

Praca Domowa 2 Raport

Jan Borowski

Contents

Analiza hclust	2
Analiza parametrów dodatkowych	3
Porównanie wydajności	4
Analiza dokładności	5
Wpływ standaryzacji	7
Zbiory	8
Podsumowanie	9

Wstęp

(W całym raporcie używam określenia mój oczywiście chodzi o moją implementację algorytmu spektralnego)
W poniższym raporcie dokonam porównania mojej implementacji algorytmu spektralnego (z parametrem $M = k, k+2, k+6$ gdzie k = liczba skupień) z następującymi algorytmami :

- Wszystkimi algorytmami hierarchicznymi z funkcji `hclust()`.
- Algorytmem *Genie* z pakietu `genie`.
- Algorytmem *hdbscan* (z parametrem `minPts=5,10,20,50,100`) z pakietu `dbscan`.

Porównania będą dotyczyły różnych aspektów działania algorytmów takich jak dokładność i czas. Dane będą przedstawiane w postaci różnego rodzaju wykresów. Surowe informacje znajdują się w przesłanych plikach .csv. Do zbierania danych używałem zbiorów benchmarkowych oraz dołączonego do nich zestawu moich zbiorów przedstawionych w pliku Testy.pdf.

Analiza hclust

Zacznę od analizy poszczególnych metod funkcji **hclust** aby wyznaczyć najlepsze metody do dalszych testów. Testy wykonywałem na 13 losowo wybranych zbiorach.

Najpierw przeprowadzę analizę średniego czasu działania poszczególnych metod :

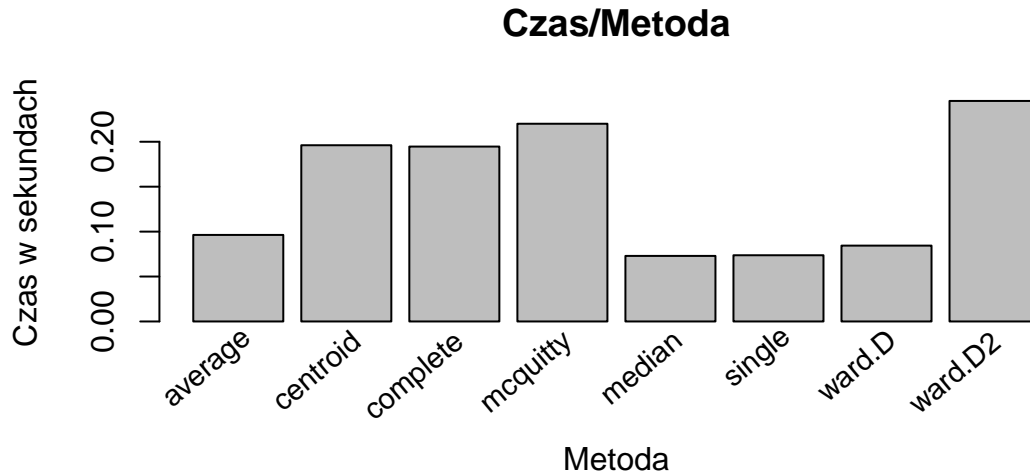


Figure 1: Zależność czasu od metody

Teraz dokonam analizy działania algorytmów porównując indeks Randa i Fowlkesa–Mallowsa :

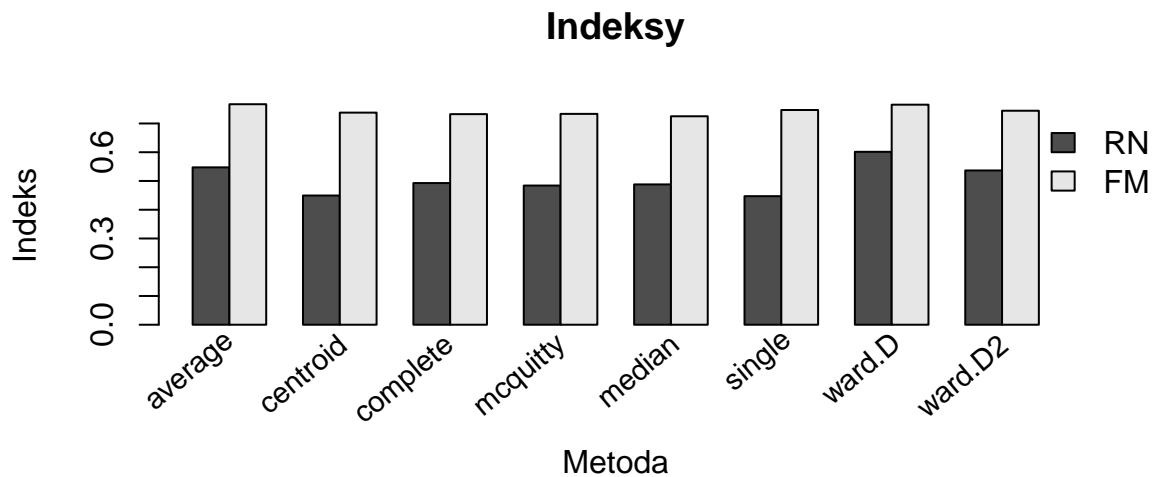


Figure 2: Zależność parametrów jakościowych od metody

Jak widać najszybsza metoda to **median** a najwolniejsza **ward.D2**. Te dwie metody wybrałem do dalszych testów wydajnościowych i porównania z innymi algorytmami. Należy jednak zauważyć, że choć w przypadku indeksu Fowlkesa–Mallowsa nie widać znaczących różnic. Metoda **ward.D** uzyskała najwyższy indeks Randa.

Analiza parametrów dodatkowych

Najpierw sprawdzę dla jakiego M moja implementacja działała najkuteoczniej (Testy przeprowadzane na 46 zbiorach):

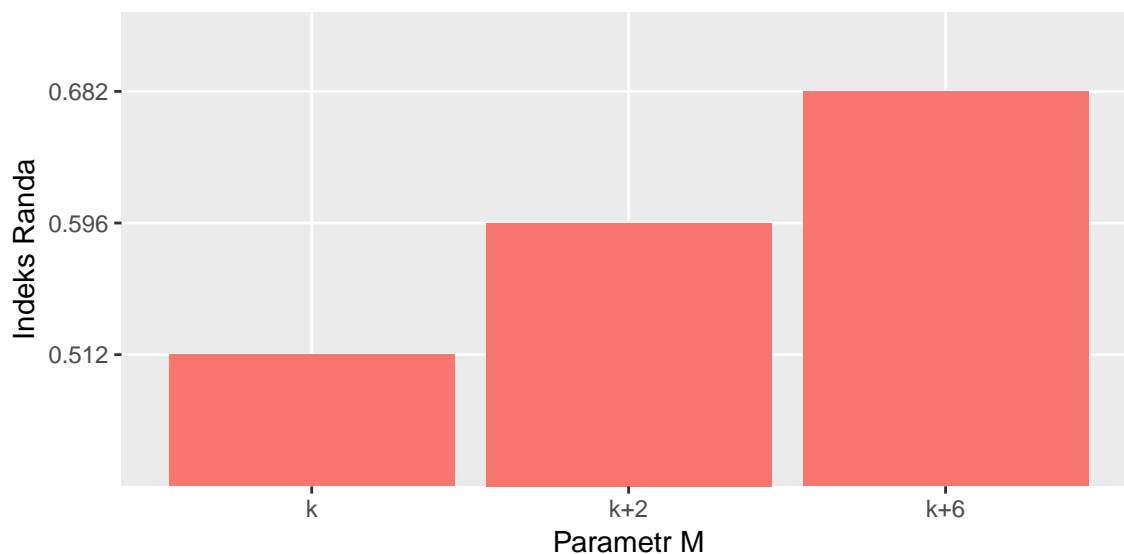


Figure 3: Zależność dokładności od parametru M

Jak widać moja implementacja uzyskała najlepszą średnią skuteczność przy $M=k+6$.
Sprawdź teraz funkcję **hdbscan**:

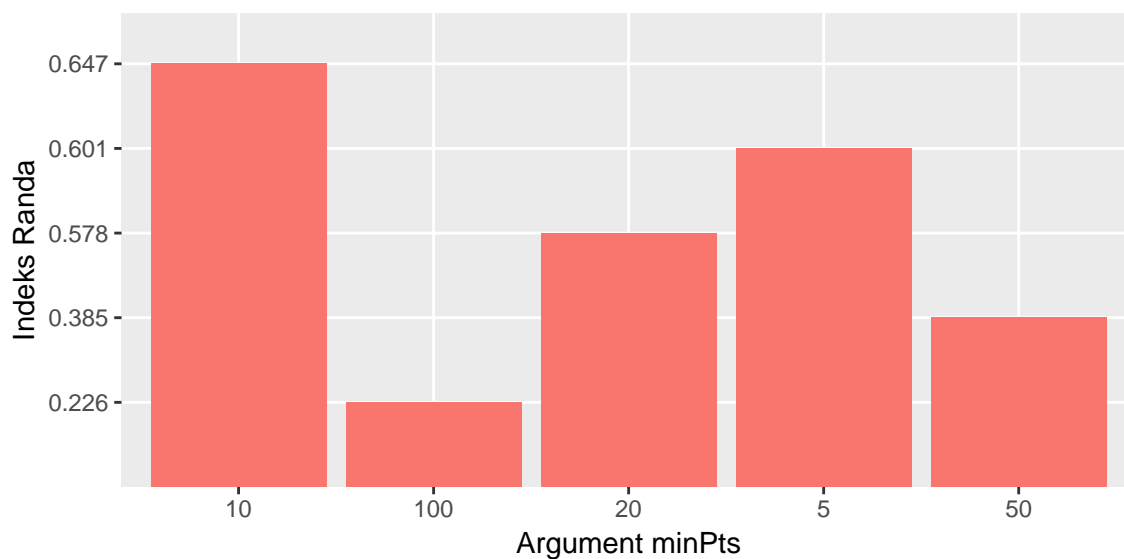


Figure 4: Zależność dokładności od parametru minPts

Jak widać funkcja **hdbscan** działa średnio najdokładniej przy parametrze $\text{minPts} = 10$.

Porównanie wydajności

Przetstawię średnie czasy wykonania algorytmów dla **Genie** z parametrem **thresholdGini** = 0.3 oraz **Genie** z **thresholdGini** = 1 , wybranych metod funkcji **hclust** ,funkcji **hdbscan** oraz mojej implementacji algorytmu spektralnego. W przypadku funkcji **hdbscan** oraz mojej implementacji gdzie występuję dodatkowy parametry wybierany jest ten dla którego czas wykonania był najmniejszy (Testy wykonywane na 46 zbiorach testowych):

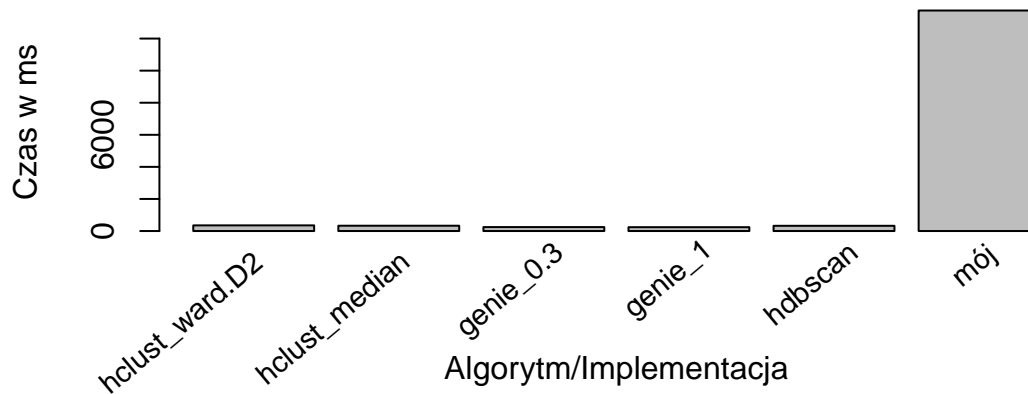


Figure 5: Zależność czasu od algorytmu/implementacji

Jak widać zdecydowanie najwolniejsza jest moja implementacja więc narazie odrzucę ją z dalszej analizy wydajności.

Przedstawię teraz analizę czasów działania pozostałych algorytmów:

Boxploty czasu od algorytmu

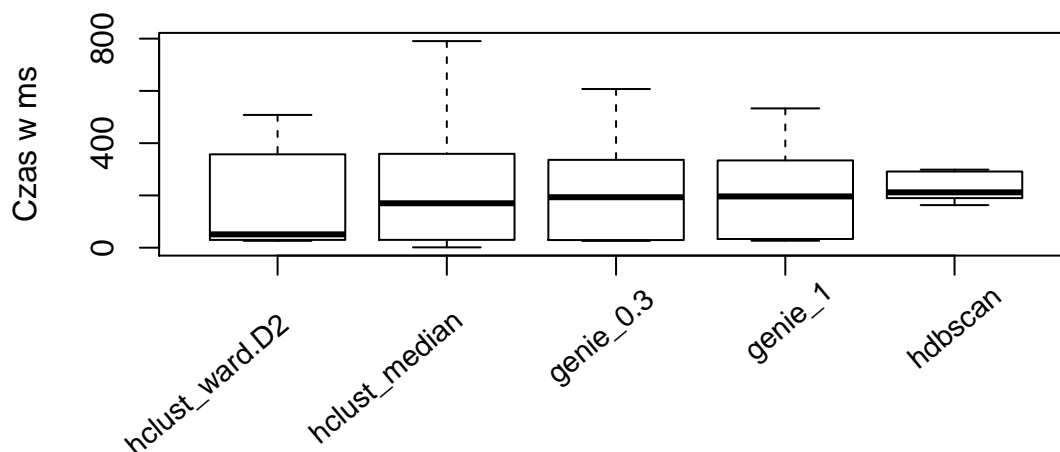


Figure 6: Zależność czasu od algorytmu

Jak widać po rozpatrzeniu danych można po pierwsze zauważyć ,że najwolniejsza metoda **hclust** okazała się najszybsza na większej próbce (benchmarki wykonywane z argumentem times = 10) . Potem mniej więcej równo jeśli chodzi o medianę **genie** i metodę **median** z funkcji **hclust** najwolniejszy okazał się algorytm **hdbscan** ale w jego wypadku mamy do czynienia z bardzo małym rozrzutem. Tym niemniej można ułożyć algorytmy w kolejności od najszybszego do najwolniejszego :

moja implementacja >hdbscan >hclust_median>genie>hclust_ward.D2

Analiza dokładności

W tej sekcji pod uwagę będzie brana metoda **ward_D** z funkcji **hclust**. Najpierw sprawdzę średni indeks Randa w zależności od algorytmu. Jak poprzednio w przypadku metod z dodatkowymi parametrami wybiorę ten o najwyższym indeksie (Testy wykonywane na 46 zbiorach testowych) :

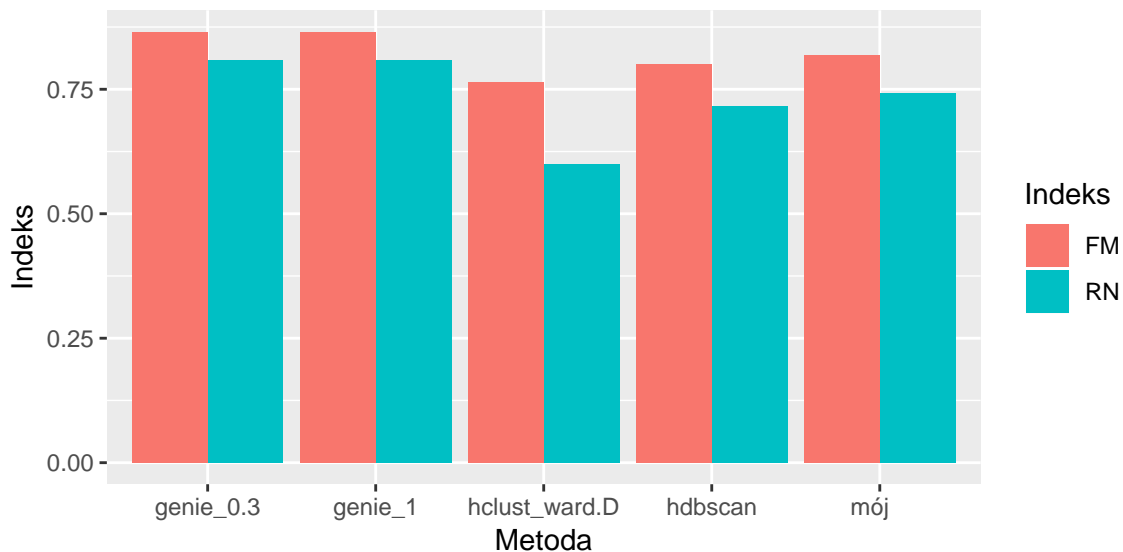


Figure 7: Zależność indeksów od metody

W tym wypadku wyraźnie widać że najlepiej działający algorytm to **genie**. Co oznacza, że dokonuje on najlepszych podziałów.

Przedstawię teraz zależność średniego indeksu Randa od średniego czasu wykonania algorytmu (czas w ms) :

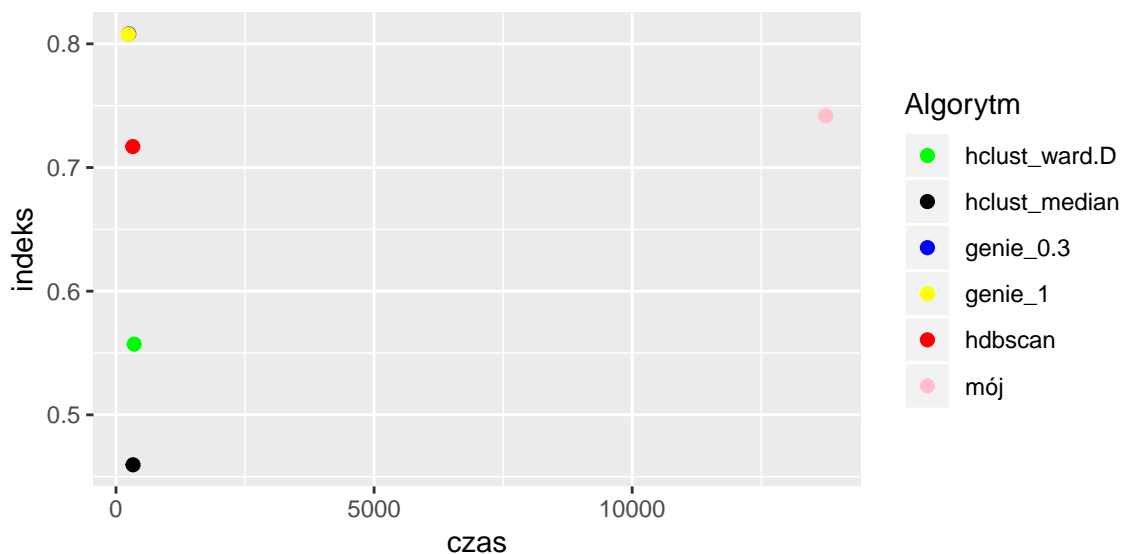


Figure 8: Zależność jakości od czasu 1

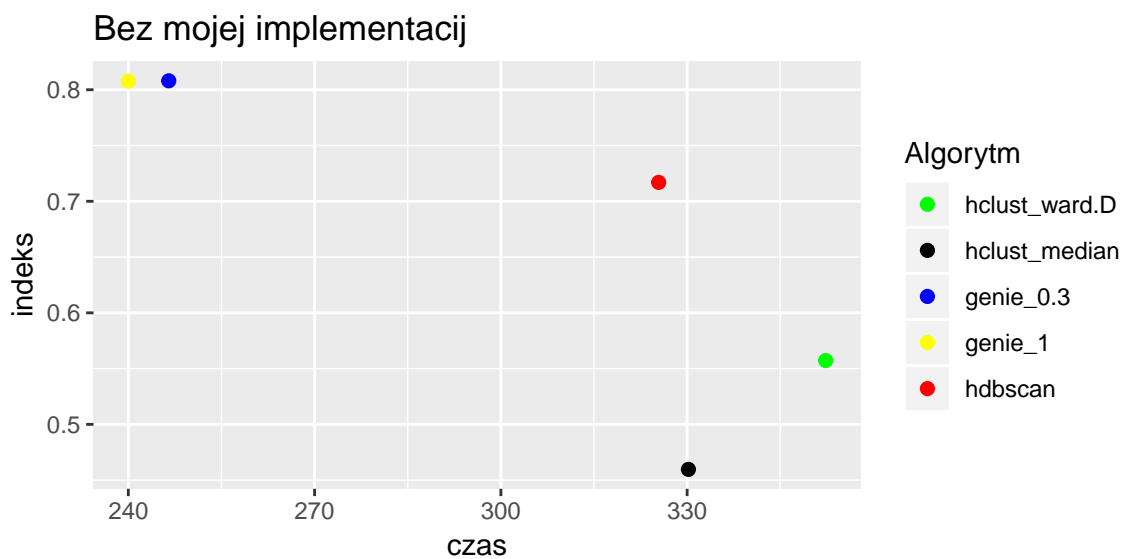


Figure 9: Zależność jakości od czasu 2

Dopiero ten wykres pokazuje skuteczność **genie** które uzyskują bardzo wysoką jakość w krótkim czasie .

Wpływ standaryzacji

W tej części sprawdzę wpływ standaryzacji na dokładność działania algorytmów (testy przeprowadane na 14 losowo wybranych zbiorach).

Najpierw sprawdzę wpływ standaryzacji na Indeks Randa:

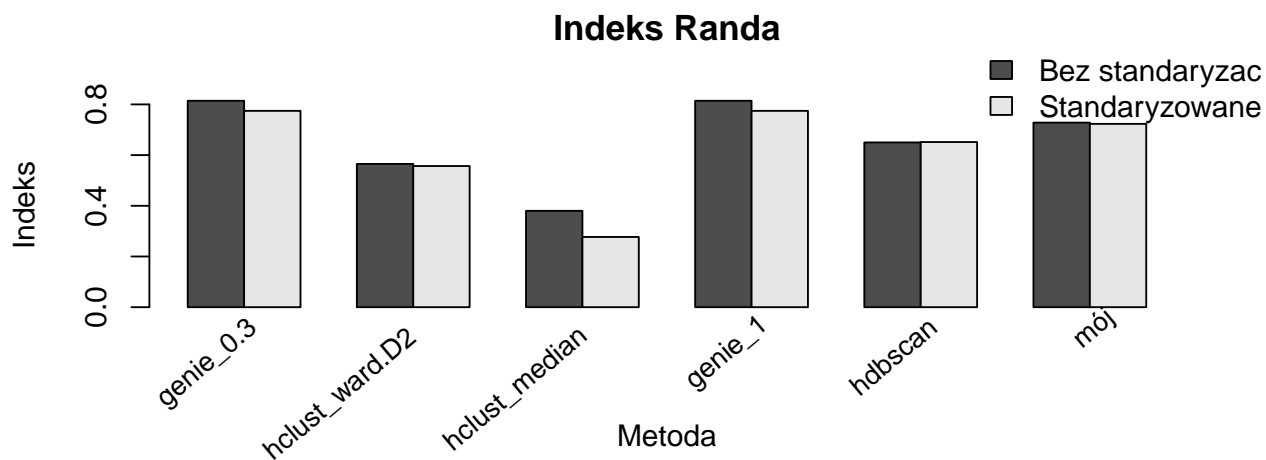


Figure 10: Zależność Indeksu Randa od Standaryzacji

Jak widać w przypadku niektórych algorytmów wpływ jest niezauważalny a w innych wypadkach negatywny. Teraz indeks Fowlkesa–Mallowsa:

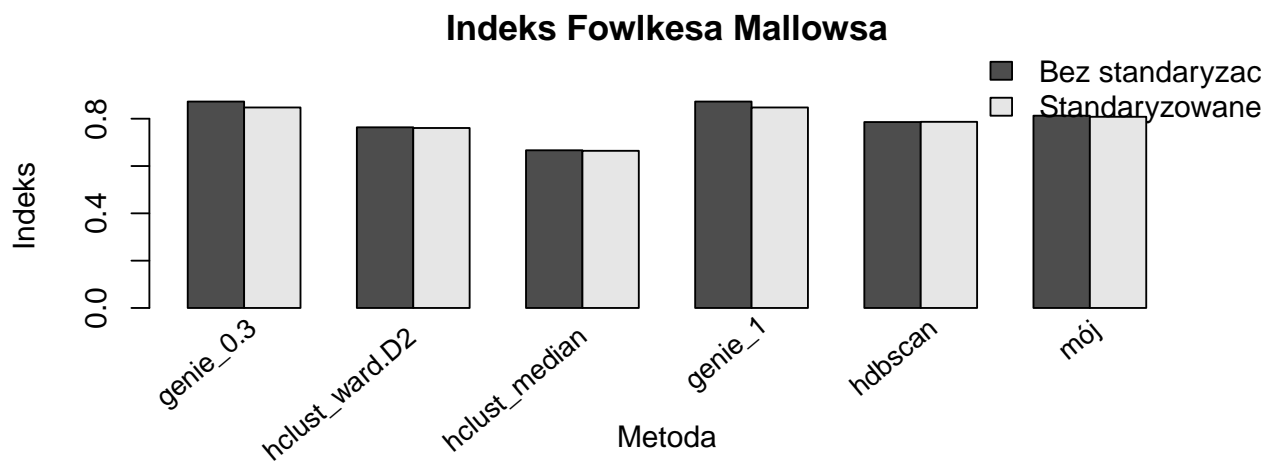


Figure 11: Zależność Indeksu Fowlkesa–Mallowsa od Standaryzacji

Jak widać drugi parametr jakościowy również wskazuje na to samo chociaż w mniejszym stopniu. Mimo to myślę, że można wyciągnąć wniosek o negatywnym wpływie standaryzacji przy analizie skupień.

Zbiory

Teraz zaprezentuję z którym zbiorem najgorzej poradził sobie dany algorytm:

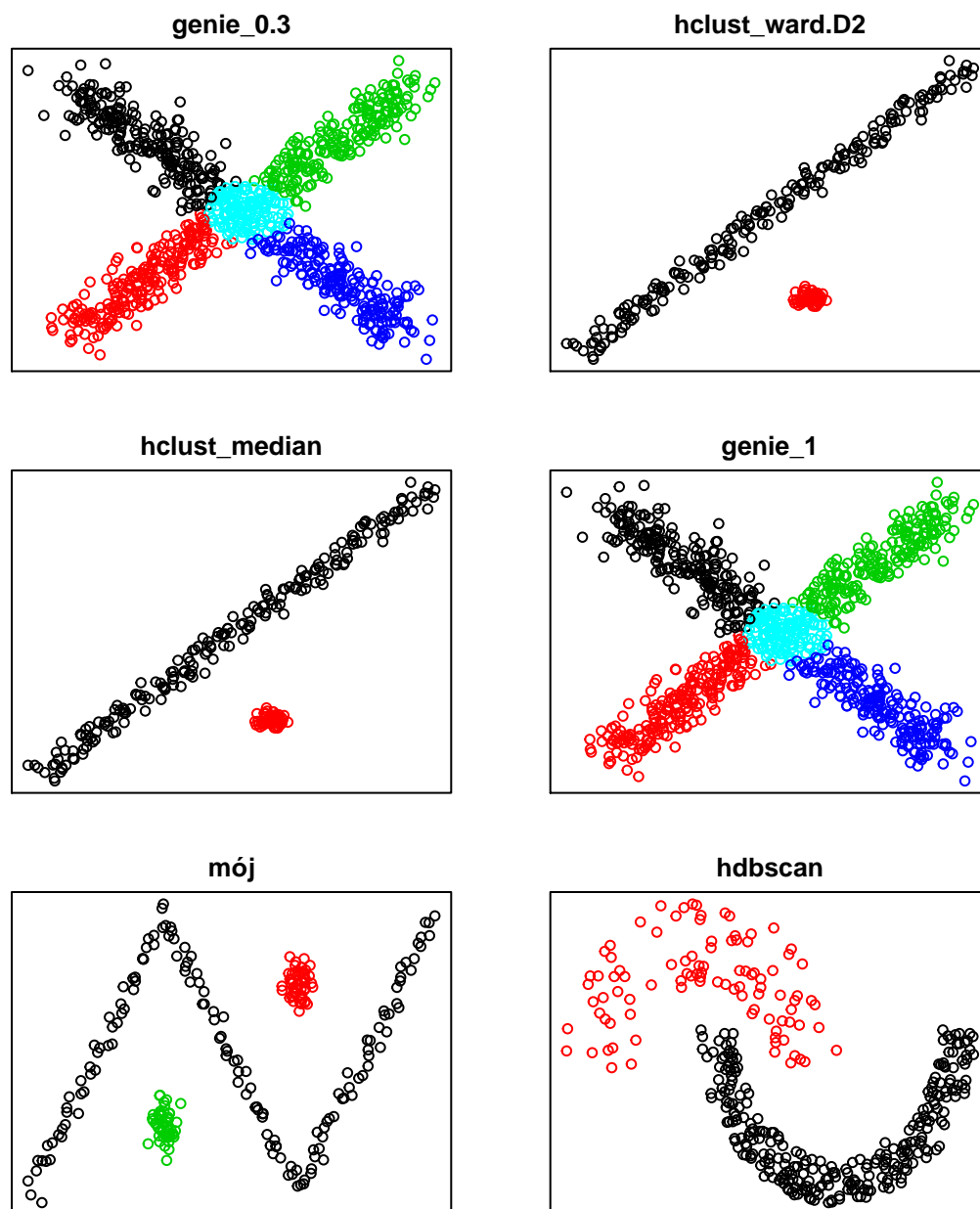


Figure 12: Najgorzej podzielone zbiory

Podsumowanie

Wszystkie przytoczone przeze mnie wykresy wskazują na wyraźną dominację algorytmu **genie** który nie okazał się najszybszy . Pomimo tego przy bardzo wysokiej prędkości działania zapewnia najlepszą skuteczność z badanych algorytmów/implementacji . Moja własna implementacja okazała się około 100 razy wolniejsza od **genie** ale nie wiele od niej mniej dokładna. Algorytm **hdbscan** pomimo działania bez podanej odgórnie liczby skupień okazał się dokładniejszy niż metody funkcji **hclust**.