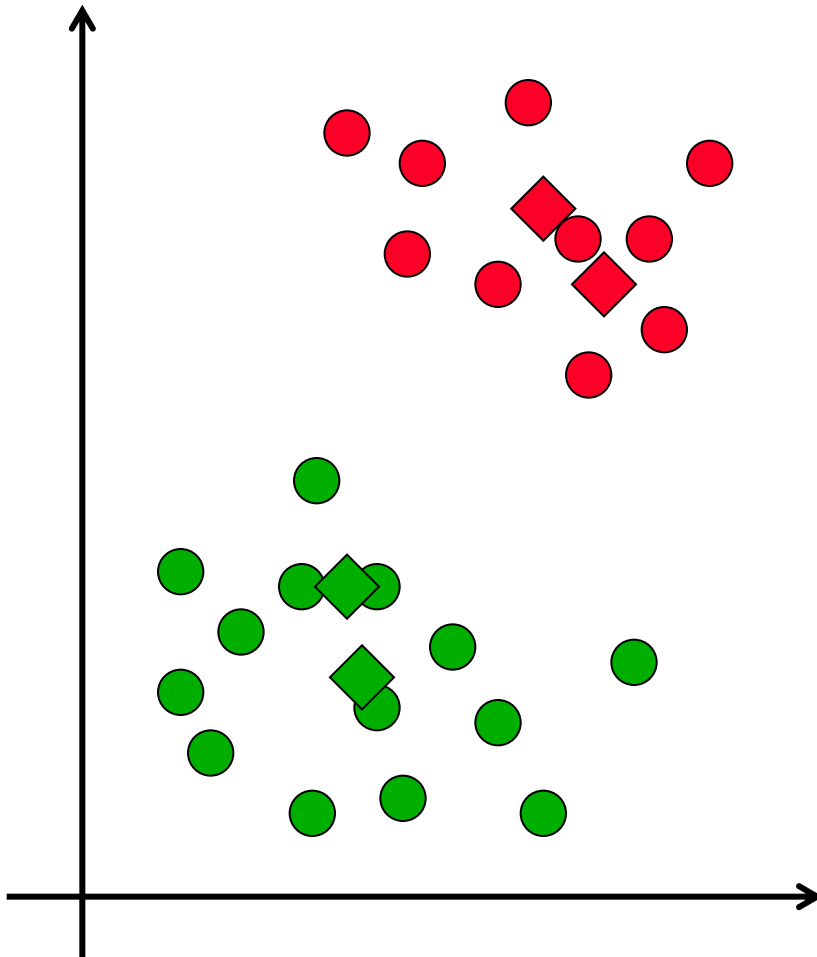
# K-means – przykład działania (5)

- Ponownie przypisujemy przykłady „centrom”



# K-means – przykład działania (6)

- I ponownie przesuwamy „centra” ...

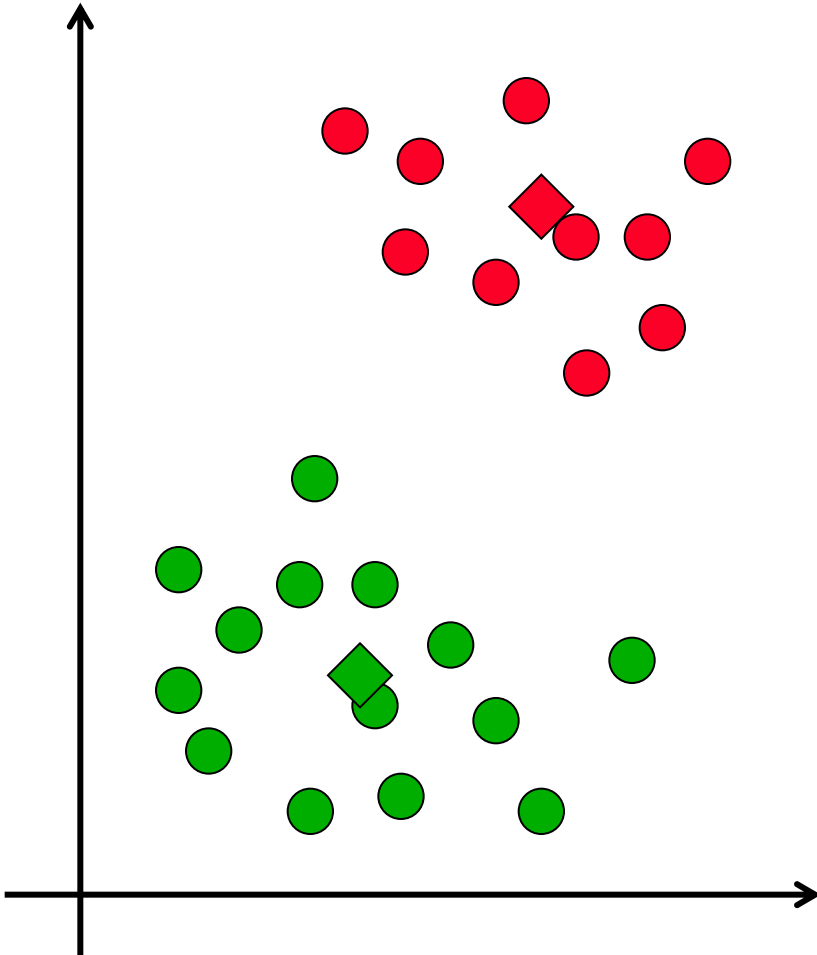


# K-means

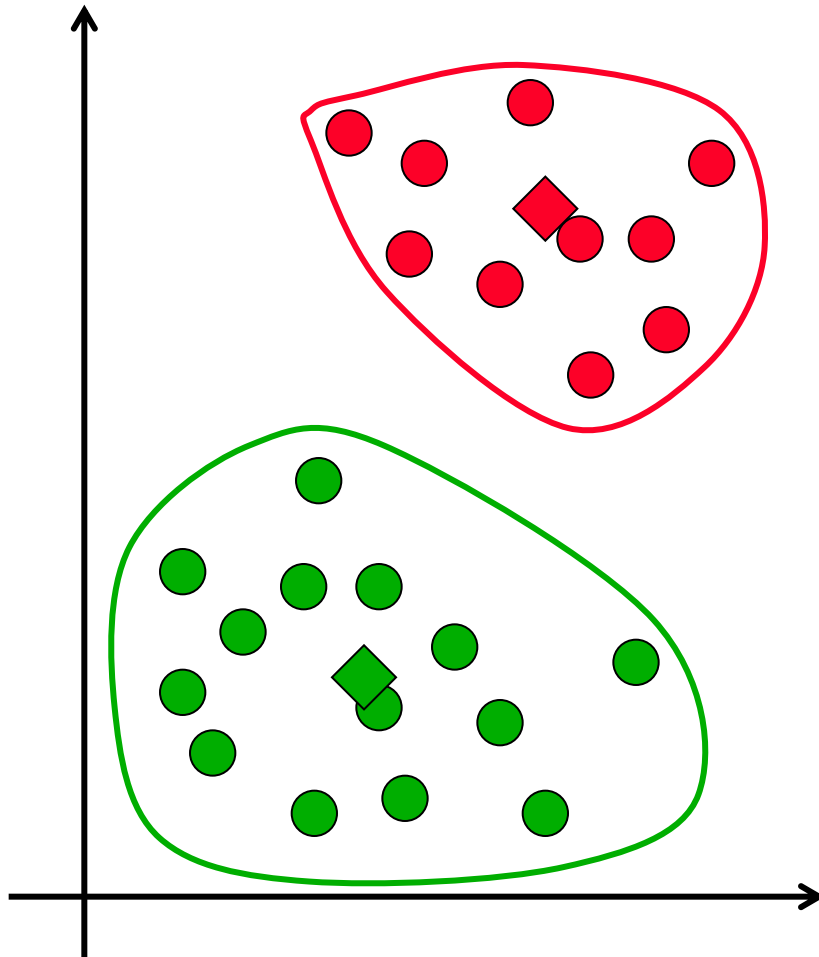
- K-means może wymagać wykonanie kilku powtórzeń całego algorytmu, jako że sam algorytm ma tendencję do znajdowania „minimum lokalnego”,
- K-means wymaga podanej z góry liczby tworzonych klastrów, zaś podział generowany przez algorytm jest płaski, co w niektórych zastosowaniach może być wadą,
- Z tego względu stosuje się hierarchiczną wersję algorytmu K-means.

# Hierarchiczny k-means

- Algorytm rozpoczyna pracę, jak zwykły k-means dla  $k = 2$



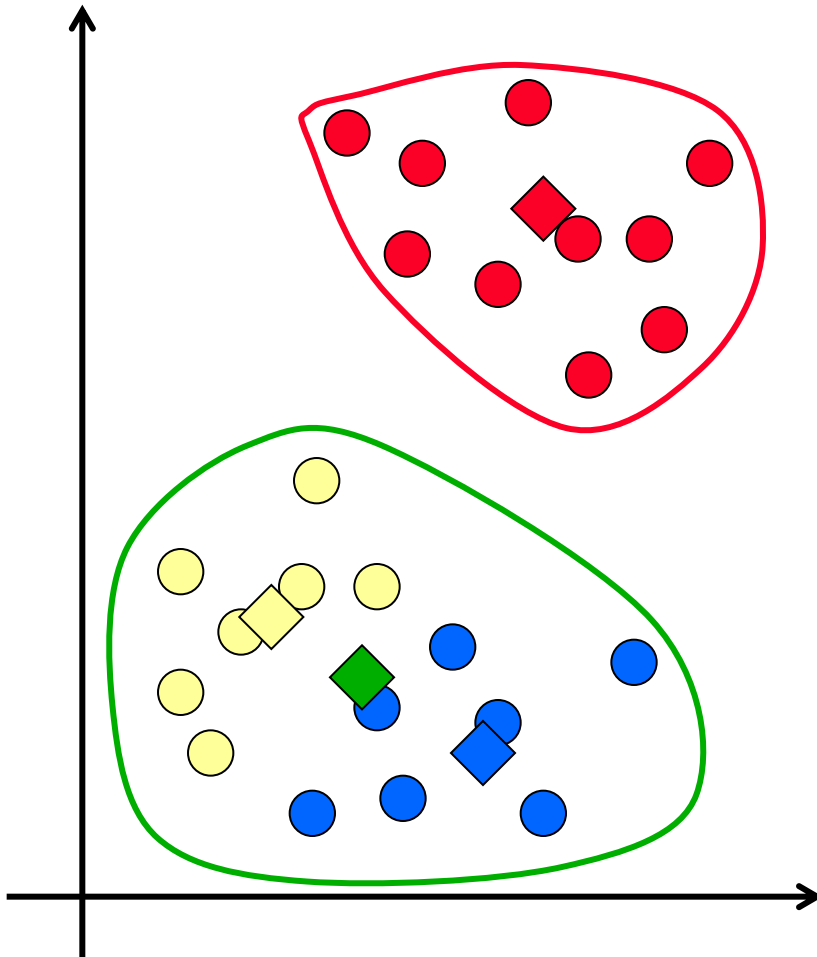
## Hierarchiczny k-means (2)



- Dla „zbyt dużych” klastrów przeprowadzany jest proces ponownego podziały na zasadzie analogicznej, jak dla algorytmu k-means

# Hierarchiczny k-means (3)

- Wynikiem działania algorytmu jest hierarchiczny podział przestrzeni przykładów na grupy

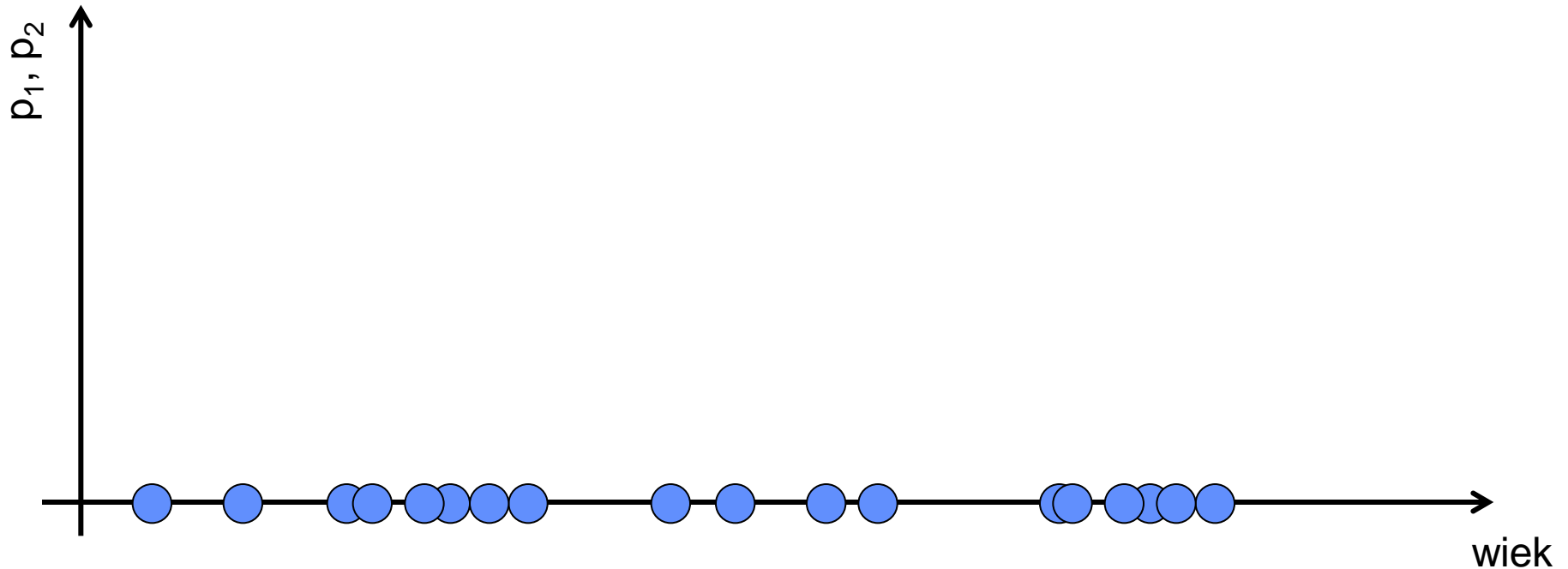


# EM (expectation-maximization)

- Algorytm EM zakłada, że zbiór uczący to „mieszanina” przykładów z  $k$  kategorii,
- EM zakłada, że każda kategoria przykładów charakteryzuje się pewnym rozkładem wartości atrybutów (najczęściej przyjmowany jest rozkład normalny),
- EM, podobnie jak  $k$ -means, działa iteracyjnie, dostosowując parametry rozkładów do przykładów obejmowanych przez kategorie na kolejnym etapie przetwarzania.

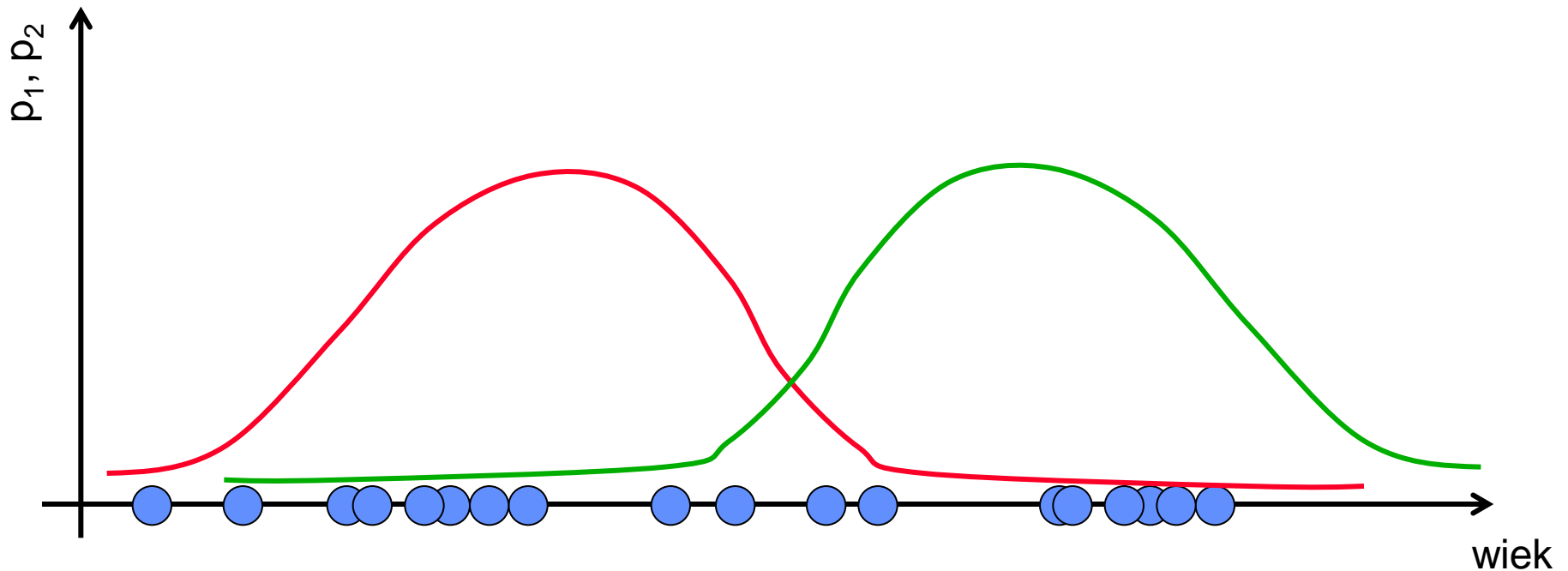


# EM – przykład działania



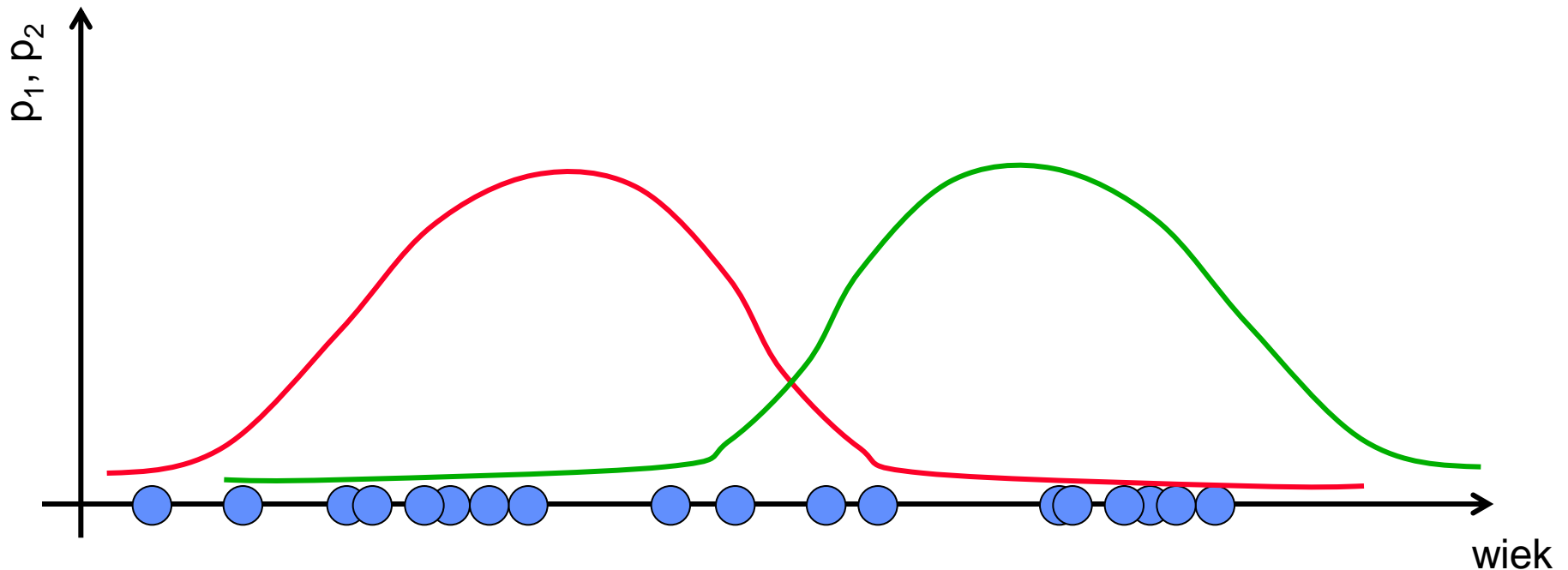
- Przestrzeń przykładów możemy rozumieć podobnie, jak w przypadku k-means, tutaj wyobraźmy sobie 1 wymiar (wartość atrybutu wiek)

## EM – przykład działania (2)



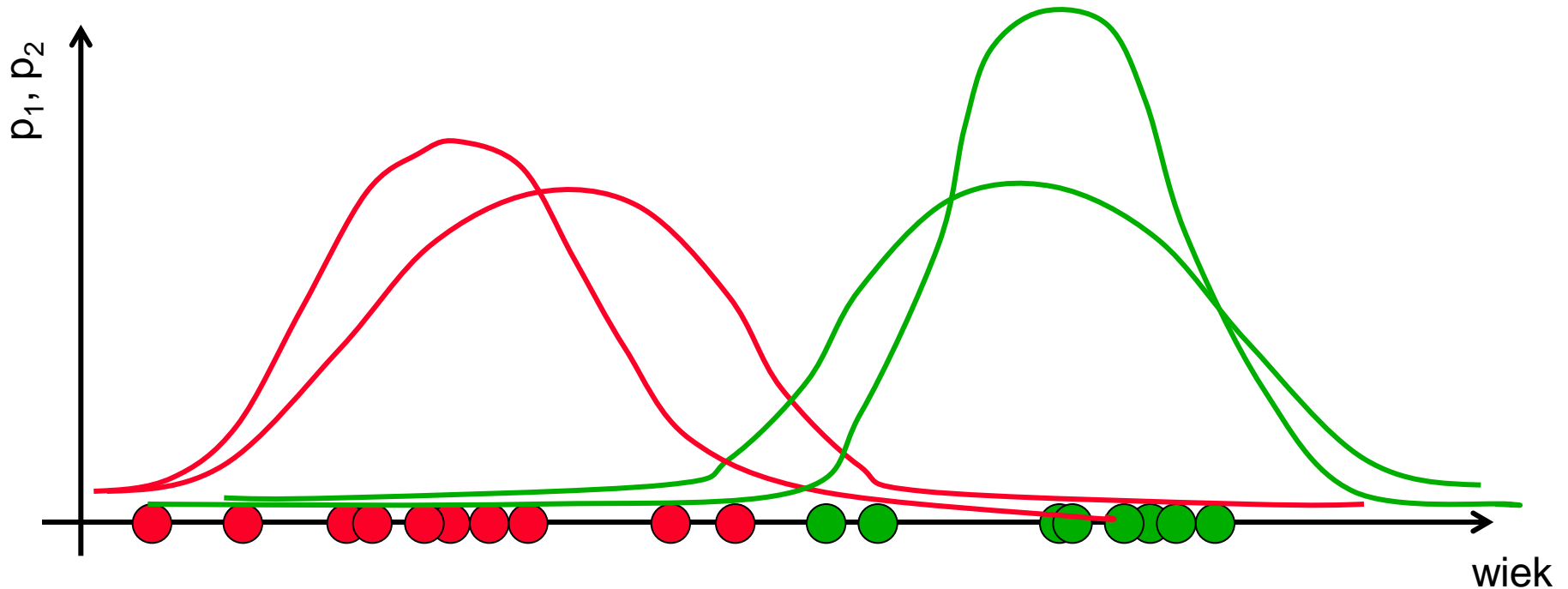
- Na początek przyjmujemy dwa „losowe” rozkłady...

## EM – przykład działania (3)



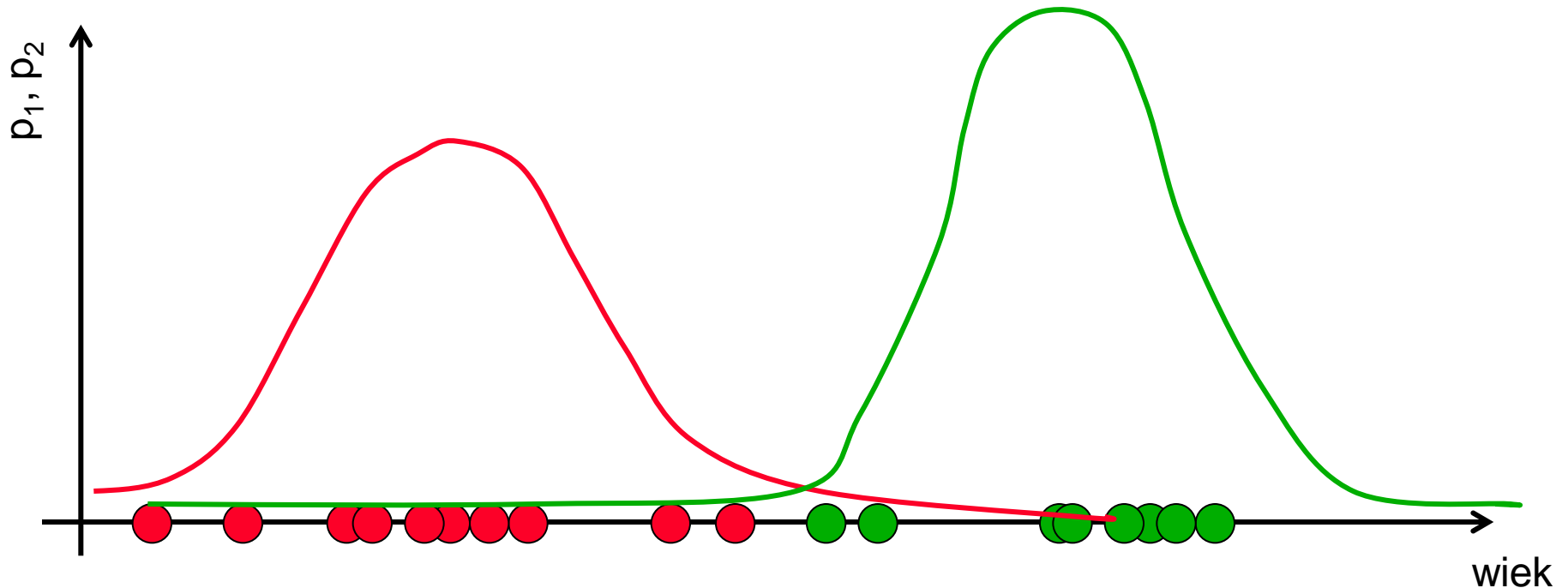
- Następnie przypisujemy przykłady do odp. kategorii, po czym modyfikujemy parametry rozkładów

## EM – przykład działania (3)



- Następnie przypisujemy przykłady do odp. kategorii, po czym modyfikujemy parametry rozkładów

## EM – przykład działania (4)



- Proces powtarzamy iteracyjnie, aż do osiągnięcia kryterium stopu (niewielka zmiana parametrów)

# EM – podsumowanie

- EM można rozszerzyć na więcej wymiarów, co jest proste przy założeniu niezależności wartości odpowiednich atrybutów,
- W przypadku odnalezienia zależności między wartościami atrybutów można zastosować rozkład dwu- (lub więcej) wymiarowy i szukać jego parametrów
- Podobnie jak k-means, EM szuka „minimum lokalnego”, stąd konieczne może być wykonanie algorytmu kilkukrotnie

# Dziękujemy za uwagę

Zapraszamy na wykład:

**REGUŁY ASOCJACYJNE**