

Rhizobium

Comamonadaceae

Grupowanie

Bradyrhizobium

Gallionellaceae

Jarosław Jasiewicz
Eksploracja danych i Uczenie maszynowe

Geoinformacja program magisterski
Specjalność Geoinformatyka

Dlaczego klasyfikujemy dane

- Klasyfikacja danych ma na celu zredukowanie złożoności danych. Zamiast wielu obiektów, każdy opisany kilku-kilkunastoma parametrami mamy kilka klas, a każdy obiekt jest przypisany do jednej (niekiedy więcej) klas
- **Klasyfikacje nadzorowane** polegają na przypisaniu nowego obiektu do już istniejącego zestawu klas
- **Klasyfikacje nienadzorowane** mają na celu wykrycie w danych ukrytych, nieoczywistych struktur. Krokiem w klasyfikacji danych jest grupowanie lub analiza skupień (*clustering*)
- Problem terminologiczny: w języku angielskim termin *classify data* jest niejednoznaczny. Może oznaczać zarówno proces klasyfikacji jak i **utajniania** danych.

Grupowanie/Analiza skupień

- Analiza skupień to proces przypisywania obiektów do niezdefiniowanych a priori grup na podstawie analizy struktury danych
- Obiekty w skupieniach wykazują tendencję do wzajemnego podobieństwa, a obiekty w różnych skupieniach wykazują tendencję do niepodobieństwa
- Podstawą analizy skupień jest koncepcja niepodobieństwa pomiędzy obiektami

Koncepcja niepodobieństwa

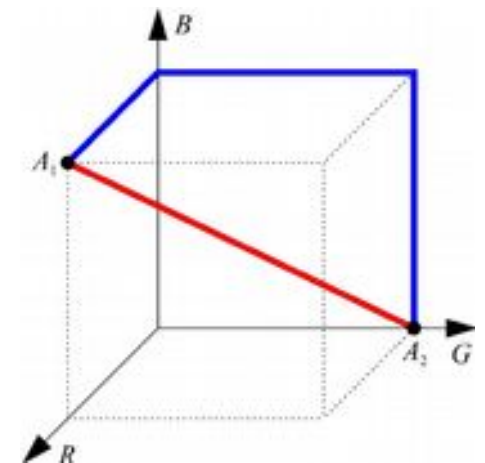
- Koncepcja podobieństwa/niepodobieństwa jest kluczowa dla analizy skupień, tak aby podobne obiekty były klasyfikowane do tych samych skupień, a niepodobne do różnych
- Pojęcie intuicyjnie zrozumiałe, ale trudne do wyrażenia matematycznie
- Niepodobieństwo jest proste to wyrażenia pomiędzy obiektami opisanymi jako punkty w przestrzeni dwu- lub trójwymiarowej poprzez pojęcie **odległości**

Koncepcja odległości

- W przestrzeni fizycznej odległość pomiędzy dwoma punktami to najkrótsza droga pomiędzy dwoma punktami
- Przy założeniu braku przeszkód – odległość liczymy jako najkrótszą możliwą odległość w przestrzeni – generalizowaną jako odległość euklidesową (*euclidean*)
- Przy istnieniu przeszkód – odległość liczymy jako najkrótszą możliwą drogę pomiędzy dwoma punktami – generalizowaną jako odległość miejską (*manhattan*).
- Odległość może być liczona w przestrzeni lub na dowolnej płaszczyźnie (np. na sferze jako tzw wielkie koło)

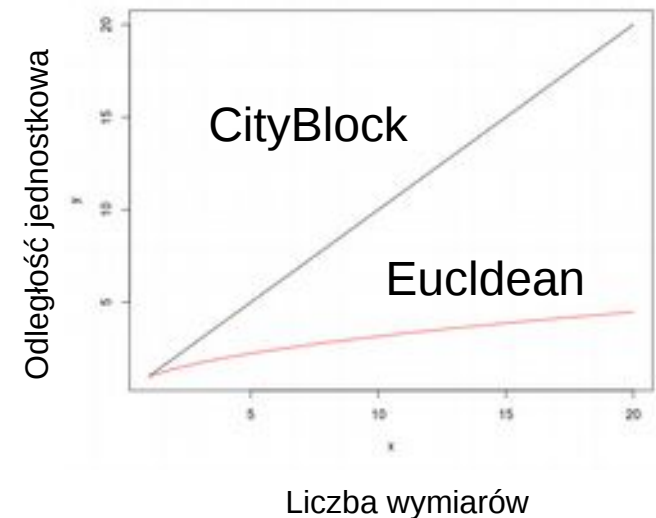
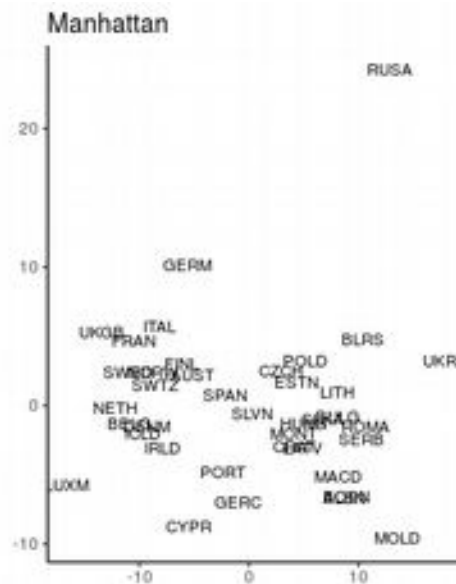
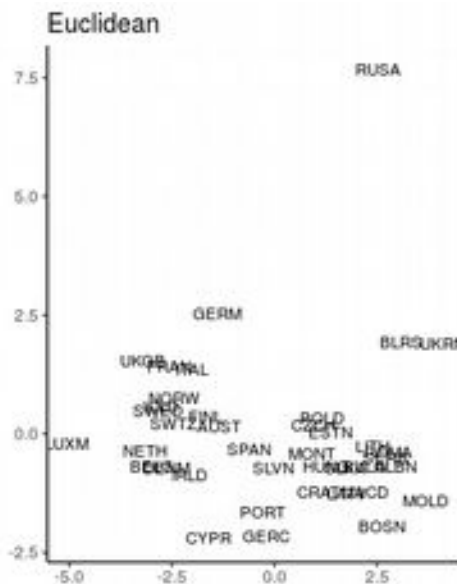
Odległość euklidesowa i miejska

- Odległość euklidesowa pomiędzy dwoma punktami to długość linii łącząca te dwa punkty
 - Dla dwóch wymiarów: $d_{eucl} = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$
 - Postać ogólna: $d_{eucl} = \sqrt{\sum (x_i - y_i)^2}$
- Odległość miejska pomiędzy dwoma punktami to suma odległości w każdym z wymiarów z osobna
 - Dla dwóch wymiarów: $d_{manh} = |x_2 - x_1| + |y_2 - y_1|$
 - Postać ogólna $d_{manh} = \sum |x_i - y_i|$



Niepodobieństwo a przekleństwo wymiarowości

- Cechy miary euklidesowej powodują że wraz ze wzrostem wymiarowości wpływ kolejnych wymiarów jest coraz mniejszy, przy dużej liczbie wymiarów należy rozważyć stosowanie odległości miejskiej



Metryka

- Przestrzeń metryczna to przestrzeń w której odległości pomiędzy wszystkimi obiektami są zdefiniowane, Zbór wszystkich odległości zwane są **metryką zbioru**
- Koncepcja metryki jest generalizacją odległości euklidesowa w 3-wymiarowej przestrzeni euklidesowej
- Każda metryka musi spełniać następujące aksjomaty:
 - Nieujemność: $d(a,b) \geq 0$
 - Identyfikacja: $d(a,b) = 0 \iff a = b$
 - Symetria: $d(a,b) = d(b,a)$
 - Nierówność trójkątna $d(a,b) \leq d(a,c) + d(b,c)$

Wektory i Normy

- Definiując obiekt jako zbiór cech opisanych wartościami możemy przedstawić go jako punkt w wielowymiarowej (n-wymiarowej) przestrzeni
- Norma to funkcja, która przypisuje długość do wektora wyznaczonego przez ten punkt (stąd **wektor cech**)
- W przestrzeni n-wymiarowej norma Euklidesowa (L2) to najkrótsza odległość pomiędzy początkiem układu a punktem:

Gdzie $x_1 \dots x_n$ kolejne wymiary

$$\|x\|_2 := \sqrt{x_1^2 + x_2^2 + \dots + x_n^2}$$

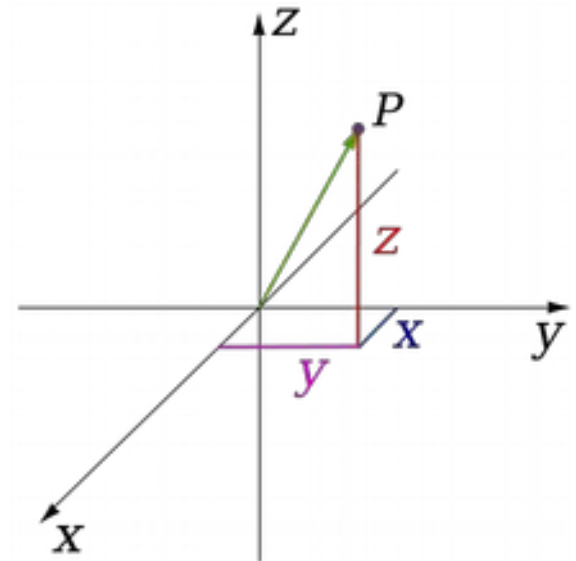
- Norma miejska - Manhattan (L1) suma współrzędnych

$$\|x\|_1 := x_1 + x_2 + \dots + x_n$$

- P-norma

$$\|x\|_p := \sum (|x_1|^p + |x_2|^p + \dots + |x_n|^p)^{1/p}$$

$$\|x\|_{inf} := \max(x_1 + x_2 + \dots + x_n)$$



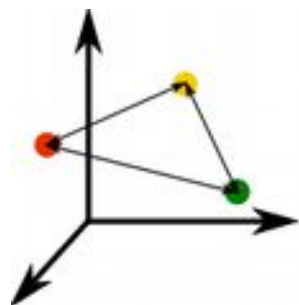
Odległość a niepodobieństwo

- Jeżeli w analizie danych obiekty opisane są poprzez ich wektory cech, niepodobieństwo pomiędzy nimi utożsamia się z odległością – metryką euklidesową. W żargonie termin *distance* używa się jako synonimu niepodobieństwa
- Nie wszystkie miary niepodobieństwa dają się wyrazić jako odległość
- W Spatial Data Science pojęcie odległości jest ambiwalentne: odległość pomiędzy obiektami to odległość w przestrzeni geograficznej czy niepodobieństwo?
- Bezpieczne terminy:

Przeźrzeń Geograficzna i kartzjańska	Odległość <i>Distance</i>	Bliskość <i>Proximity</i>
Przeźrzeń Informacyjna	Niepodobieństwo <i>Dissimilarity</i>	Podobieństwo <i>Similarity</i>

Inne miary niepodobieństwa

- Istnieje ponad 200 miar podobieństwa i niepodobieństwa.
- Dobór miar zależy od:
 - Ilości wymiarów
 - Rodzaju atrybutów (komplementarne, binarne itp.)
 - Rodzaju problemu
- Źródłem miar jest:
 - Norma
 - Przecięcie zbiorów
 - Ilość informacji (entropia)
 - Iloczyn skalarny
 - Test statystyczny
- Jeżeli wartość miary niepodobieństwa może być interpretowana ilościowo, mówimy że miara ma semantykę (np. odległość euklidesowa da się interpretować ilościowo)



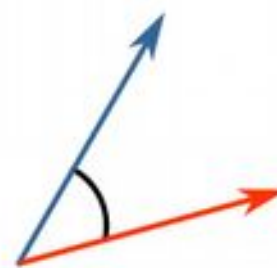
1



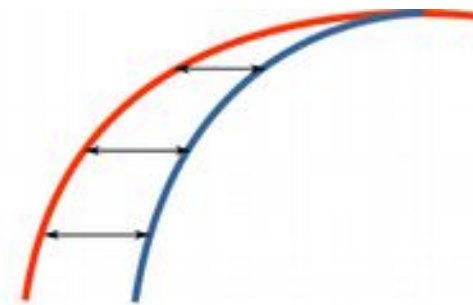
2



3



4



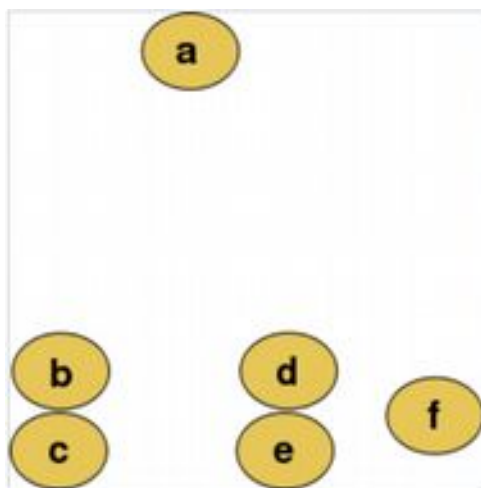
5

Wybrane miary

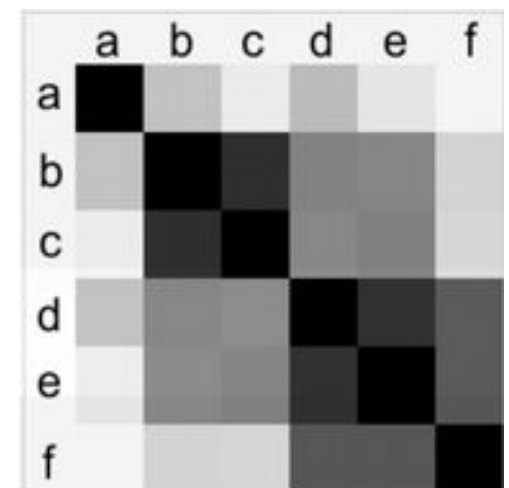
Miara	Rodzaj	Opis	Wzór
Współczynnik korelacji	S	Współczynnik korelacji Pearsona	
Kosinusowa	D	Kąt pomiędzy wektorami o niezerowej długości. Pokazuje orientację a nie natężenie. Stosuje się do wielowymiarowych zbiorów np. tekstów	$\text{cosine}(x, y) = \frac{\sum x_i, y_i}{\sqrt{\sum X^2} \sqrt{\sum Y^2}}$
Mahalanobis	D	Miara dostosowuje się do liniowej kombinacji wymiarów, określa ile odchyłeń standardowych jest obiekt od średniej dla każdego z wymiarów	$d_{mah} = \sqrt{(X - Y) S^{-1} (X - Y)^T}$
Canberra	D	Standaryzowana [0,1] odmiana odległości miejskiej	$d_{canb} = \sum \frac{ x_i - y_i }{(x_i + y_i)}$
Trójkątna	D	Standaryzowana [0,1] odmiana odległości euklidesowej, wysoka zgodność z JSD	$d_{tri} = \sqrt{\frac{1}{2} \sum \frac{(x_i - y_i)^2}{(x_i + y_i)}}$
Jensen-Shannon	D	Współdzielona ilość informacji dla dwóch rozkładów zmiennej kategoryzowanej, entropia wzajemna	$d_{jsd} = \sqrt{H\left(\frac{X+Y}{2}\right) - \frac{1}{2}[H(X) + H(Y)]}$
Jaccard	S	Miara wielkości przecięcia dwóch zbiorów, stosowana dla atrybutów binarnych	$s_{jaccard} = \sum \frac{ X \cdot Y }{X} + Y - X \cdot Y $
Rużicka	D	Miara niezgodności rozkładów	$s_{roz} = \frac{\sum \min(X, Y)}{\sum \max(X, Y)}$

Macierz niepodobieństwa

- Macierz niepodobieństwa – zestawienie każdy z każdym wartości niepodobieństwa pomiędzy obiektami. W praktyce macierz dwu- wymiarowa, na przekątnej wartości 0 (aksjomat identyczności) i symetryczna (aksjomat symetryczności metryki)
- W przypadku nawet niewielkich zbiorów macierze przedstawia się w formie wizualizacji graficznej zamiast zbioru liczb
 - Skalowanie wielowymiarowe
 - Mapy ciepła
 - Grafy

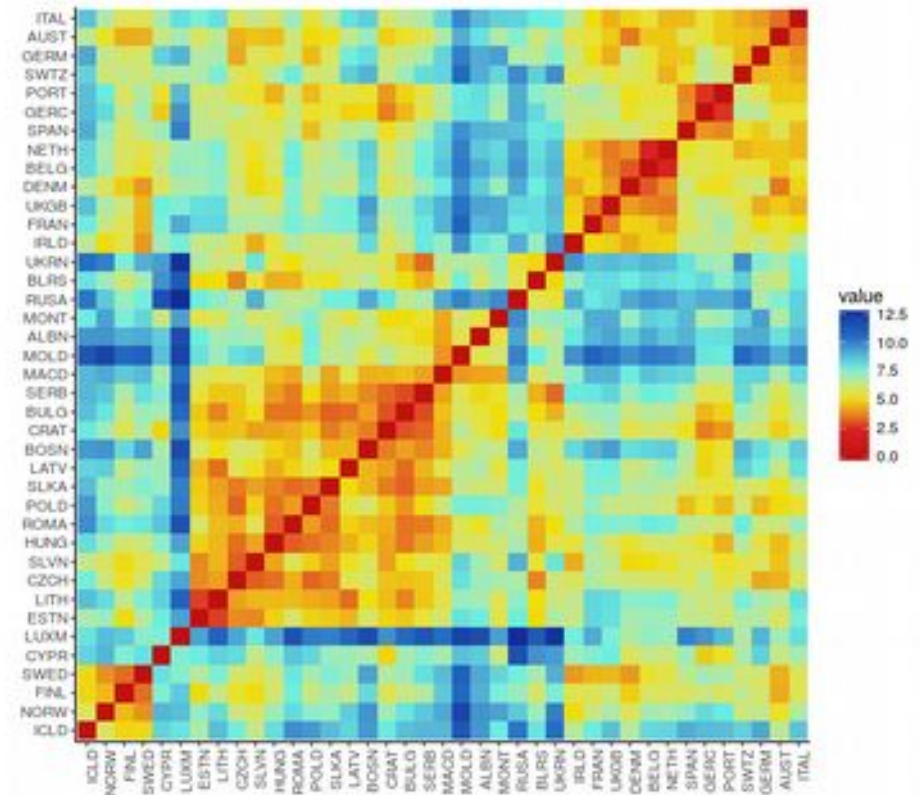
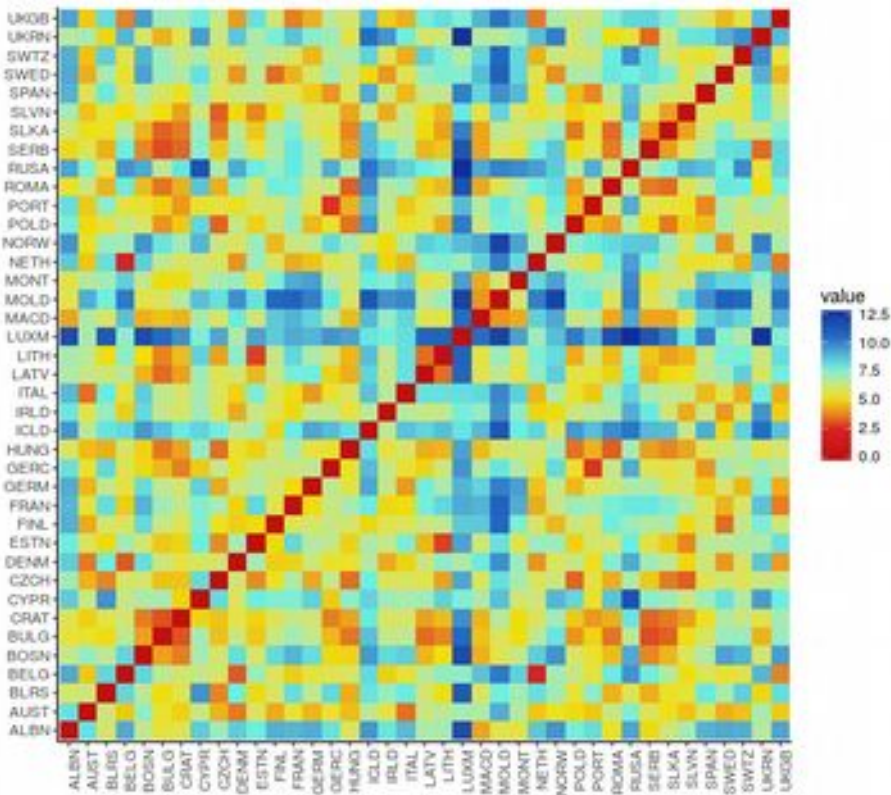


	a	b	c	d	e	f
a	0	184	222	177	216	231
b	184	0	45	123	128	200
c	222	45	0	129	121	203
d	177	123	129	0	46	83
e	216	128	121	46	0	83
f	231	200	203	83	83	0



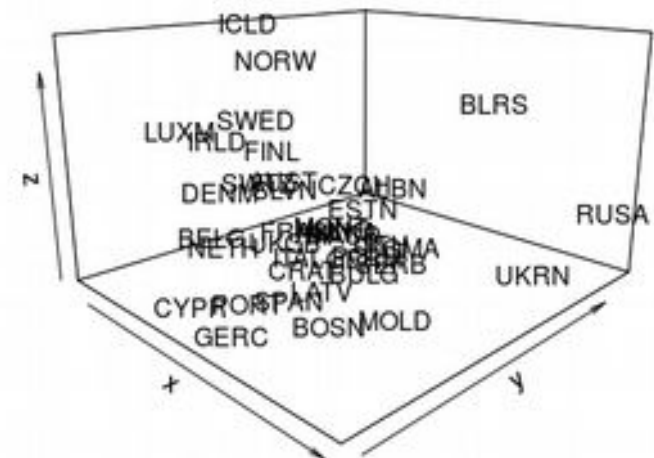
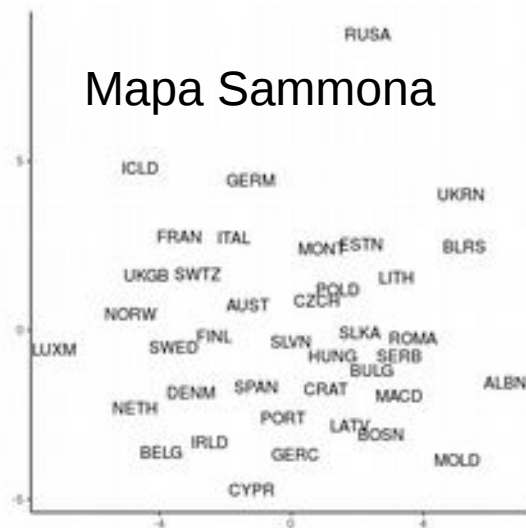
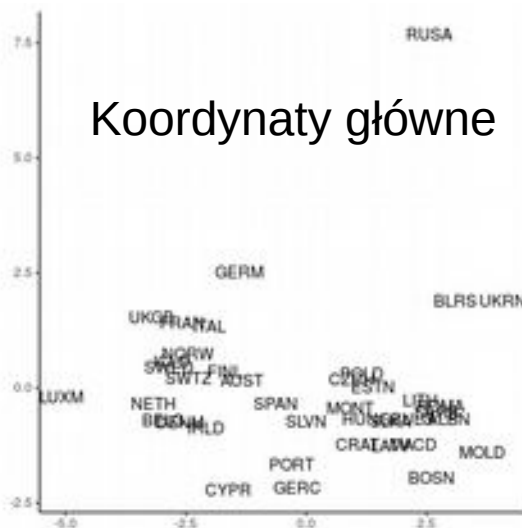
Mapy ciepła

- Mapy ciepła to wizualizacja macierzy gdzie niepodobieństwo wyrażone jest kolorem. Uporządkowanie mapy wg niepodobieństwa ciepła pozwala wykryć struktury w danych



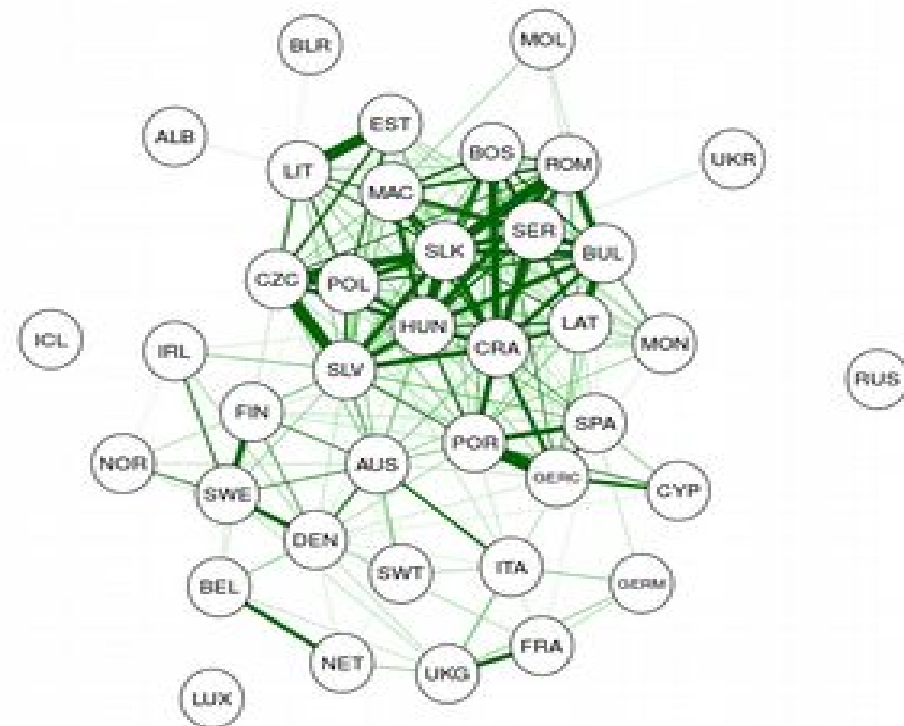
Skalowanie wielowymiarowe

- To forma prezentowania niepodobieństwa pomiędzy obiektami poprzez rzutowanie ich do przestrzeni 2-u lub trójwymiarowej, w taki sposób aby minimalizować różnice niepodobieństwa pomiędzy wartościami z oryginalnej wielowymiarowej przestrzeni a nowej przestrzeni zredukowanej
- Skalowanie wielowymiarowe stosuje się również do konwersji pomiędzy atrybutami komplementarnymi a wektorami cech



Grafy

- Grafy to forma prezentacji (nie)podobieństwa w formie obiektów rozmieszczonych w przestrzeni (wierzchołki albo węzły) a łączących je linii (krawędzie), których waga reprezentuje podobieństwo między obiektami

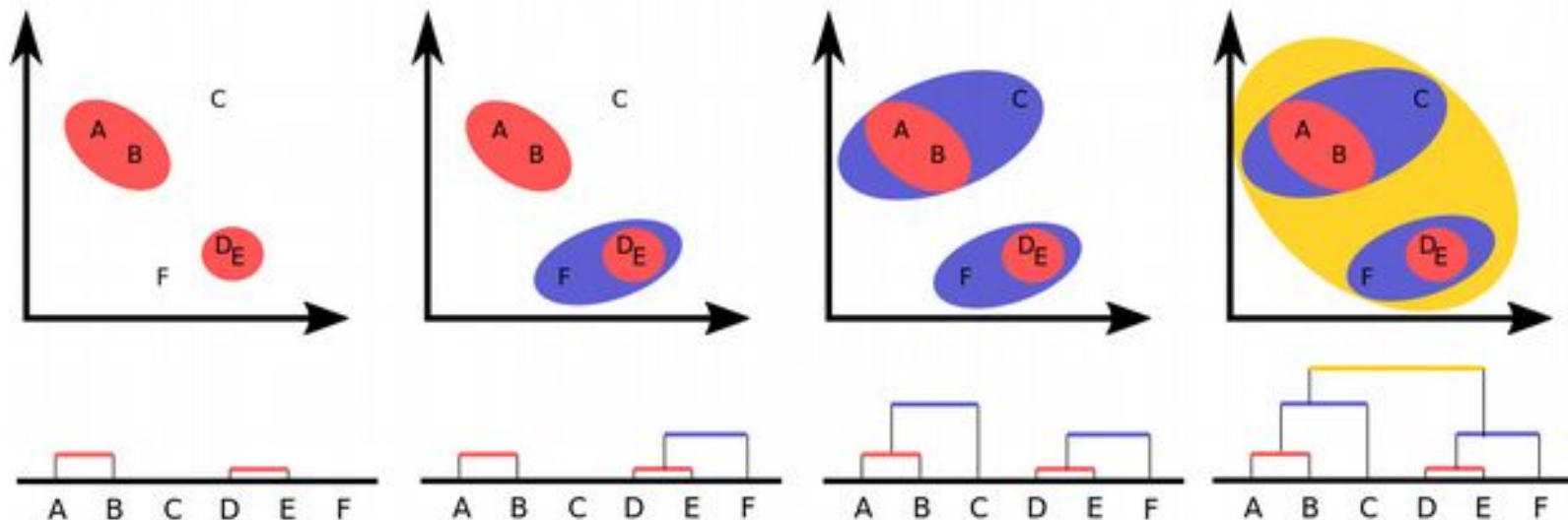


Algorytmy grupowania

- Grupowanie to podział zbioru danych na grupy rozłączne i wewnętrznie spójne
- Stosuje się kilka różnych metod
 - **Metody hierarchiczne**
 - **Metody rozdzielające (partitioning)**
 - **Metody rozmyte**
 - **Grupowanie probabilistyczne**
 - Metody gęstościowe
 - Ulepszony hierarchiczny (BIRCH)

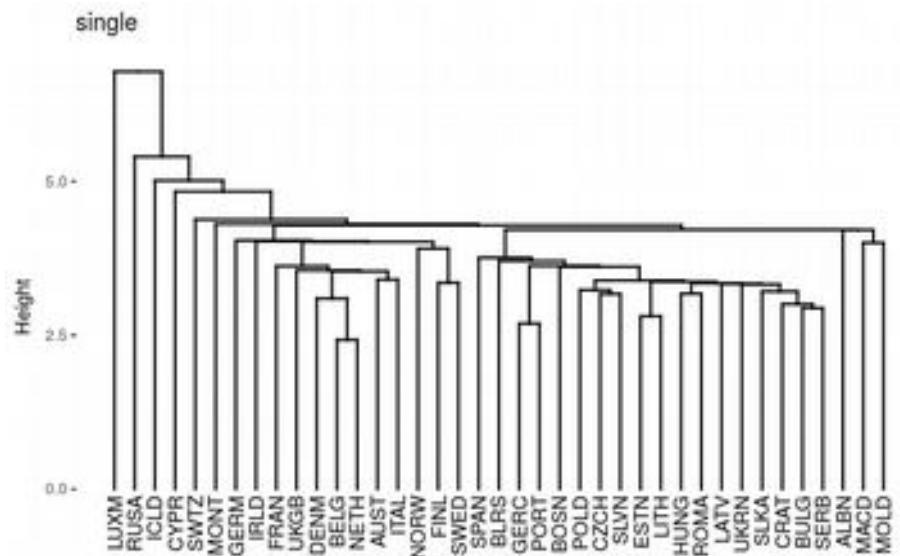
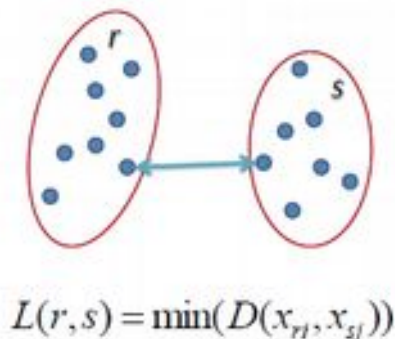
Grupowanie hierarchiczne

- Metoda analizy skupień, której celem jest zbudowanie hierarchii grup
- Stosuje metody aglomeracyjne lub dzielące (rzadziej)
- Algorytm zachłanny, szybki ale niepotymalny globalnie
- Nadaje się do małych zbiorów danych, gdzie struktura (hierarchia) jest ważniejsza niż same skupienia
- Strategie łączenia: pojedyncze, całkowite, średnie



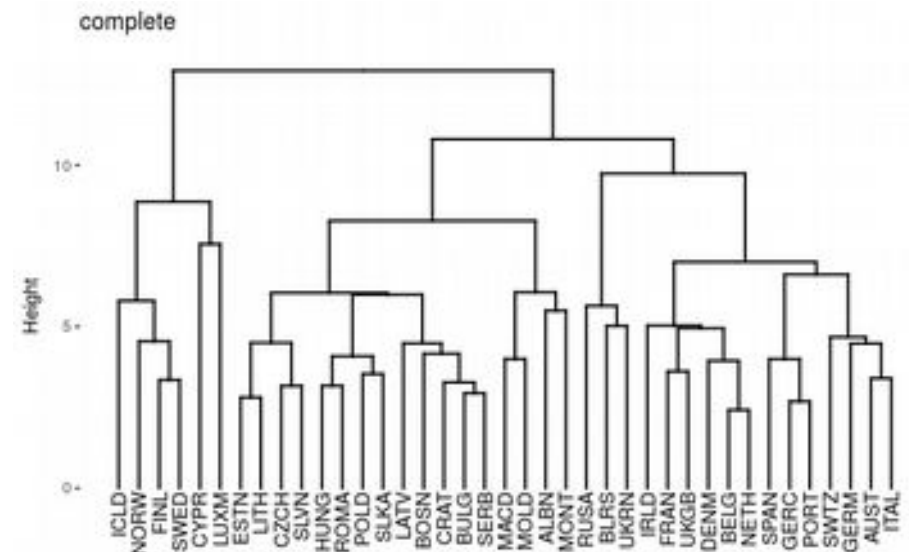
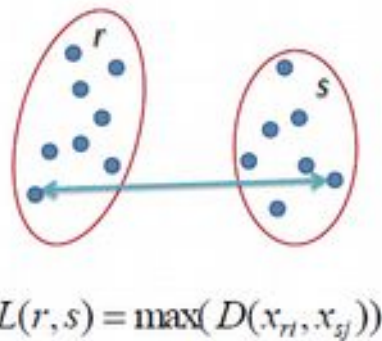
Pojedyncze łączenie

- Jako pierwsze zostaną połączone dwa obiekty o najmniejszym niepodobieństwie
- W następnych krokach będą łączone te obiekty lub grupy, gdzie niepodobieństwo pomiędzy dwoma **najbardziej podobnymi** obiektami **jest najmniejsze**
- Metoda wykrywania obiektów odstających



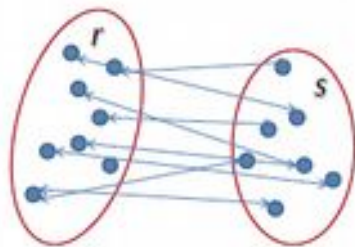
Całkowite łączenie

- Jako pierwsze zostaną połączone dwa obiekty o najmniejszym niepodobieństwie
- W następnych krokach będą łączone te obiekty lub grupy, gdzie niepodobieństwo pomiędzy dwoma **najmniej podobnymi** obiektami **jest najmniejsze**
- Klasyczna metoda budowania hierarchii, nie wykrywa obiektów odstających

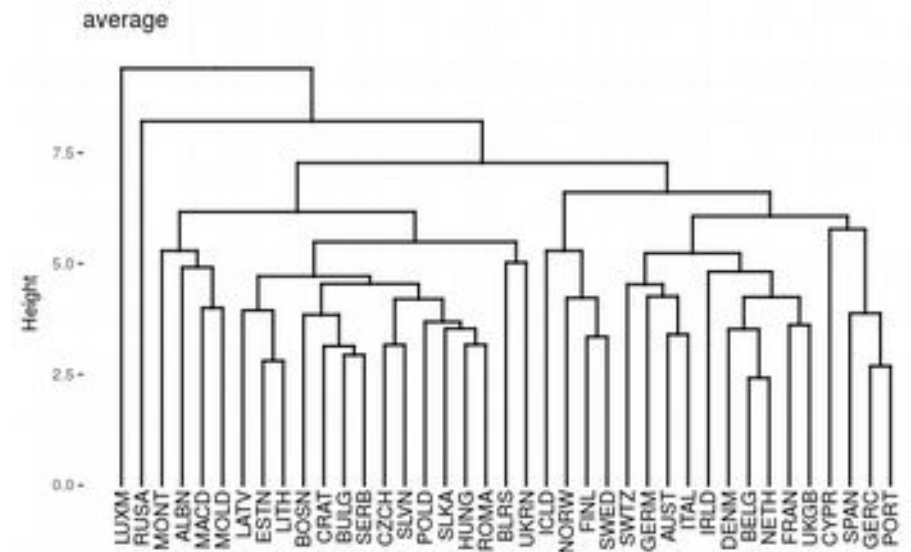


Średnie Łączenie

- Jako pierwsze zostaną połączone dwa obiekty o najmniejszym niepodobieństwie
- W następnych krokach będą łączone te obiekty lub grupy, gdzie **średnie** niepodobieństwo pomiędzy obiektami **jest najmniejsze**
- Klasyczna metoda budowania hierarchii, z wykrywaniem obiektów odstających

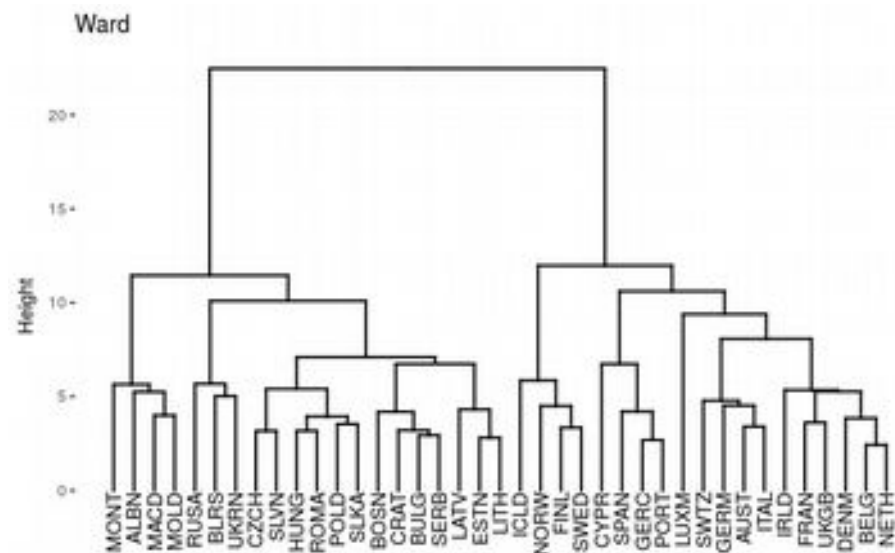


$$L(r, s) = \frac{1}{n_r n_s} \sum_{i=1}^{n_r} \sum_{j=1}^{n_s} D(x_{ri}, x_{sj})$$



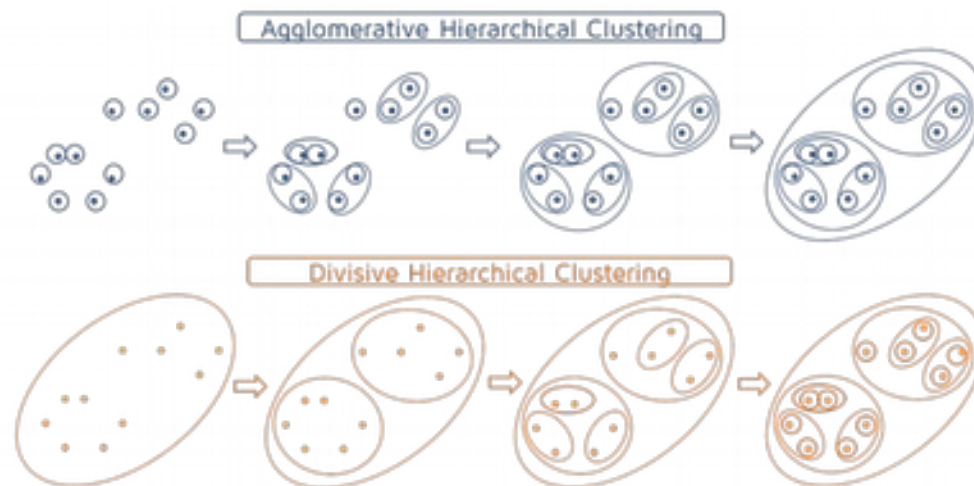
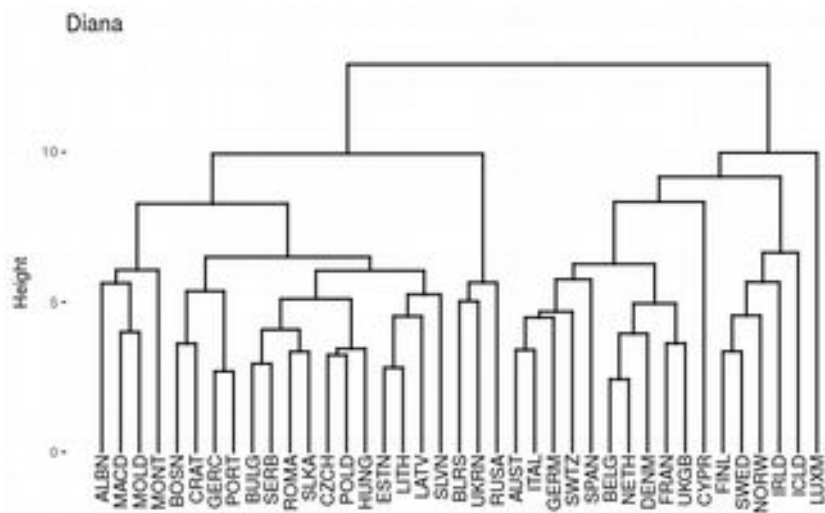
Grupowanie hierarchiczne Warda

- Minimalizacja wariacji wewnątrz skupień, maksymalizacja wariacji między skupieniami
- Nie buduje rzeczywistych hierarchii ale pozwala określić naturalną liczbę skupień
- Nie wykrywa obiektów odstających



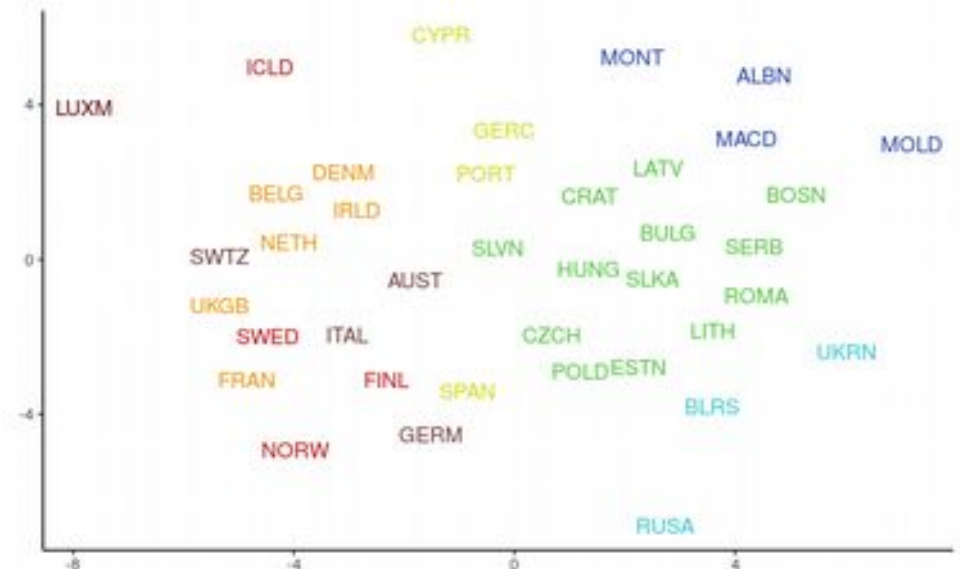
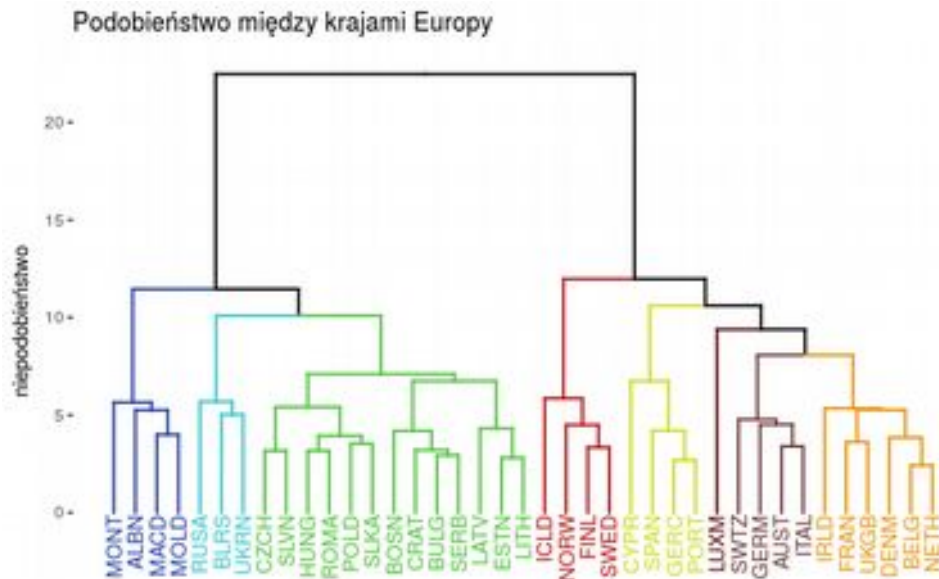
Metoda hierarchicznego rozdzielania

- Rozpoczyna od jednego skupienia obejmującego wszystkie obiekty
- Rozdziela skupienie tak aby maksymalizować wariancję między nimi
- Kontynuuje proces aż do końca
- Nie wykrywa obiektów odstających



Zalety i wady metod hierarchicznych

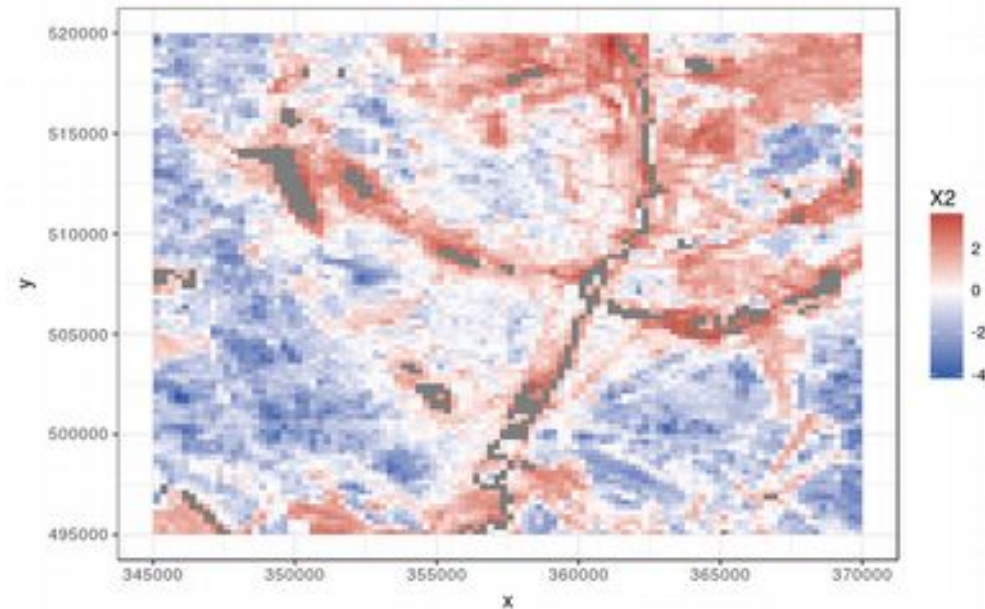
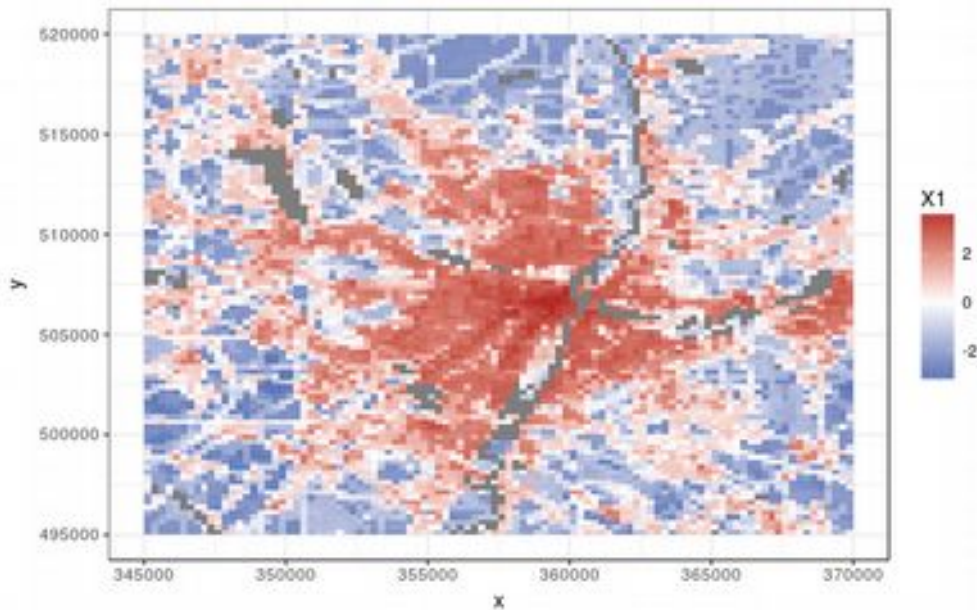
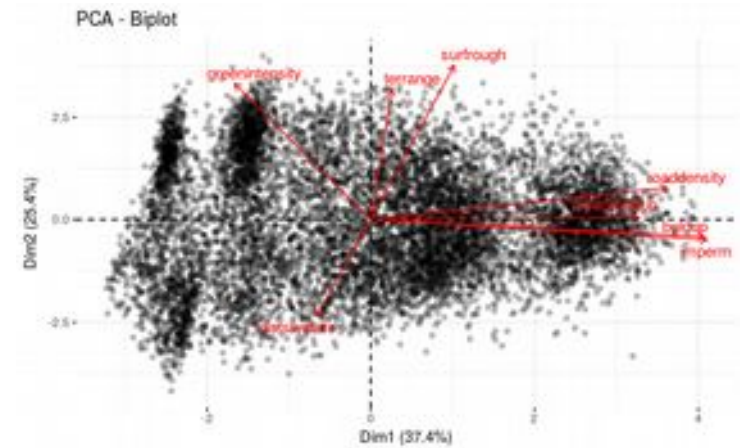
- Zalety
 - Szybki algorytm
 - Deterministyczny algorytm (powtarzalność wyników)
 - Buduje intuicyjnie zrozumiałą hierarchię
- Wady
 - Algorytm zachłanny, optymalizowany na poziomie kroku a nie całości wyniku
 - Każda decyzja nie może być zmieniona
 - Skupienia rozmieszczone są liniowo, tracimy informację o relacjach pomiędzy skupieniami



Metody partycjonujące

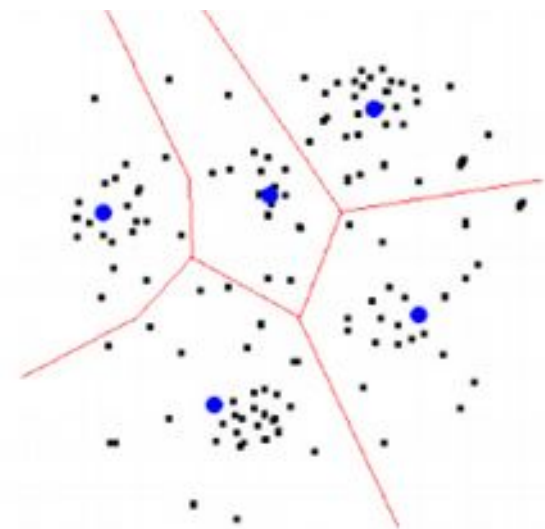
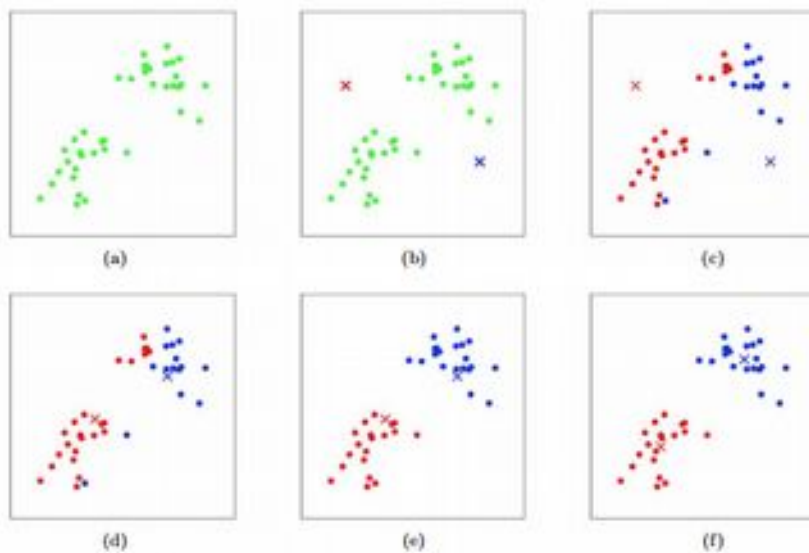
- K-średnie
- K-medoidy
- Propagacja afiniczności
- Rozmyte k-średnie

Jako przykład zostaną użyte dane z poprzedniego wykładu



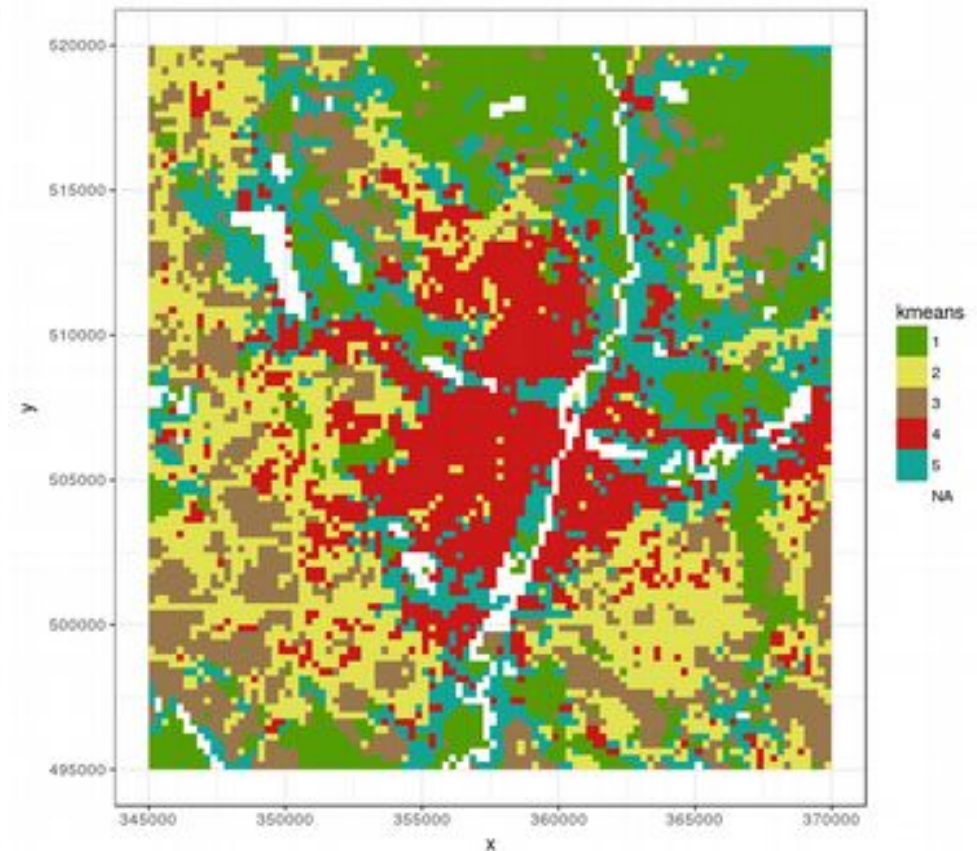
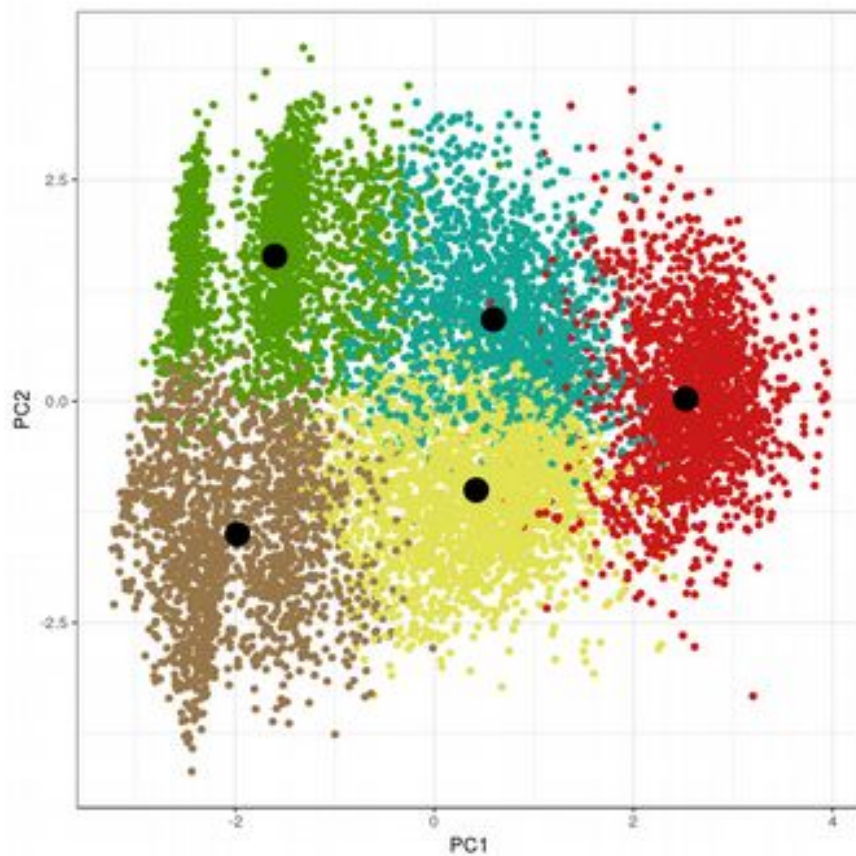
Metoda k-średnich

- 1) Algorytm stochastyczny rozpoczyna losowo położonymi punktami (centroidami)
 - 2) Przypisuje obiekty do centroidów na zasadzie minimalnego niepodobieństwa
 - 3) Wyznacza nową lokalizację na podstawie zasięgu skupienia
 - 4) Powtarza (2) aż do momentu gdy położenie centroidów nie zmieni się
- Ze względu na duży wpływ początkowej konfiguracji algorytm rozpoczyna proces wielokrotnie, wybierając najbardziej powtarzalne wyniki



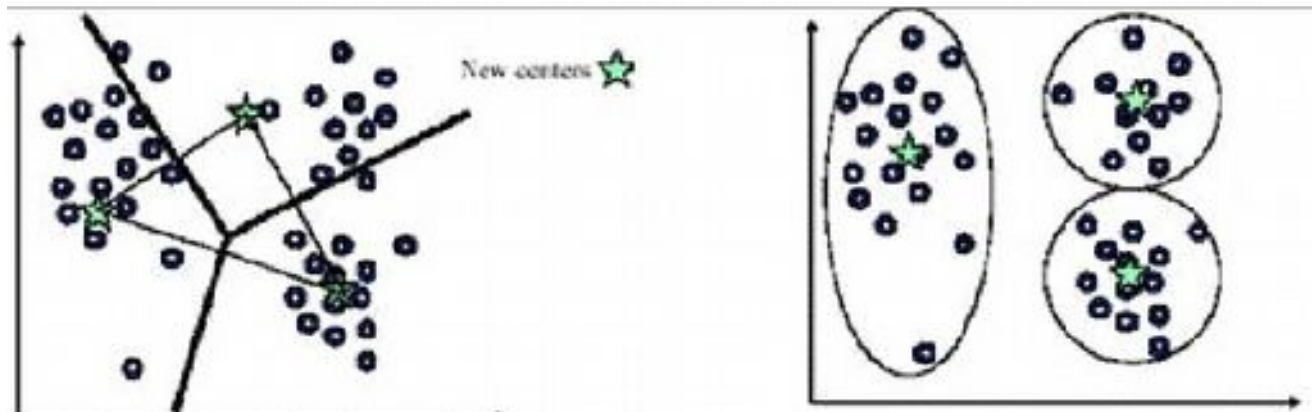
Zastosowanie dla danych geoprzestrzennych

- Metoda wymaga podania liczby skupień
- Wynik podziału jest zgodny z kryterium Voronoi
- Zagęszczenia w rozkładzie nie mają znaczenia dla procesu wyznaczania skupień
- Niepewność przynależności nie jest brana pod uwagę



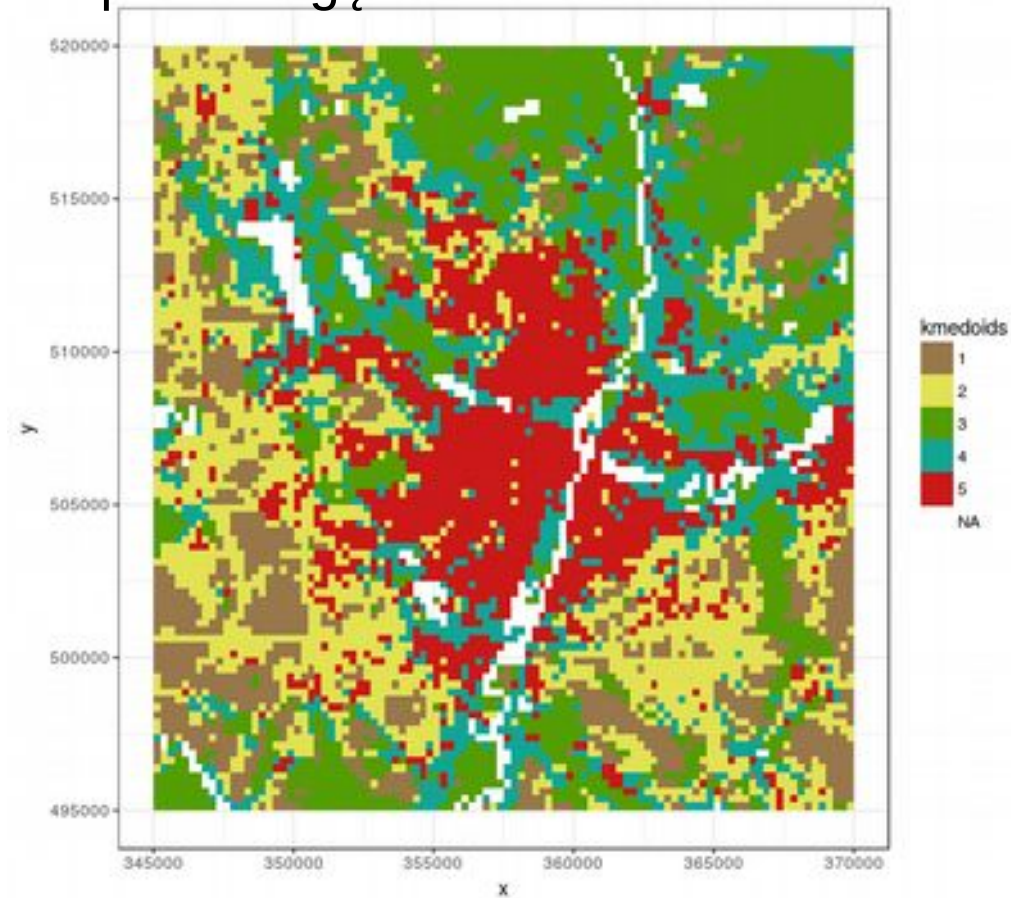
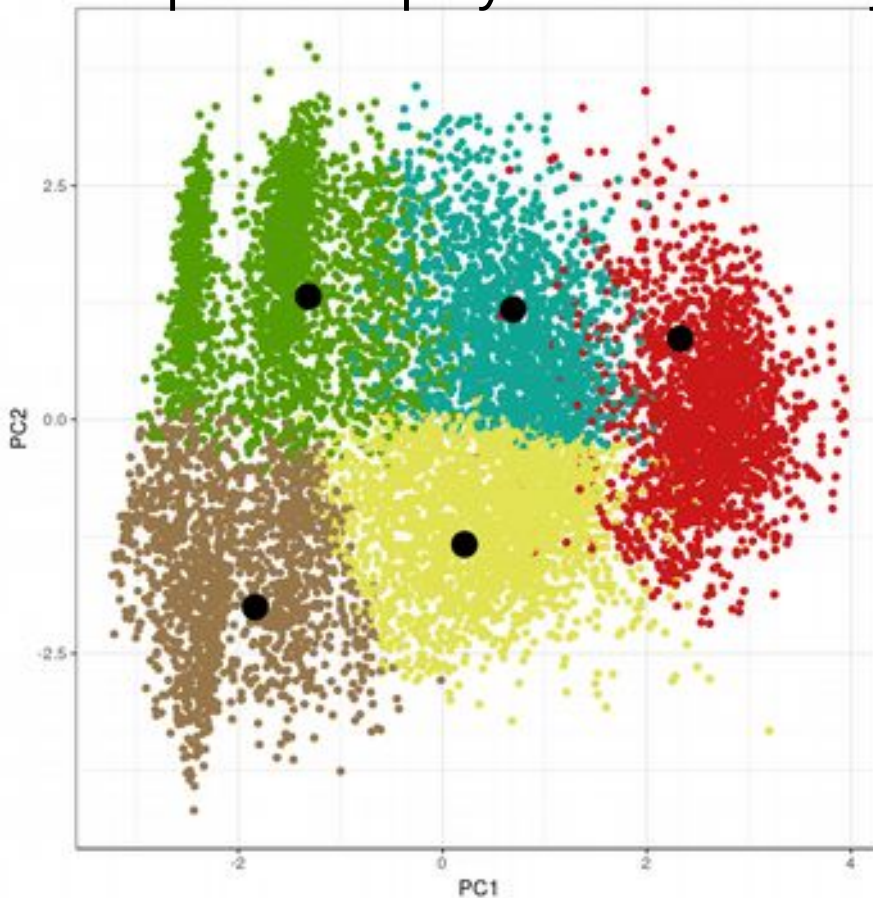
Metoda k-medoidów

- Algorytm podobny do k-średnich. Nie używa abstrakcyjnych centroidów ale rzeczywiste obiekty ze zbioru (medoidy)
- 1) Algorytm stochastyczny rozpoczyna losowo wybranymi obiektami (medoidami)
- 2) Pozostałe kroki jak w k-means
- Ze względu na duży wpływ początkowej konfiguracji algorytm rozpoczyna proces wielokrotnie, wybierając najbardziej powtarzalne wyniki
- W przeciwieństwie do kmeans dużo bardziej odporny na obiekty odstające – jeżeli występują



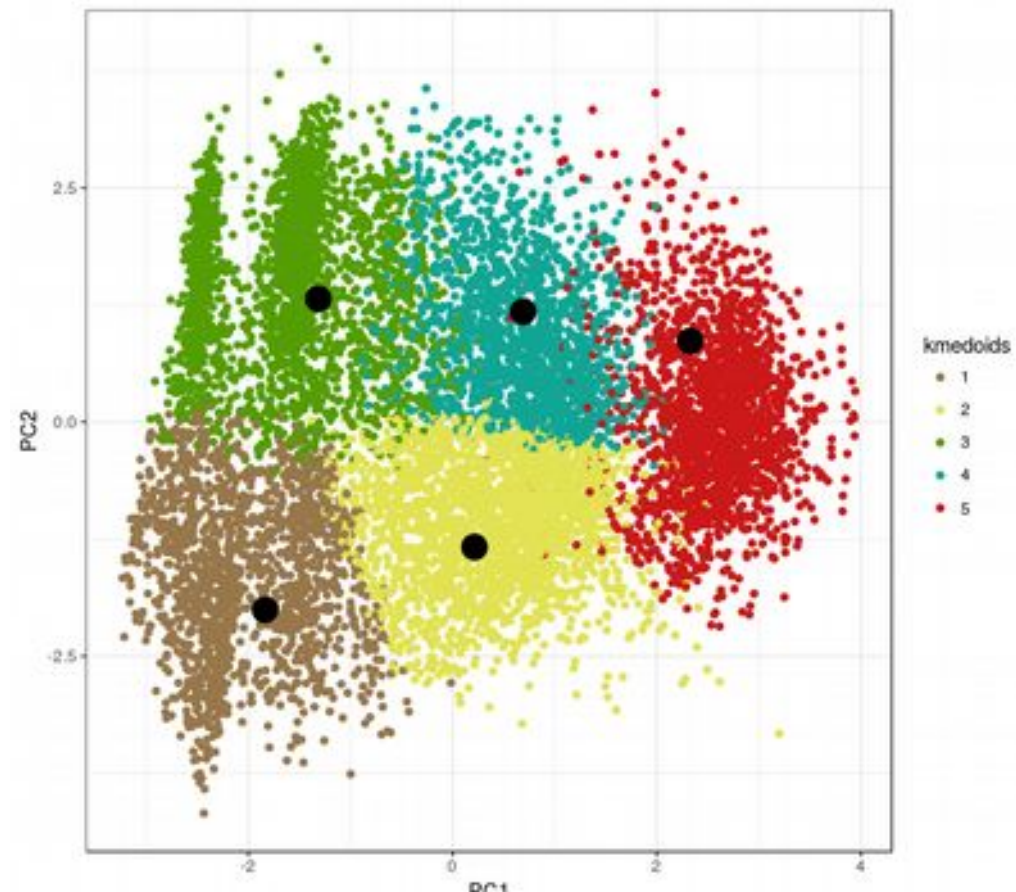
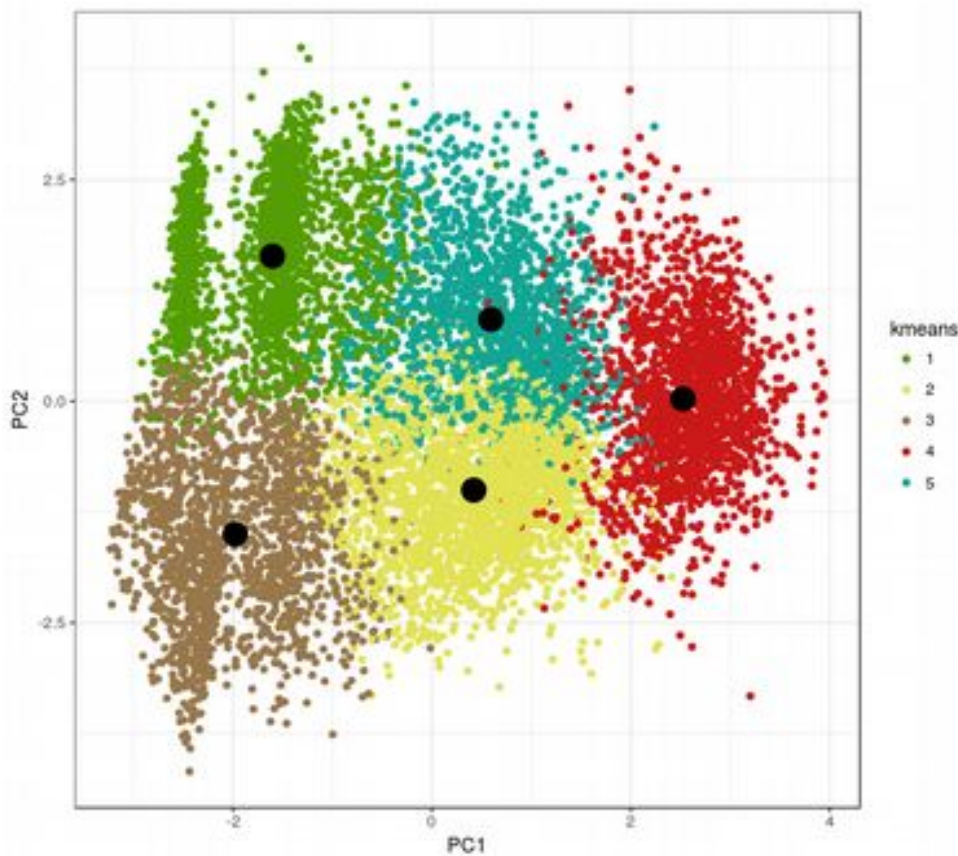
Zastosowanie dla danych geoprzestrzennych

- Metoda wymaga podania liczby skupień
- Zagęszczenia w rozkładzie mają znaczenie dla wyznaczania skupień
- Mały wpływ obiektów odstających
- Niepewność przynależności nie jest brana pod uwagę



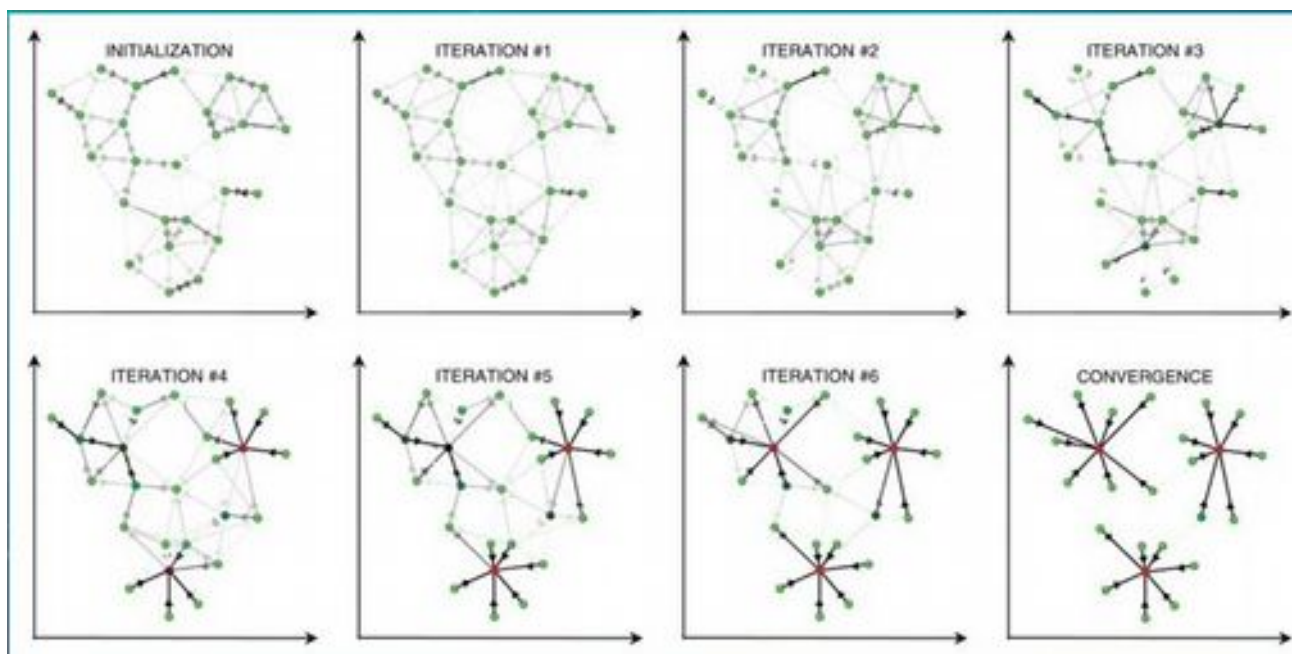
Porównanie skupień

- Wyniki bardzo podobne
- K – medoids daje wyraźniejsze skupienia
- Obie metody zaliczane są do suboptymalnych tj wynik nie jest najlepszy z możliwych ale akceptowalny



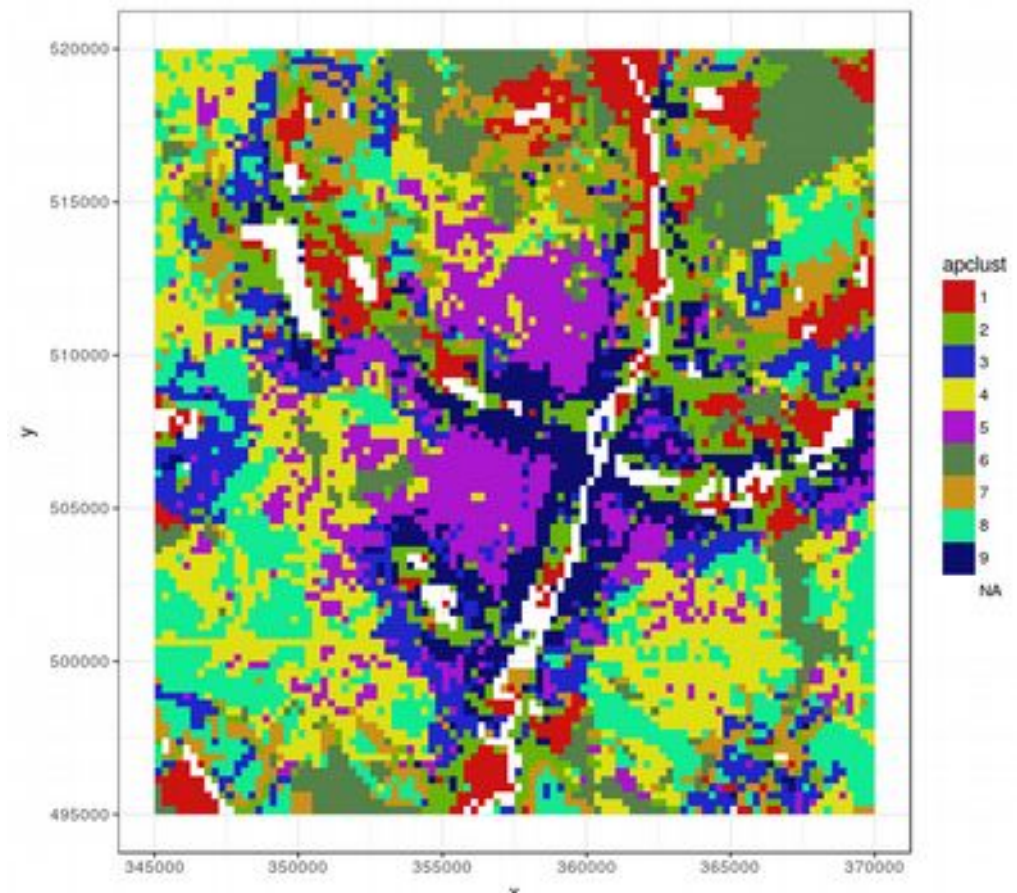
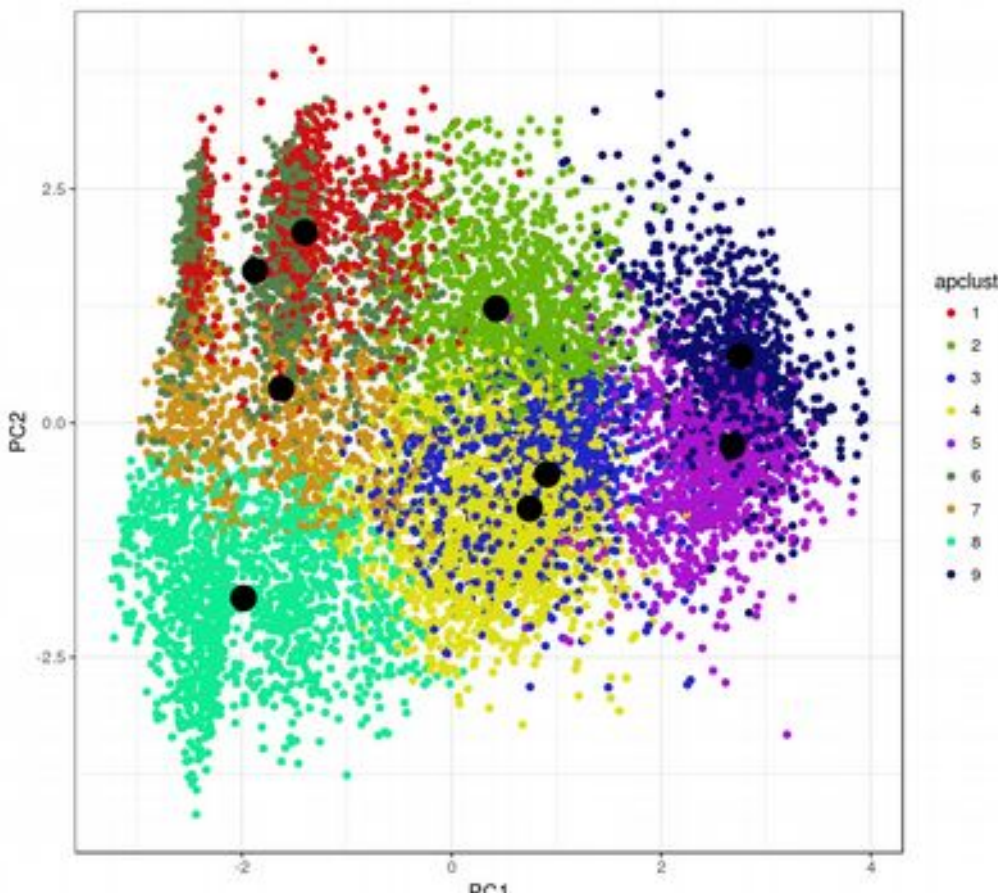
Propagacja powinowatości

- Affinity propagation – metoda polegająca na iteracyjnym „przekazywaniu wiadomości” pomiędzy obiektami. Ma na celu wybór obiektów – egzemplarów
- Jako egzemplary wskazywane są te obiekty, które wykazują dodatni bilans pomiędzy byciem egzemplarem, a posiadaniem egzemplara
- Metoda wyszukiwania naturalnych liderów. Bardziej „pasuje do wszystkich” niż „mistrz w jednym”



Zastosowanie dla danych geoprzestrzennych

- Metoda nie wymaga podania liczby skupień jedynie kryterium selekcji egzemplarów
- Zagęszczenia w rozkładzie mają znaczenie dla wyznaczania skupień
- Obiekty odstające tworzą nowe skupienia
- Niepewność przynależności nie jest brana pod uwagę

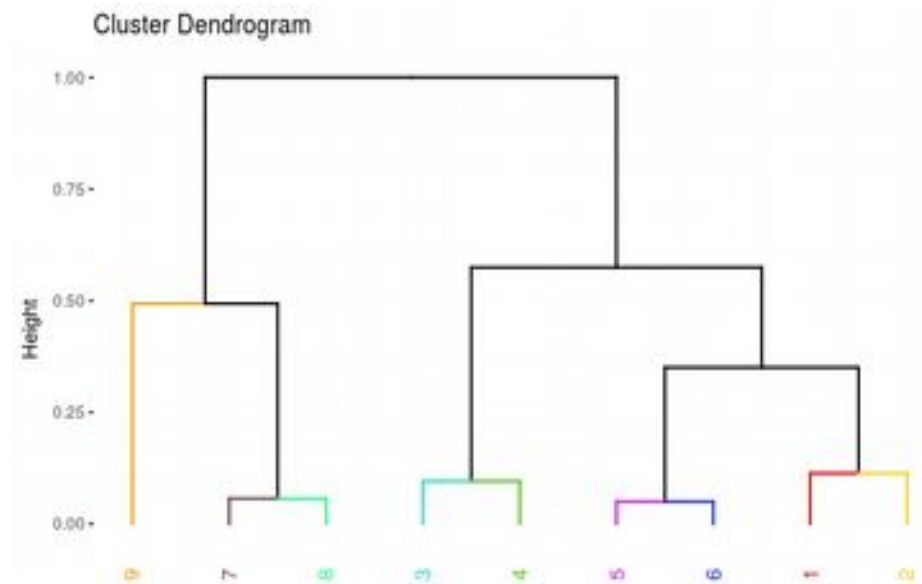
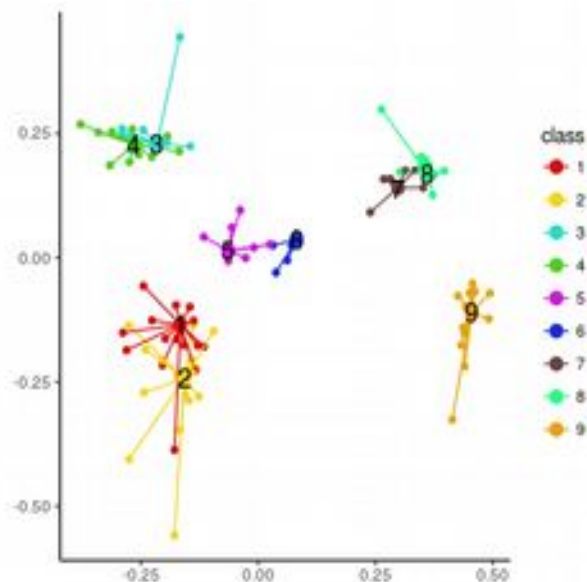


Centroid – Medoid - Egzemplar

- **Centroid:** współrzędne w przestrzeni wielowymiarowej oznaczające geometryczny środek skupienia. Nie jest to fizyczny obiekt. Może być poza obszarem skupienia
- **Medoid:** obiekt najbardziej podobny do innych obiektów. Z reguły występuje w największym zagęszczeniu skupienia.
- **Egzemplar:** naturalny przedstawiciel skupienia, najbardziej reprezentatywny dla innych obiektów

Łączenie metod partycjonujące i hierarchicznych

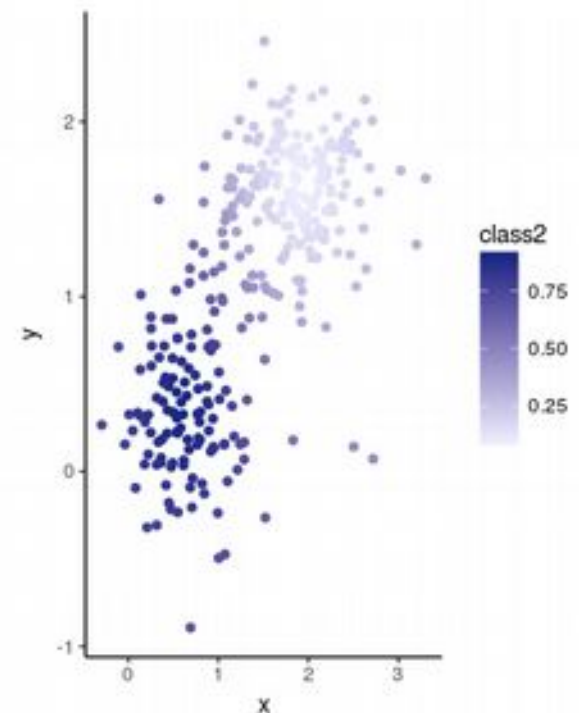
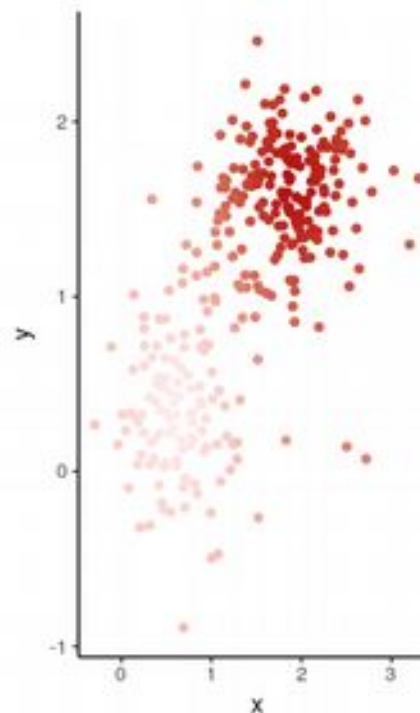
- Metody partycjonujące i hierarchiczne można łączyć, ale nie wszystkie implementacje używają tej możliwości
- Metody hierarchiczne jako szybsze, ale mniej dokładne używa się do wstępnego podziału zbioru na skupienia, przed uruchomieniem metody k-means/k-medoids w celu uniknięcia losowej konfiguracji startowej
- W przypadku dużych zbiorów danych metody partycjonujące używa się do wyznaczenia dużej liczby małych zwartych skupień, a następnie małe skupienia łączy się w hierarchie





Metody rozmyte (k- means i k-medoids)

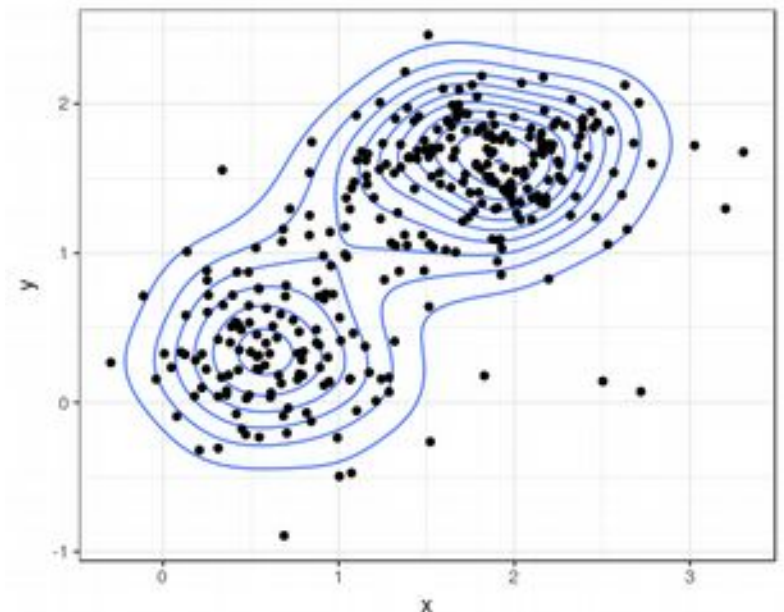
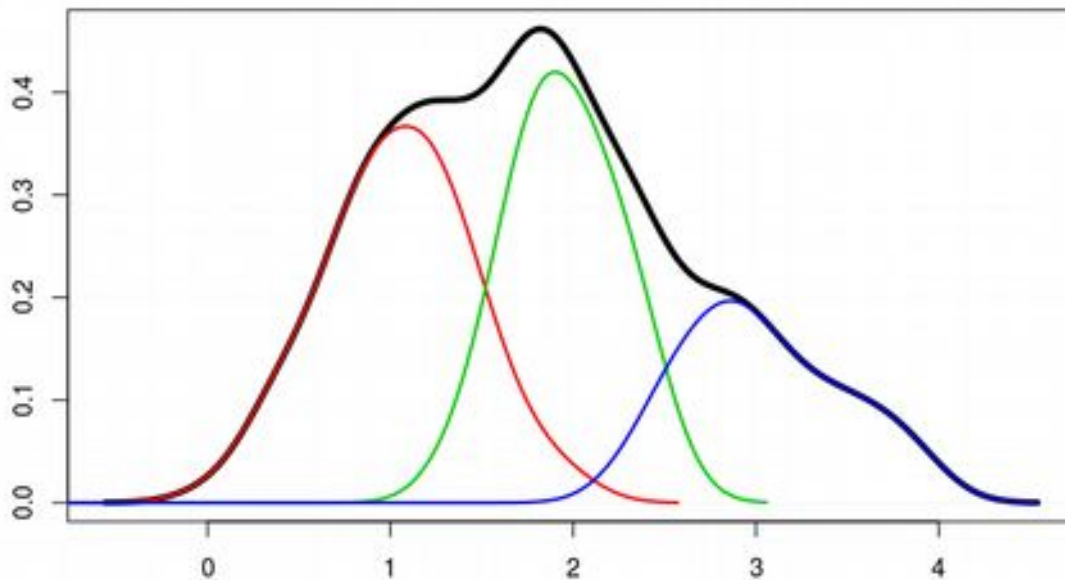
- Metoda bierze pod uwagę niepewność przynależności
- Każdy obiekt jest przypisywany do więcej niż jednej z klas
- zastosowaniem metod rozmytych jest sytuacja, gdy interesują nas jedynie wybrane skupienia i chcemy określić dla nich tolerancję przynależności kosztem innych skupień.
- Koncepcja krytykowana:
 - w ostateczności obiekt musi przynależć to jakiegoś skupienia,
 - do określenia niepewności przynależności służą inne metody
 - metody rozmyte mają problemy z prawidłowym klasyfikowaniem punktów na obrzeżach.
 - problemy z wizualizacją skupień, gdyż wymagają osobnego diagramu dla każdego skupienia.



Metody probabilistyczne

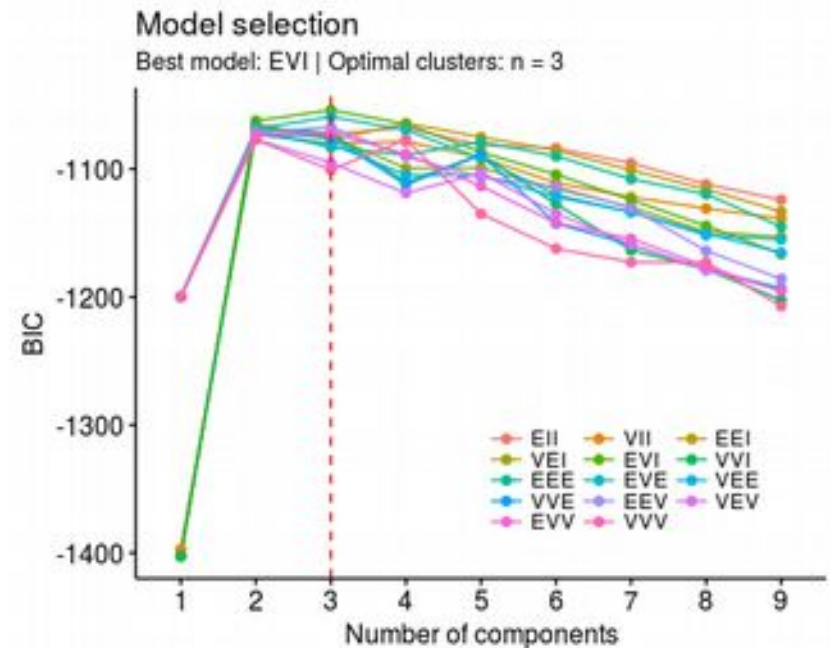
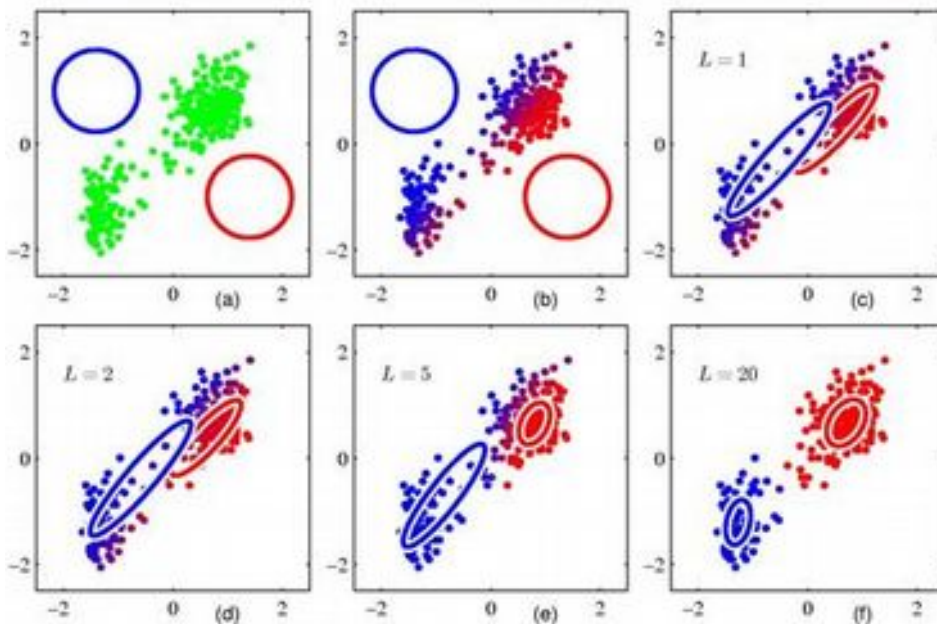
- Gaussowskie modele mieszane
- Jeżeli rozkład gęstości nie ma jakiejś konkretnej postaci można przyjąć założenie że jest sumą wielu rozkładów normalnych
- Znalezienie rozwiązania jest problemem optymalizacyjnym, szuka się optymalnej liczny skupeń oraz właściwych dla nich rozkładów

Rozkład gęstościowy



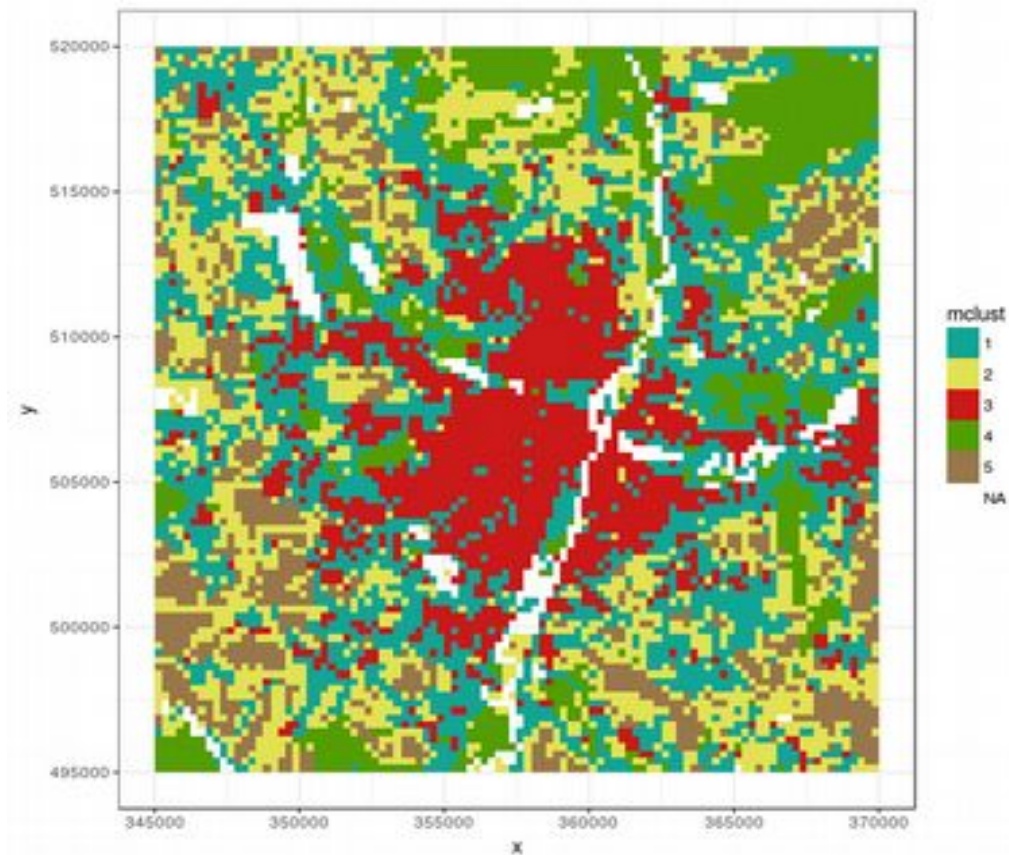
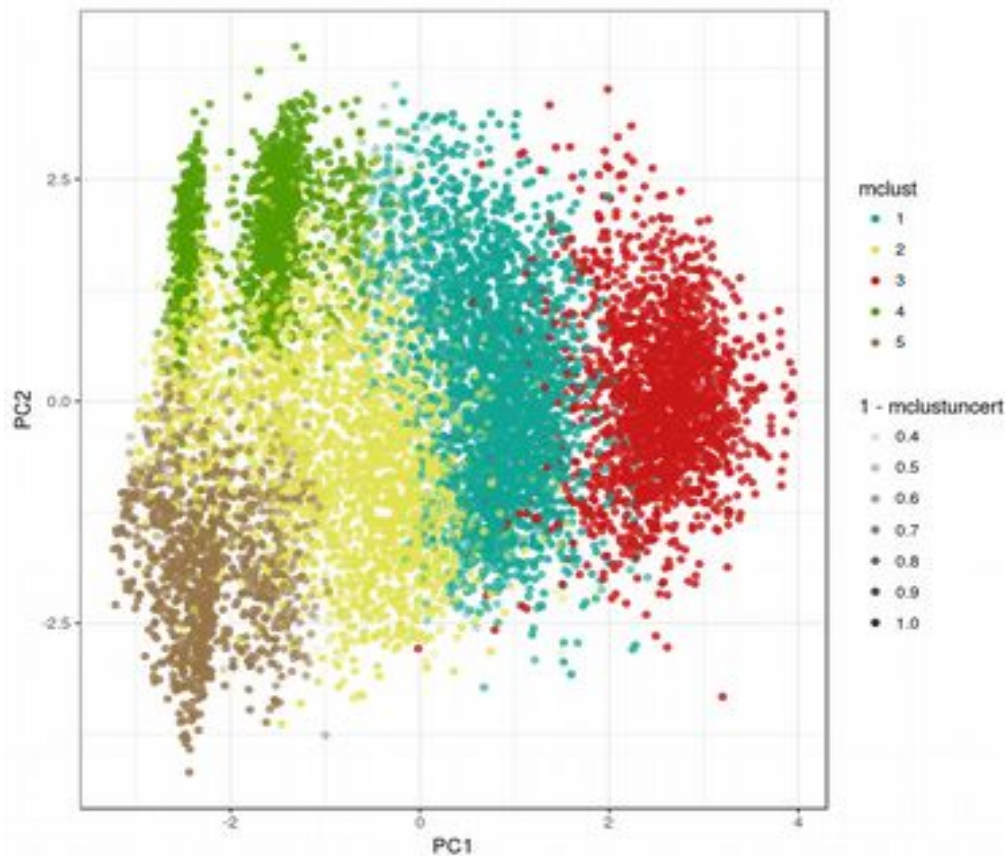
Expectation-maximisation

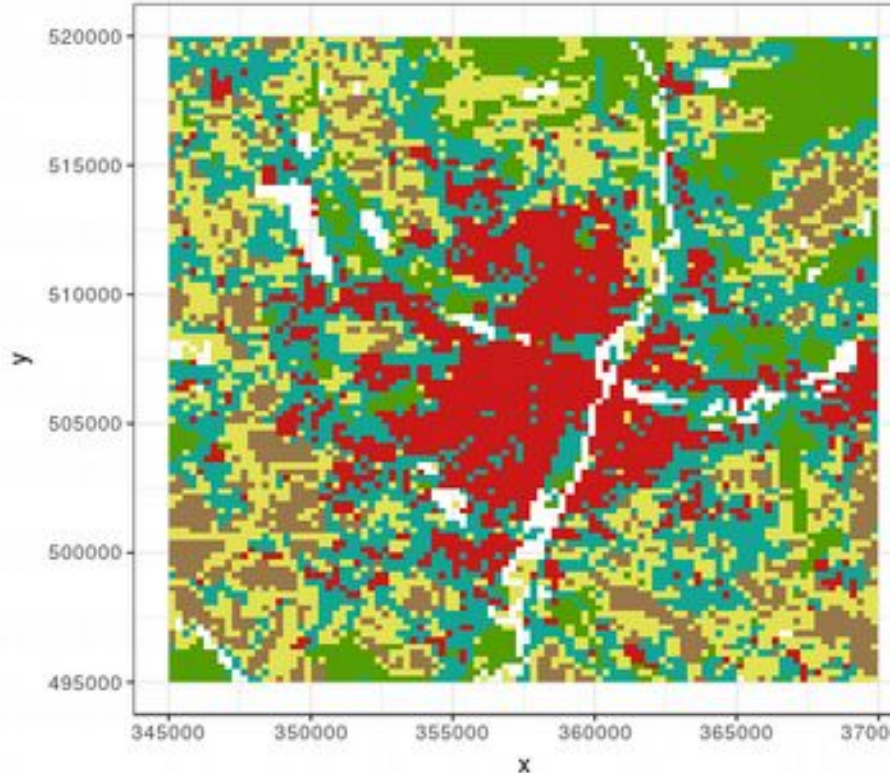
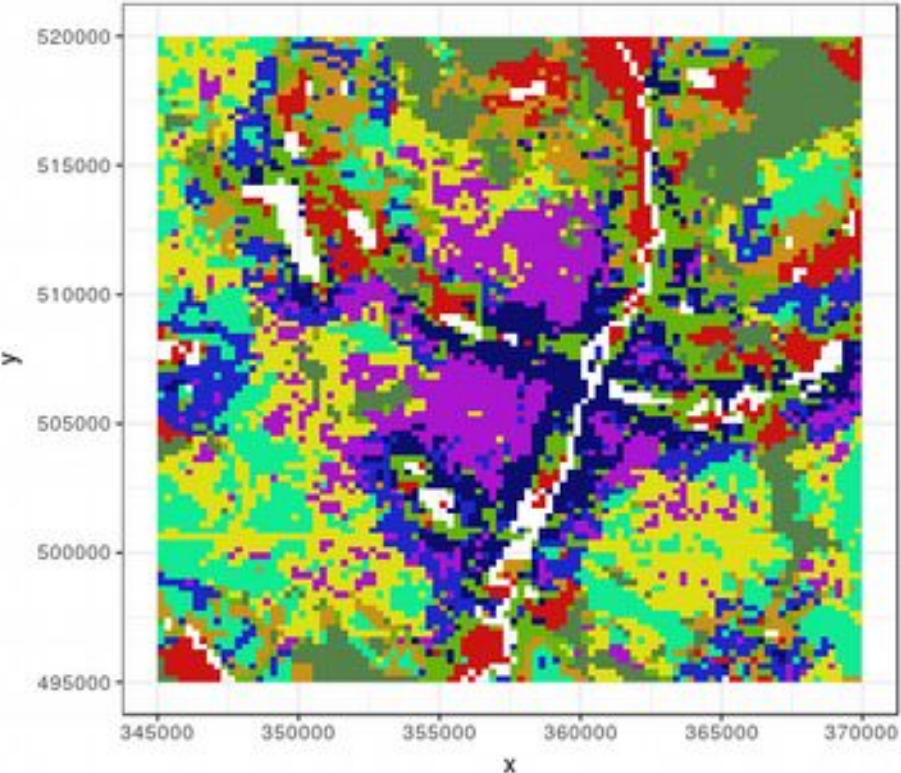
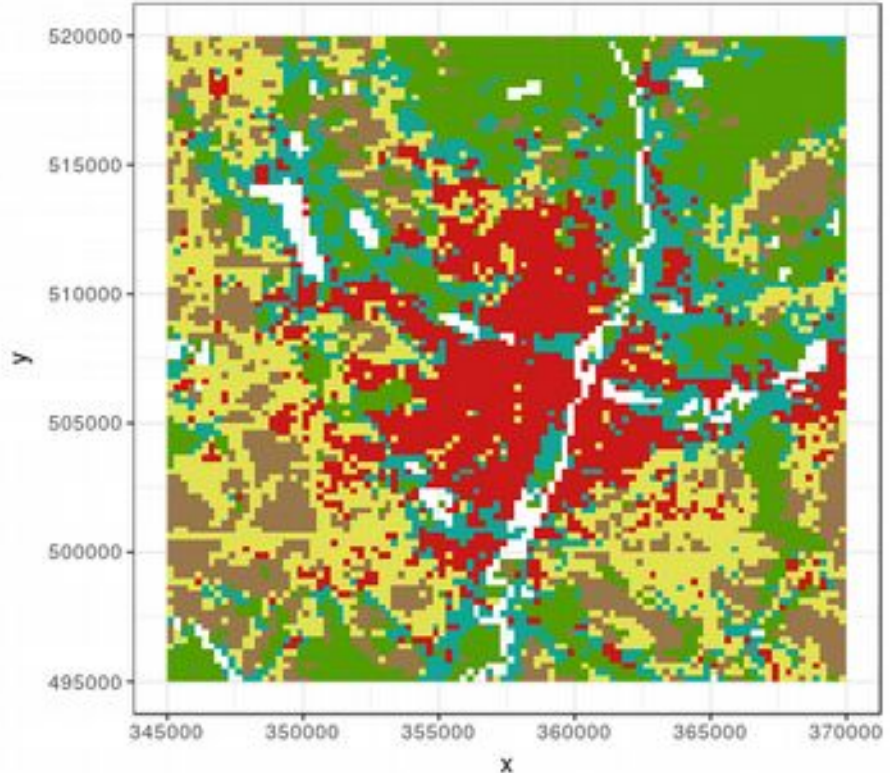
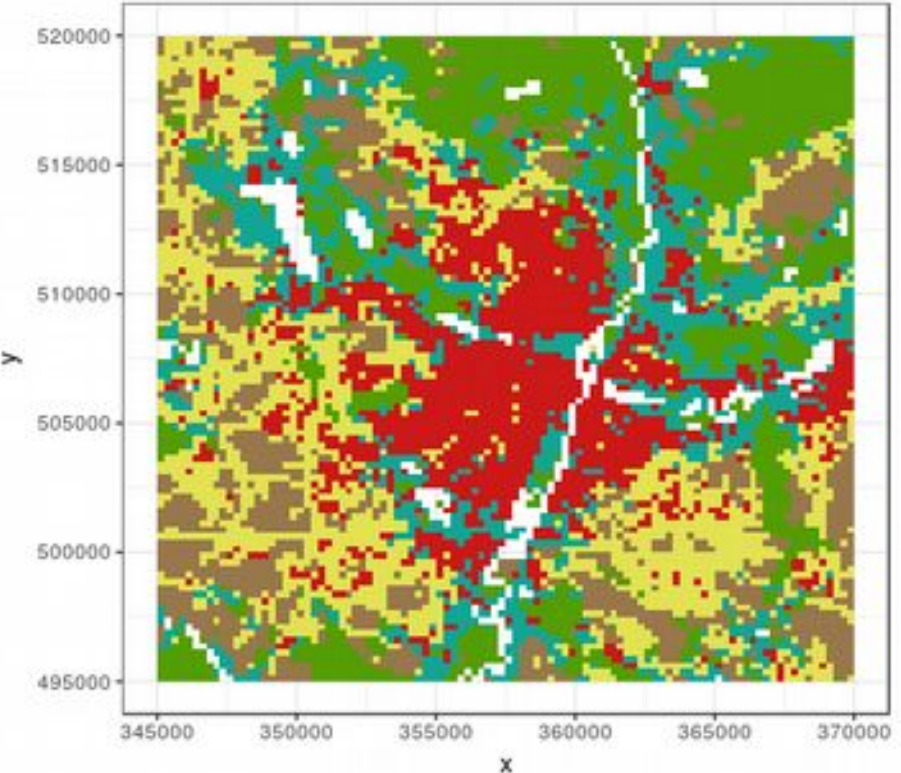
- Jest to proces iteracyjny, składający się z dwóch kroków: E (*expectation*) czyli znalezienia najlepszego rozkładu oraz M (*maximisation*) polegającego na uaktualnieniu parametrów modelu poprzez maksymalizację funkcji wiarygodności (*likelihood*)
- Wybór rozwiązania opiera się na minimalizacji parametru BIC (Bayes inf. criterion), które powinno być najmniejsze



Zastosowanie dla danych geoprzestrzennych

- Zagęszczenia w rozkładzie mają duży wpływ na ostateczny wynik
- Niepewność jest brana pod uwagę
- Znikomy wpływ obiektów odstających
- Bardzo wolny czas obliczeń
- Wynik jest optymalny dla podzbioru danych





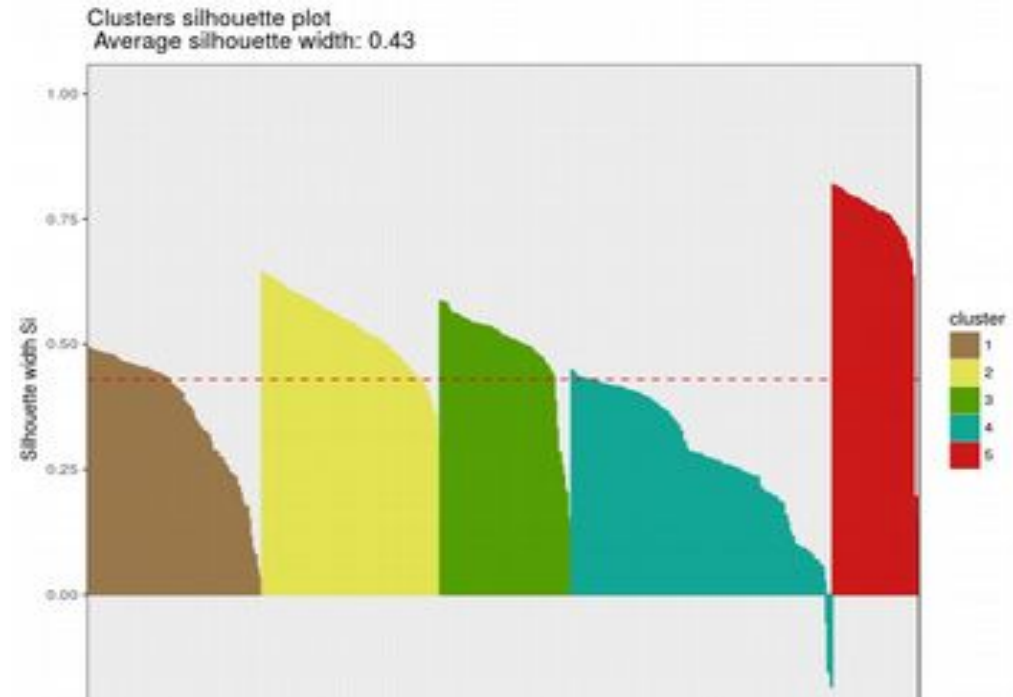
Parametry skupień

- Ocena jakości skupień służy do określenia na ile prawidłowo dobrano liczbę skupień oraz do jakiego stopnia obiekty zostały zakwalifikowane do właściwych skupień
- Najpopularniejsze wskaźniki to:
 - **Zwartość** (*compactness*) - jak podobne względem siebie są obiekty w tym samym skupieniu – wzajemne średnie/maksymalne niepodobieństwo pomiędzy obiektami
 - **Oddzielność** (*separation*) – jak niepodobne są obiekty w różnych skupieniach: wzajemne minimalne/średnie niepodobieństwo obiektów w różnych skupieniach
 - **Łączność** (*connectivity*) – do jakiego stopnia obiekty położone blisko siebie znajdują się w tych samych skupieniach

Diagramy sylwetkowe

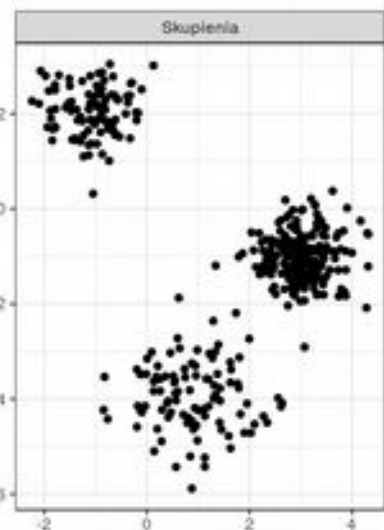
- **Sylwetki** – ocena w jak bardzo obiekty w skupieniu są podobne do pozostałych obiektów w skupieniu względem obiektów w innym (najbardziej podobnym) skupieniu, im większa wartość parametry sylwetki tym lepsze skupienia

- S – bliskie 1; dobre skupienie
- S – bliskie 0; przynależność niejasna
- $S < 0$ błędna przynależność, zmiana przynależności podniesie jakość skupień

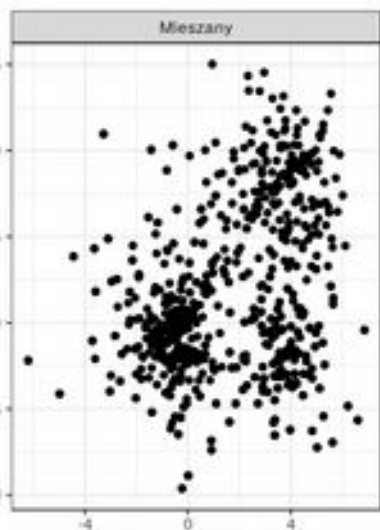


Wybór algorytmu grupowania

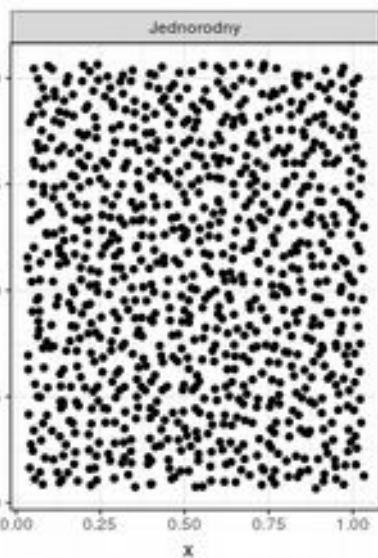
- Brak jednolitej reguły



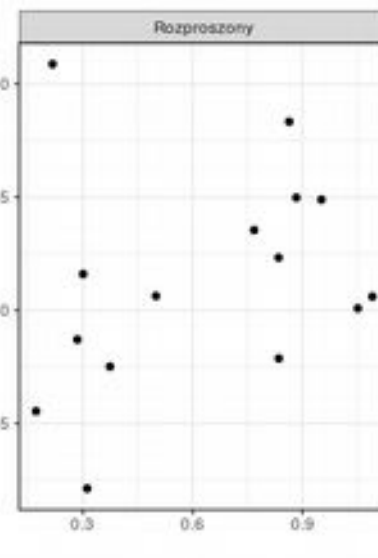
K-medoids
K-means
Ward



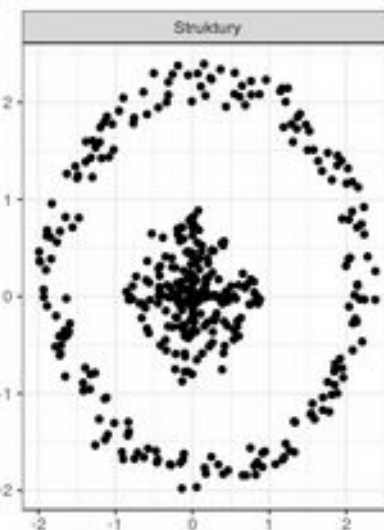
GMM



AP



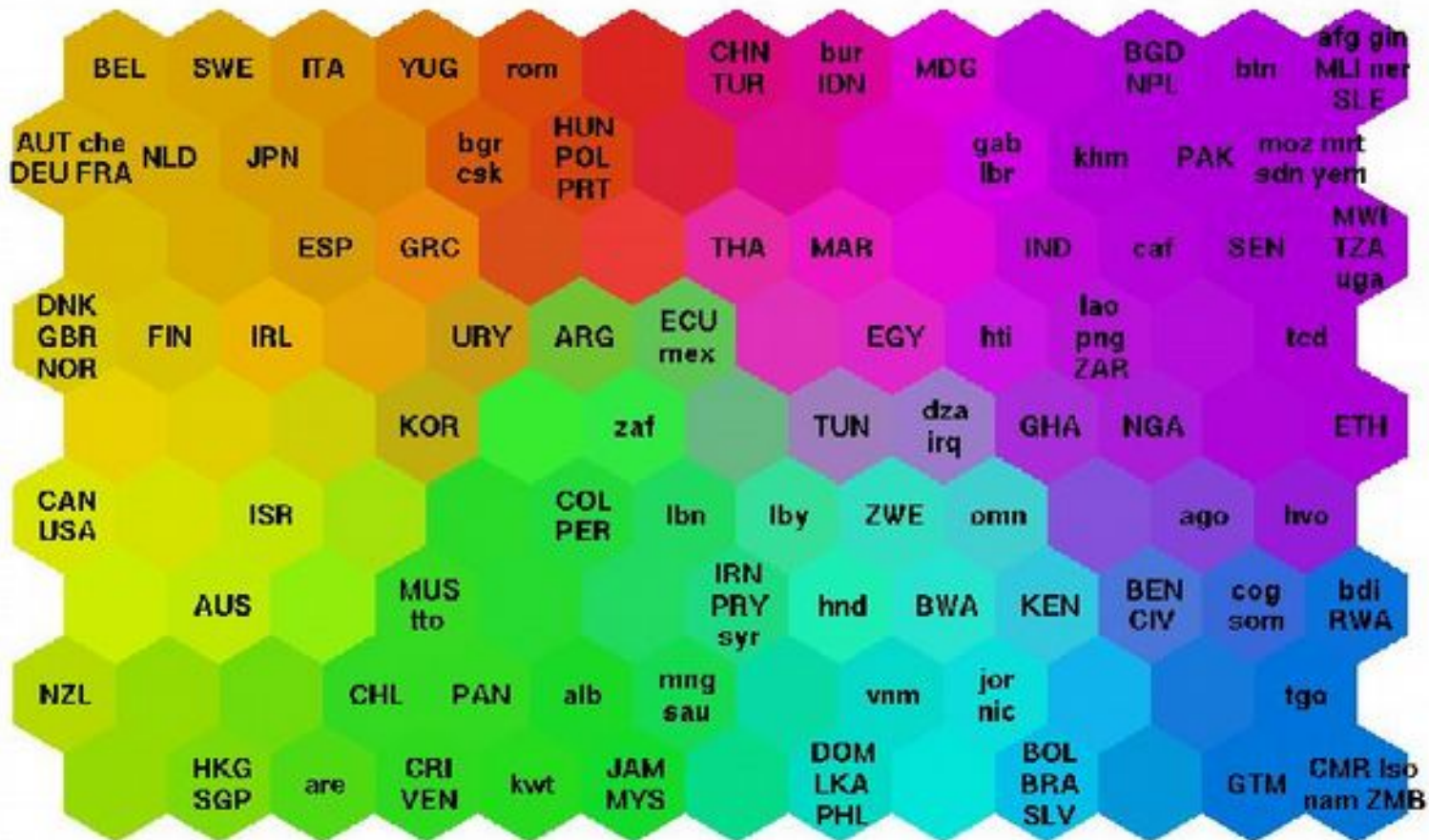
Complete
linkage



dbscan

Samo-organizujące się mapy

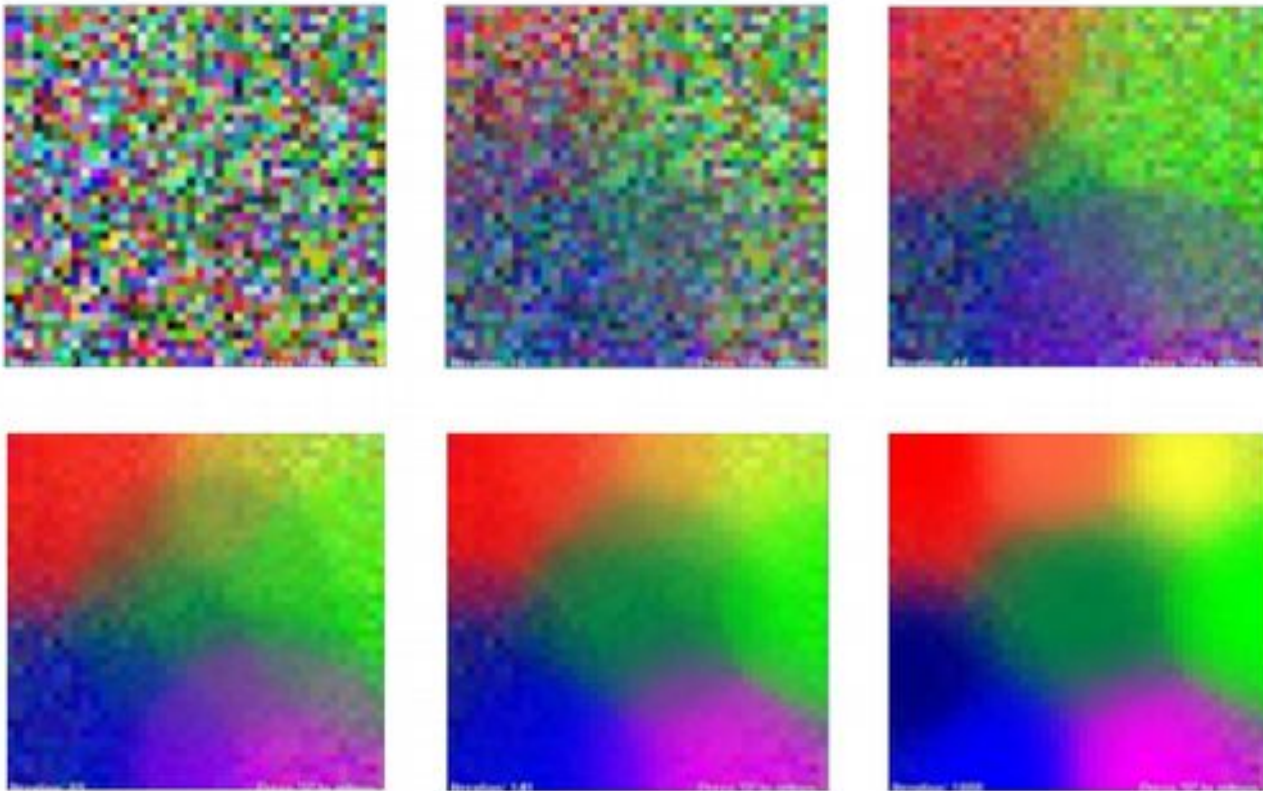
Countries organized on a self-organizing map based on indicators related to poverty:



Koncepcja SOM

- Sieć neuronowa, narzędzie wizualizacji danych wielowymiarowych w postaci mapy topologicznej
- Metoda nienadzorowana, nie wymaga wzorców (w przeciwieństwie do klasycznych sieci neuronowych)
- Polega na obliczaniu odległości pomiędzy wektorami docelowymi a próbkami i przypisywaniu próbek do wektorów docelowych, jednocześnie je zmieniając. Stąd pojęcie „samoorganizujące”

Wektor kodowy



SOM: Components

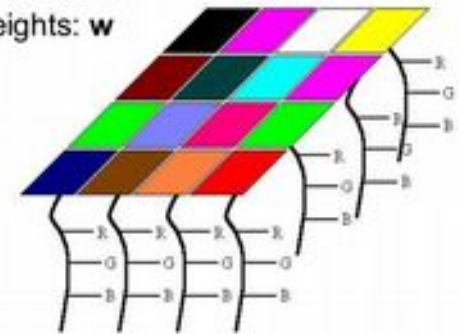
Inputs: x



$X=(R,G,B)$ is a vector!
Of which we have six here.

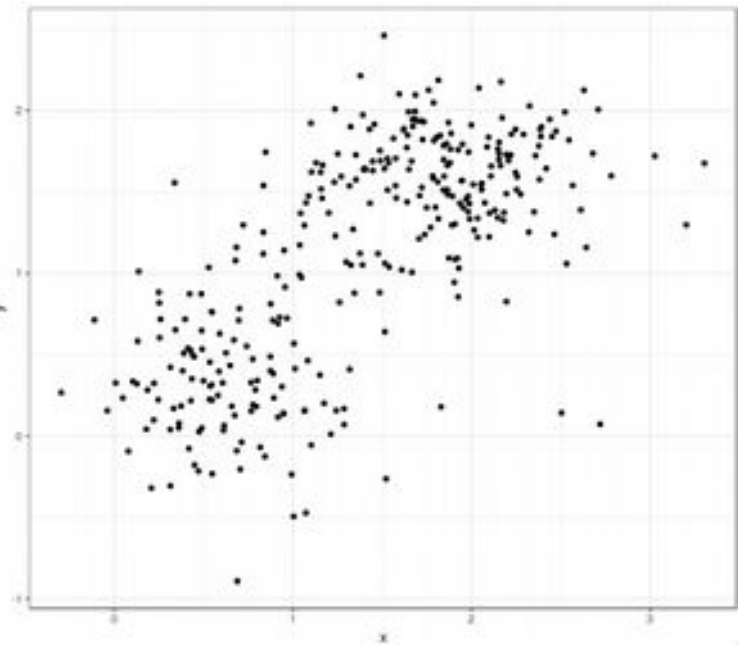
Weights: w

We use 16
codebook vectors
(you can choose
how many!)

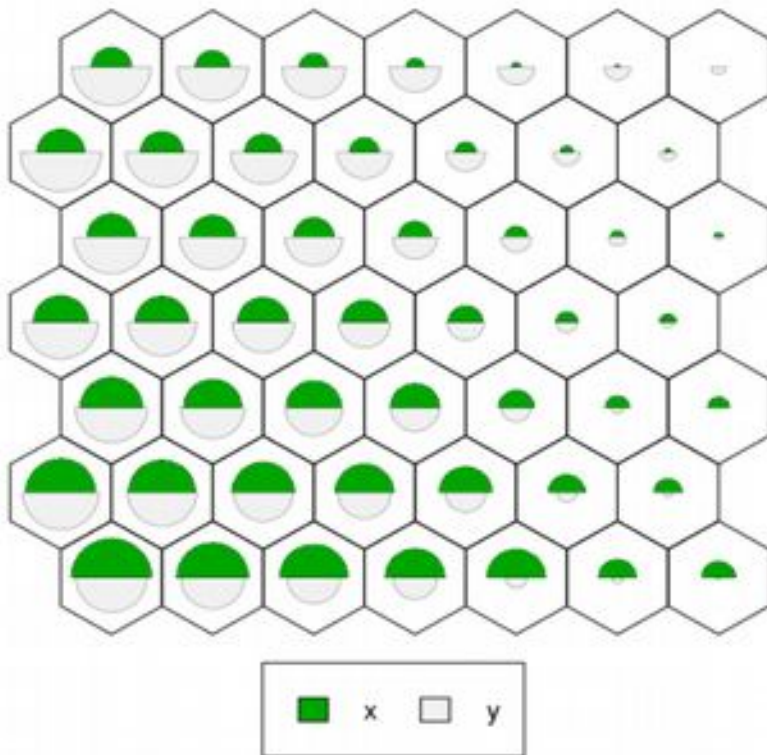


Porządkowanie kolorów RGB

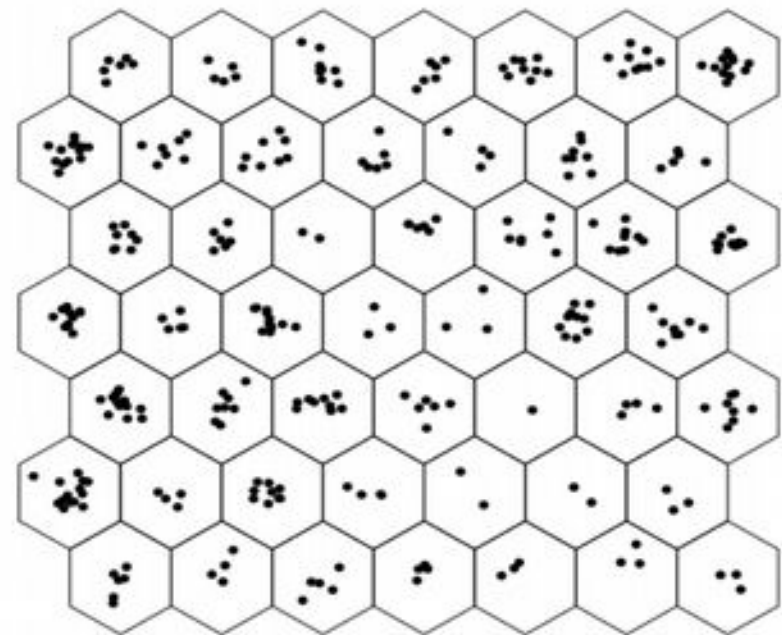
SOM jako narzędzie redukcji wymiarowości



Codes plot

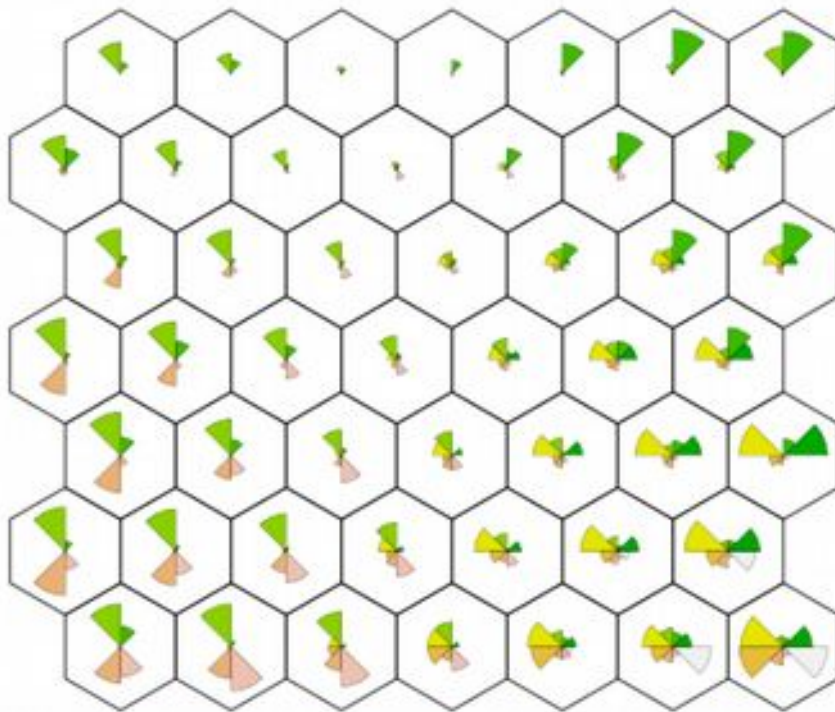


Mapping plot

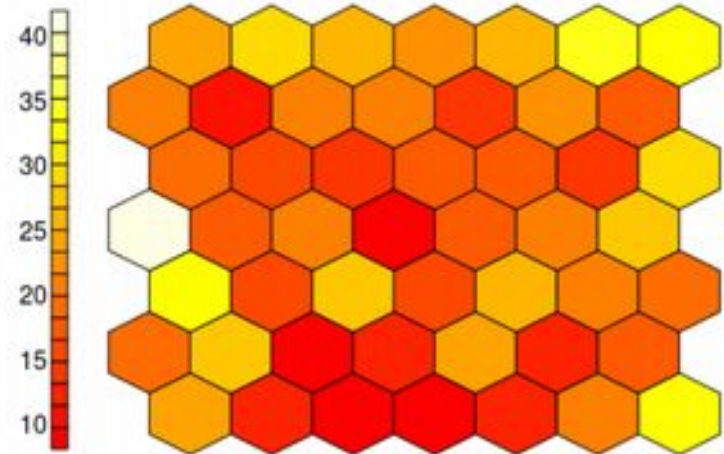


Dane miejskie

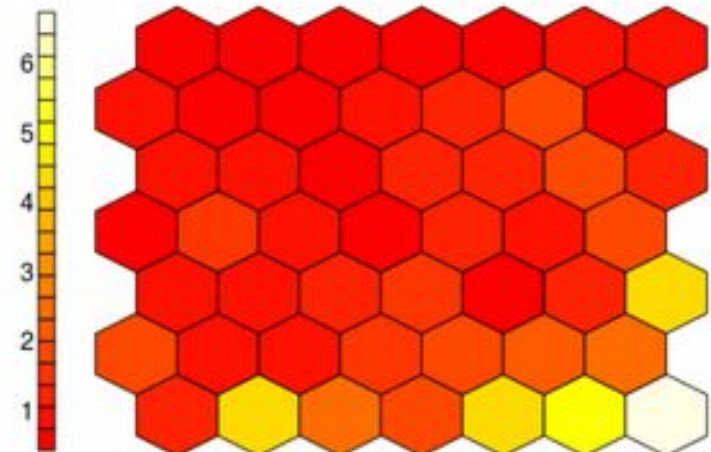
Codes plot



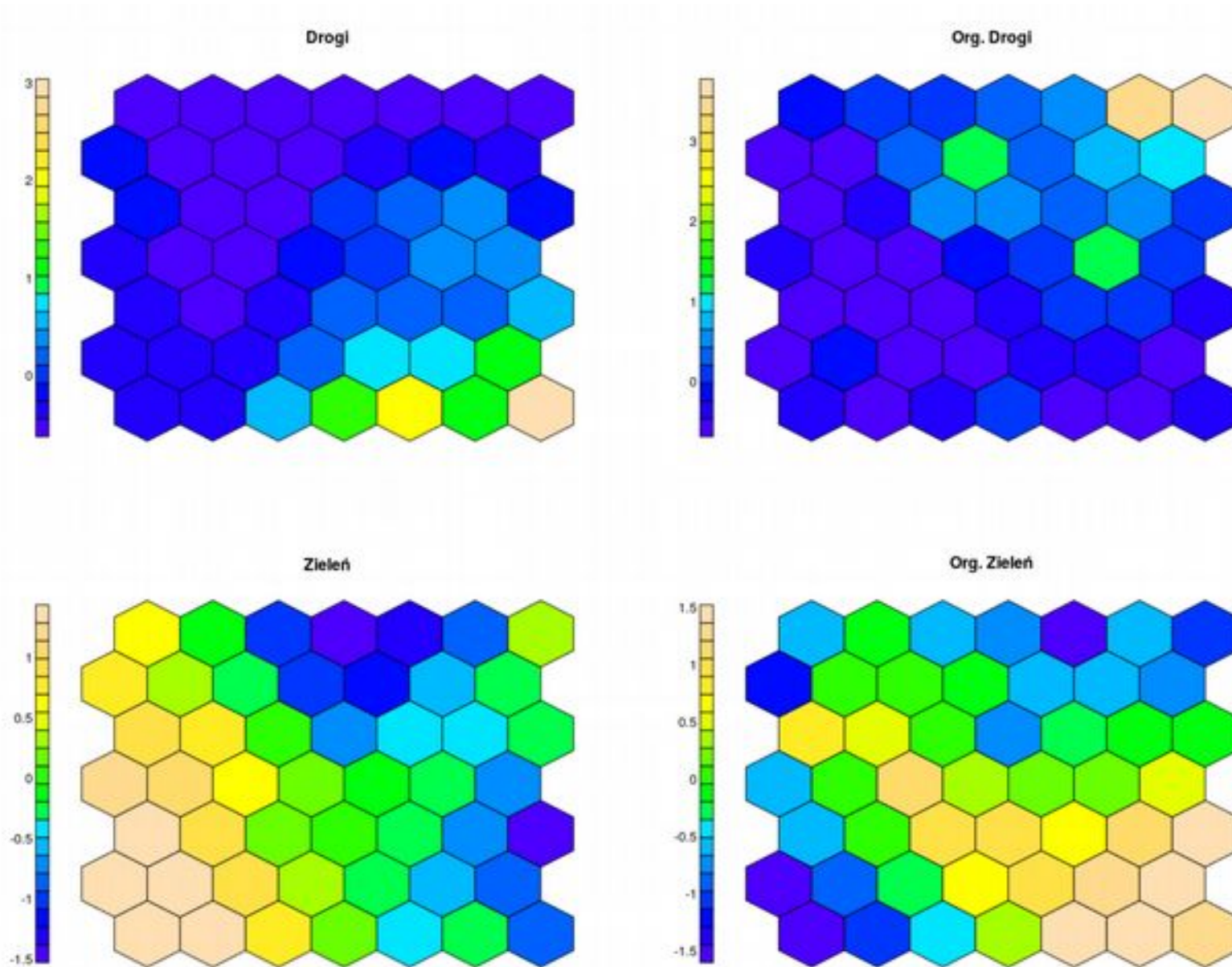
Counts plot



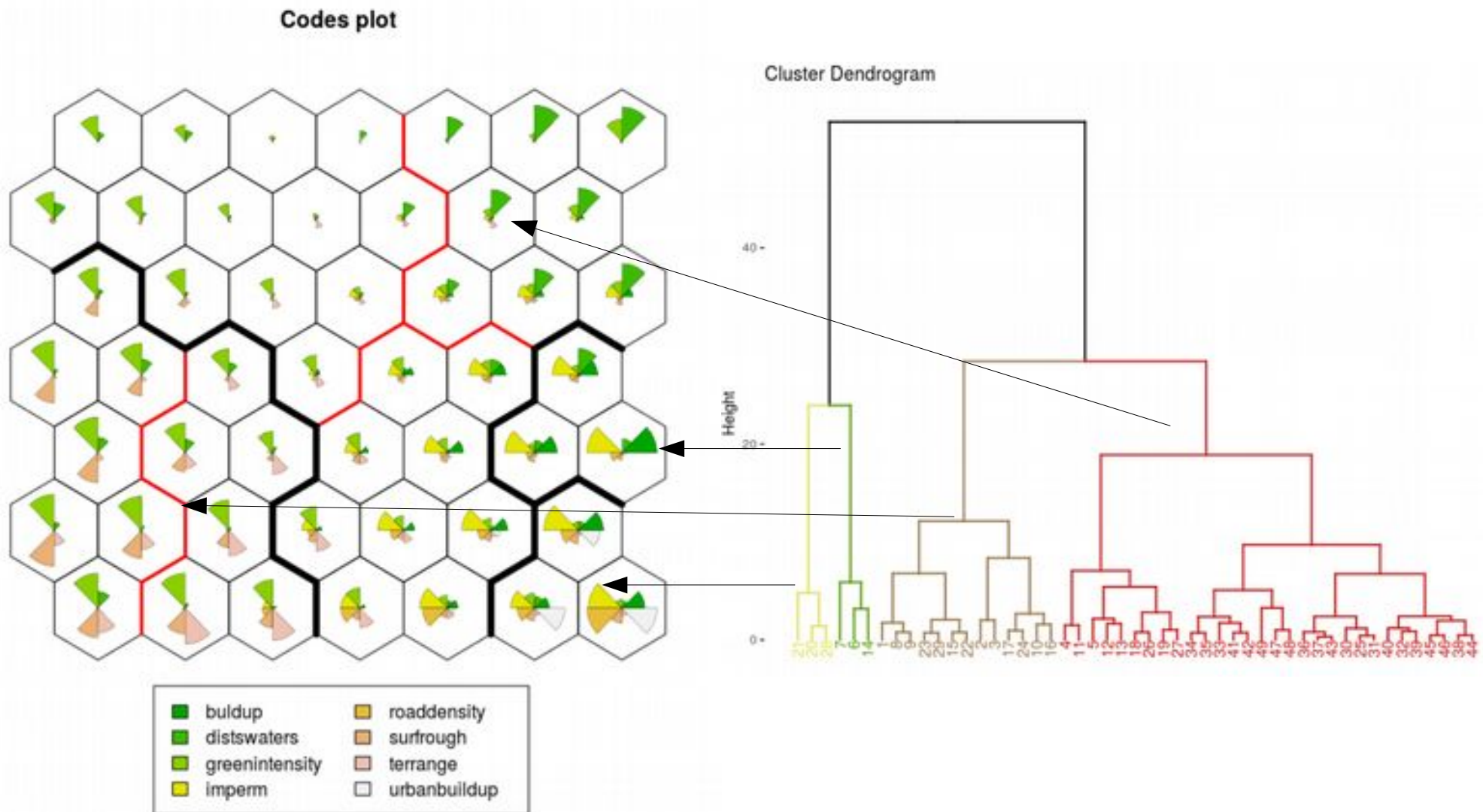
Quality plot



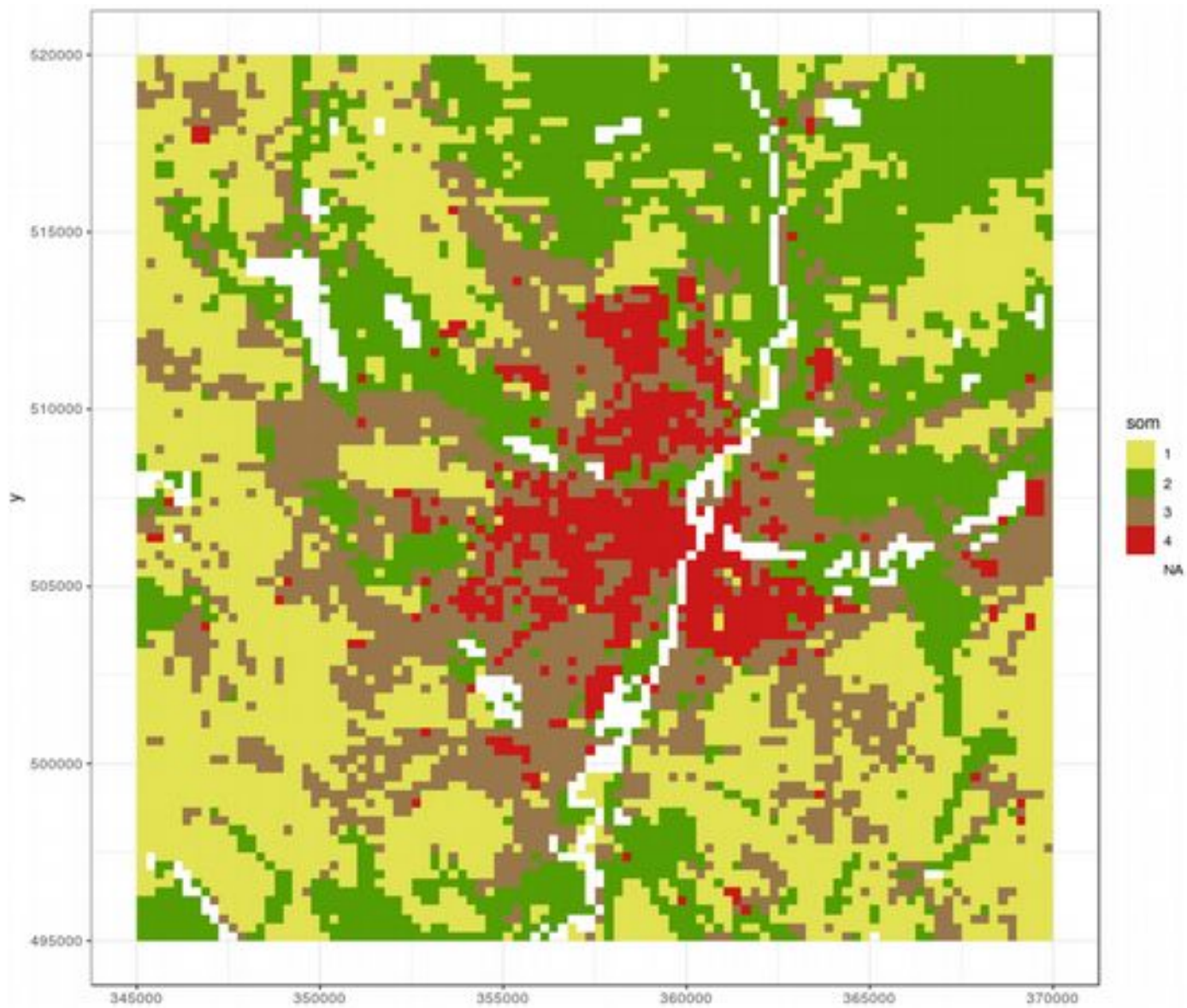
Codebook a dane oryginalne



SOM w klasyfikacjach nienadzorowanych



SOM a dane geoprzestrzenne



Grupowanie a klasyfikacja

- Grupowanie jest procesem budowania optymalnych skupień, proces klasyfikacji to nadawanie skupieniom znaczenia (*labeling*) *a posteriori*
- Optymalne skupienia nie muszą odpowiadać optymalnym klasom – klasy są pochodnymi badanego problemu: np. różnicowanie pokrywy roślinnej jest statystycznie większe; ale mniej istotne punktu widzenia człowieka niż np. różnicowanie pokrycia zabudowy
- W przypadku klasyfikacji danych uporządkowanych (np. geoprzestrzennych) do interpretacji klasy ma znaczenie nie tylko charakterystyka obiektów ale również ich położenie

Najczęstsze błędy w klasyfikacjach nienadzorowanych

- Brak transformacji danych (standaryzacji/normalizacji)
- Nieodpowiednie miary niepodobieństwa (nadużywanie metryki euklidesowej)
- Brak redukcji wymiarów i obiektów odstających
- Wymuszanie skupień w jednorodnych danych
- Stosowanie metod hierarchicznych do dużych zbiorów danych
- Sugerowanie się klasami *a priori* (zamiast metod nadzorowanych)

