

Wstęp

W analizie skupień która przeprowadziłem zdecydowałem się na użycie zgodnie z poleceniem własnego algorytmu analizy skupień, wszystkich algorytmów z pakietu stats, algorytmu genie oraz znalezionej przeze mnie algorytmu cmeans z pakietu e1071 dostępnego na CRANie. O ostatnim z wymienionych algorytmów rozpiszę się troszeczkę później gdy będę prezentował wyniki przez niego uzyskane, gdyż z pewnością zasada jego działania jest interesująca.

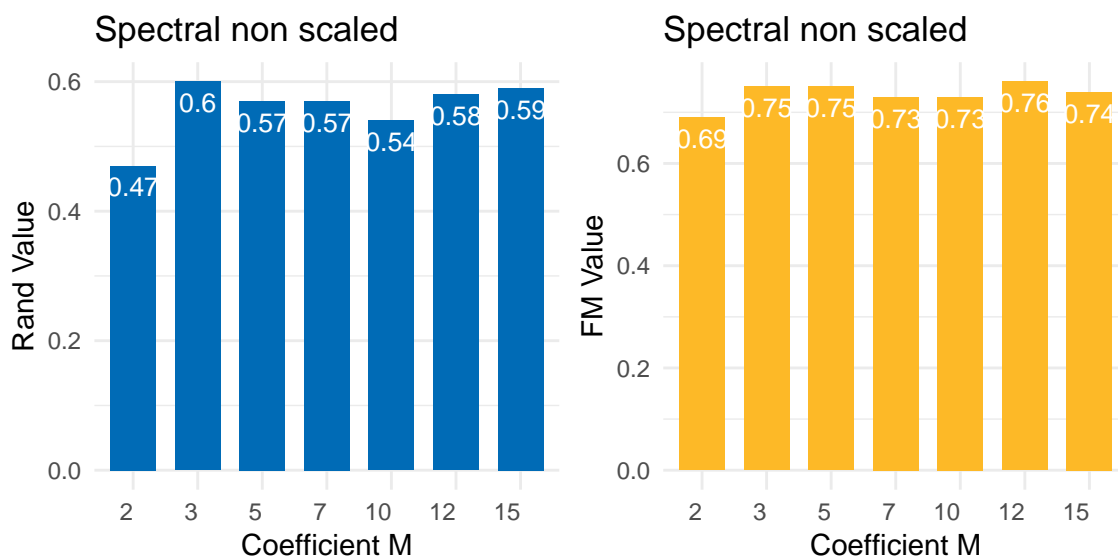
Z racji na dużą ilość zbiorów (łącznie z dodanymi przeze mnie jest ich 46) uznałem iż znaczną część swojej analizy oprę o pokazywanie wyników średnich uzyskiwanych przez algorytmy. Zgodnie z wytycznymi testy przeprowadziłem najpierw na zbiorach normalnych, a następnie na zbiorach ustandaryzowanych.

To po czym będziemy decydować czy algorytm jest dobry bądź zły będą etykiety referencyjne, przygotowane przez ekspertów. Do porównania zgodności wyników zwracanych przez algorytm i danych z góry etykiet będziemy stosować: indeks Fowlkesa-Mallowsa (`dendextend::FM_index()`) oraz skorygowany indeks Randa (`mclust::adjustedRandIndex()`). Każdy z powyższych indeksów zwraca wartość równą 1, jeśli dane podziały są równoważne. Im wartość jest mniejsza od 1 tym gorzej dany algorytm się sprawuje i tym mniejszą skutecznością się on charakteryzuje.

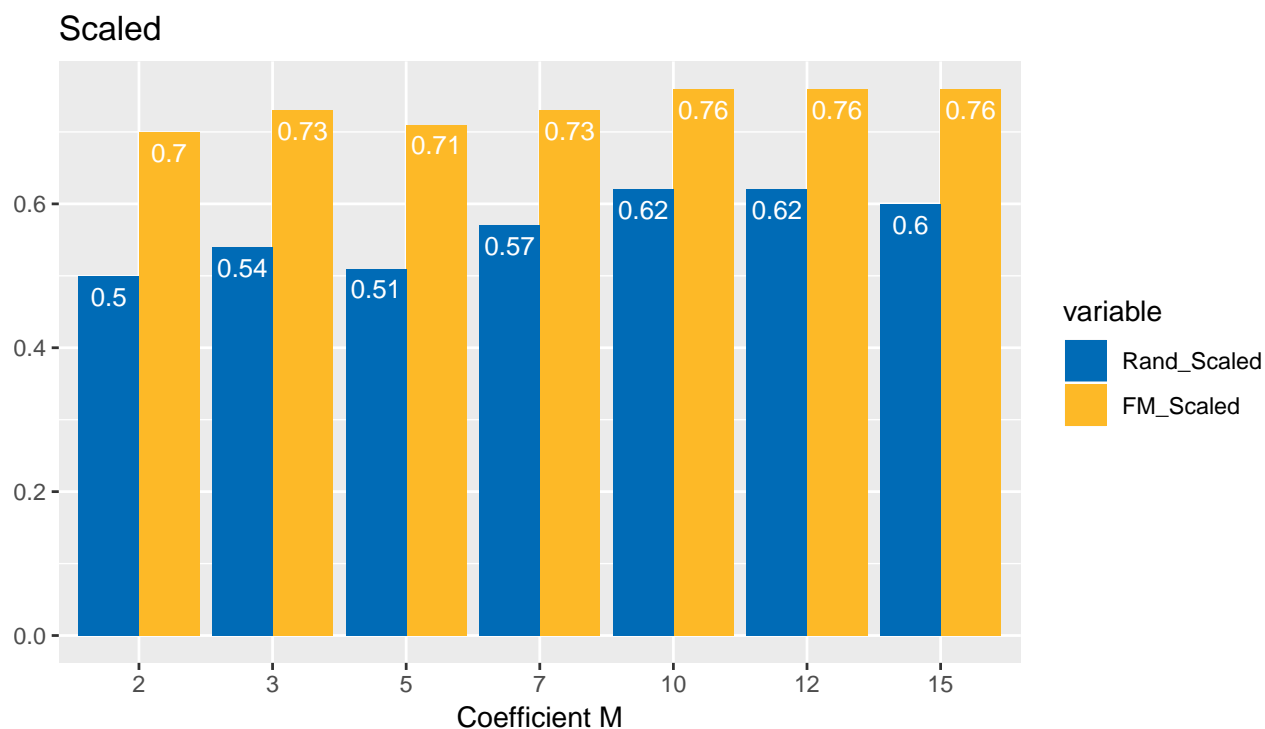
Algorytm spektralny

Prezentacje rozpocznę od pokazania działania algorytmu zaimplementowanego przeze mnie w zależności od różnych wartości parametru M . Parametr ten określa sposób generowania macierzy sąsiedztwa dla punktów analizowanego zbioru. Ponieważ działanie tego algorytmu opiera się na wyłonieniu k losowo wybranych centroidów, za każdym razem, gdy algorytm jest wywoływany, można uzyskać inny rezultat. Ponadto proces clusteringu jest wrażliwy na początkowy wybór środków. Funkcja `kmeans()` ma parametr `nstart`, który zaczyna od wielu początkowych konfiguracji i wybiera tę najkorzystniejszą. Uznałem ze ustawię opcję `nstart=25`, aby zwiększyć szansę na otrzymanie dobrych wyników.

Oto wartości odpowiednio współczynników indeksów Randa oraz FM.



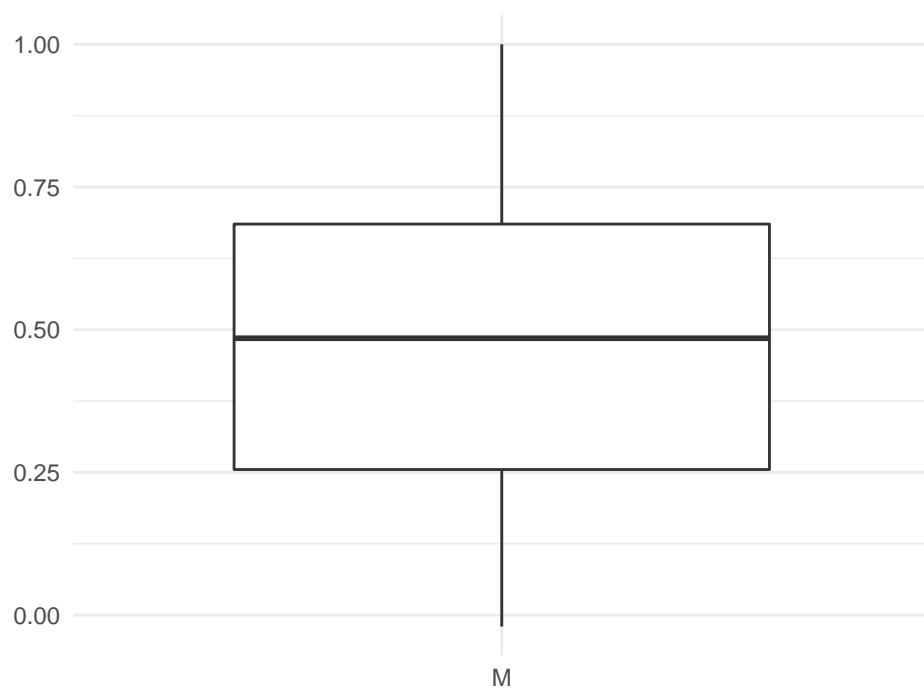
Jak widać na wykresach najlepiej algorytm prezentuje się dla parametru M równego 3 a także dla $M=12-15$. Jednakże różnice pomiędzy poszczególnymi wybranymi przeze mnie wartościami M nie są zbyt duże i nie powodują one znaczącej różnicy w wynikach działania. Sprawdźmy jak algorytm sprawuje się dla danych ustandaryzowanych.



Standaryzacja w niewielki sposób wpłynęła na wyniki, pozwoliła jedynie podnieść wynik indeksu Randa 0.02, co jest znikomą poprawą.

Zdecydowałem się również sprawdzić jak będzie wyglądało rozłożenie wyników osiągniętych przez algorytm dla poszczególnych zbiorów, dla najlepszego M. Chcę w ten sposób sprawdzić czy algorytm generuje stabilne wyniki, czyli czy można ufać jego wynikom. Dla większej czytelności użyję jedynie indeksu Randa, aby lepiej móc zobaczyć sprawowanie się algorytmu.

Rand index on best one



Widać ze wyniki są dość mocno rozbieżne dla różnych zbiorów co widać poprzez dość szerokie pudełko na wykresie pudełkowym. Widać tutaj, iż algorytm nie do końca dobrze radzi sobie z bardziej skomplikowanymi zbiorami.

Jeszcze jedna rzecz która chciałbym sprawdzić jest dokładnie ile wyników wpada do jakiego przedziału dla pokazanego powyżej algorytmu puszczanego dla jednego konkretnego M . Uważam iż w ten sposób będzie można dokładnie zobaczyć jak prezentuje się rozkład danych, które algorytm daje.

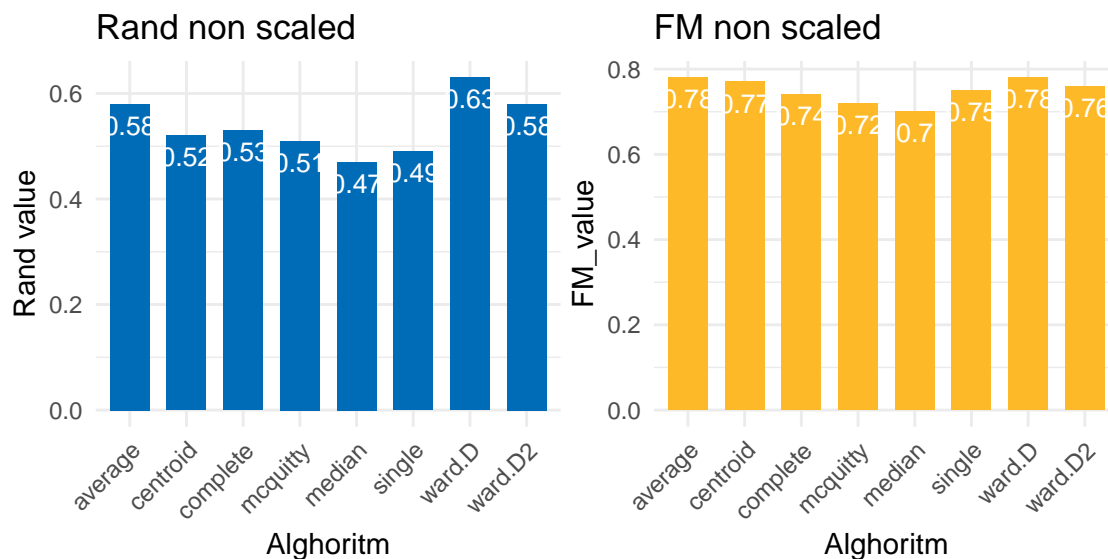
Przedziały	Wartości
0	1
(0,0.25]	11
(0.26,0.5]	14
(0.51,0.75]	11
(0.76,1]	9

Tak jak widać na wykresie pudełkowym największa część danych wpadła do środkowych kubełków. Jednocześnie widzimy teraz jednoznacznie że do najlepszego kubełka wpadło 9 obserwacji. Jest to wartość która na pewno należy zapamiętać do porównania względem następnych algorytmów.

Nawiazuając do tego co napisałem w raporcie podsumowującym stworzone przeze mnie zbiory, widać że algorytm osiąga niższe średnie wyniki niż te które miał dla moich zbiorów. Potwierdziło się więc moje przypuszczenie że to “prostość” moich zbiorów miała wpływ na wynik algorytmu. Podsumowując ten algorytm można powiedzieć że działa on w miarę poprawnie, jednakże jego wyniki nie są oszałamiające. W następnej części raportu pozwolę sobie porównać jego wyniki z innymi algorytmami które testowałem.

Stats

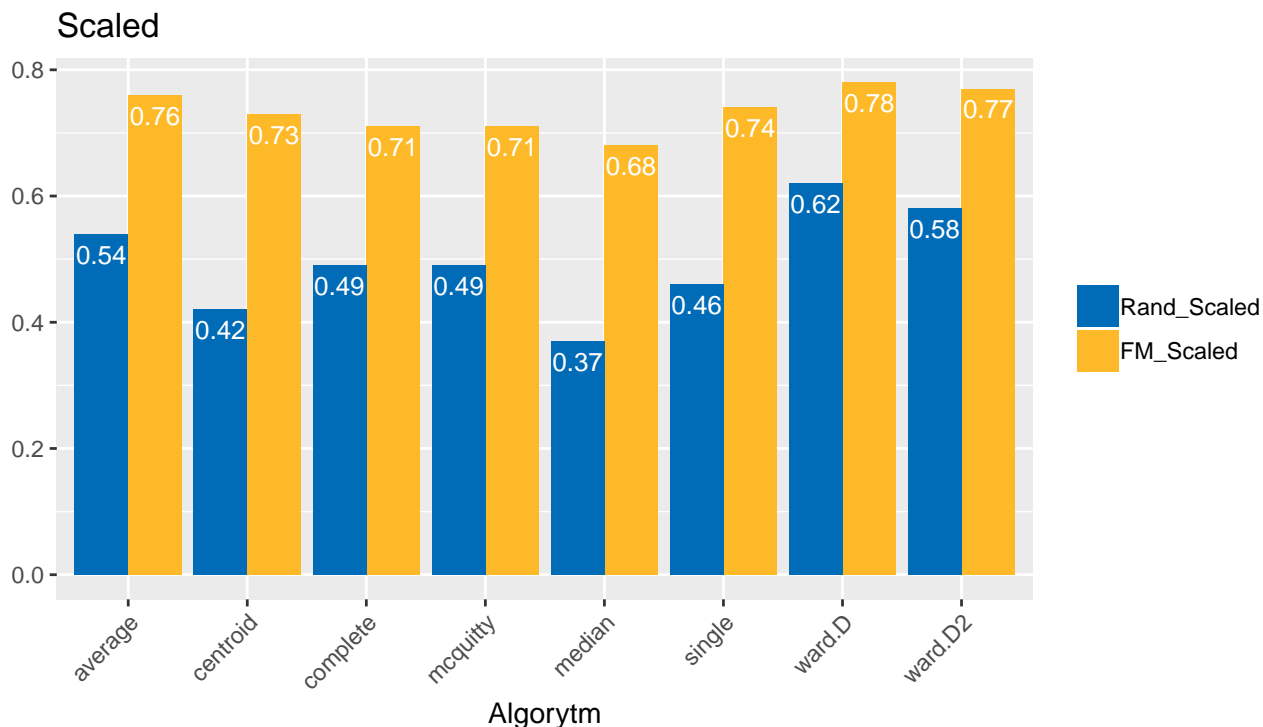
W tej czesci raportu zaprezentuje dzialanie algorytmow z pakietu stats. Zawarte jest w nim 8 algorytmow grupowania hierarchicznego. Przeanalizowalem dzialanie kazdego z nich co przedstawie teraz na wykresach.



W tym przypadku widac juz roznice w stosunku do poprzedniego algorytmu, gdyz jest ‘lider’ ktory osiaga najlepsze wyniki. Liderem tym jest algorytm “ward.D”.

Rowniez wyniki indeksu FM zdaja sie to potwierdzac, jednakze inne algorytmy rowniez osiagaja dobre wyniki. Widac jednak ze algorytmy zawarte w pakiecie stats sprawuja sie trozke lepiej niz algorytm napisany przeze mnie. Z pewnoscia wynika to z jakosci napisanego kodu, jak i innego podejscia do problemu.

Sprawdzmy czy tym razem standaryzacja pozwoli poprawic wyniki.



Jak widac standaryzacja nie poprawila wynikow, a miejscami wręcz niepokojaco pogarsza wyniki.

Tak samo jak przy algorytmie spektralnym wykonam wykres pudełkowy, aby zobaczyć jak rozkładają się wyniki najlepszego algorytmu z pakietu stats. Dzięki temu będziemy mogli wyciągnąć już ciekawe wnioski dotyczące zachowania oraz stabilności algorytmów.

Rand index on best one



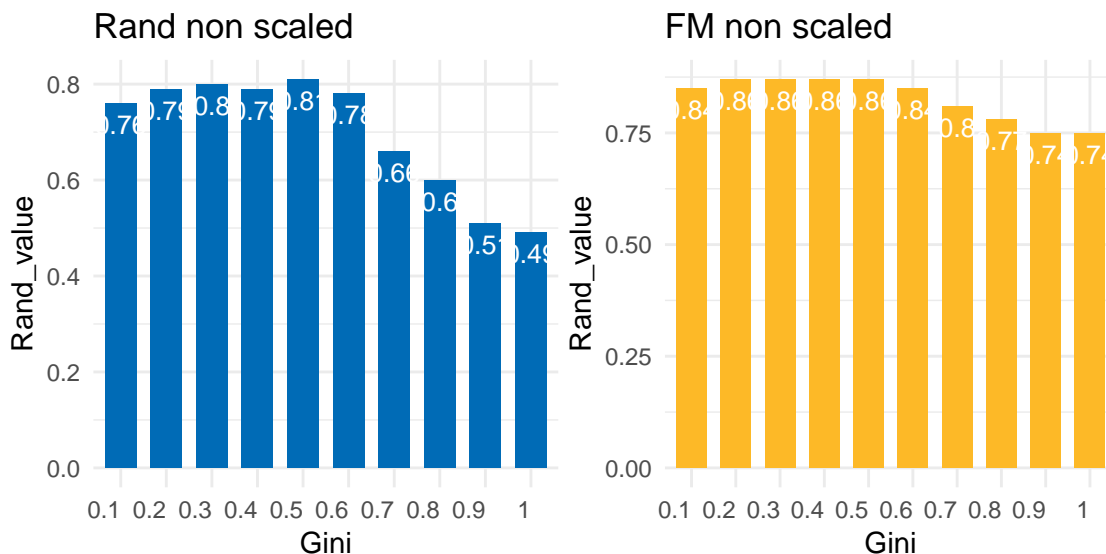
Widać tutaj jeszcze większy rozrzut wyników niż w przypadku algorytmu spektralnego. Pozwala to sądzić że algorytm hclust jest jeszcze mniej odporny na różnego rodzaju odchyły przy różnych zbiorach, jednakże trzeba mu oddać że prezentuje znacznie lepsze wyniki niż stworzony przeze mnie algorytm, jest po prostu skuteczniejszy.

Przedziały	Wartosci
0	1
(0,0.25]	8
(0.26,0.5]	7
(0.51,0.75]	9
(0.76,1]	21

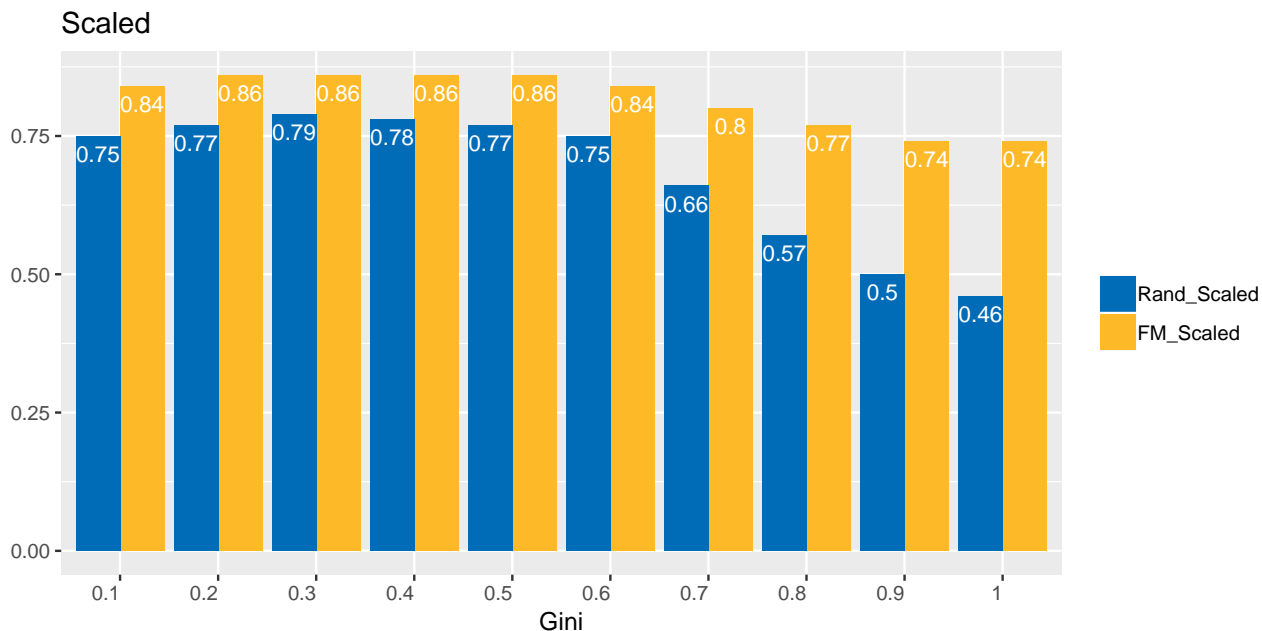
Widać znaczną poprawę względem poprzedniego algorytmu. Mamy znacznie więcej obserwacji w najwyższym kubku, co oczywiście oznacza większą skuteczność działania.

Genie

Kolejnym algorytmem który przetestowałem jest hclust2 z pakietu Genie. Pozwala on na różne dostosowanie parametru thresholdGini w zależności od którego osiąga on rozbieżne wyniki. Zdecydowałem się sprawdzić jak ten algorytm będzie sprawował w zależności od tego parametru.



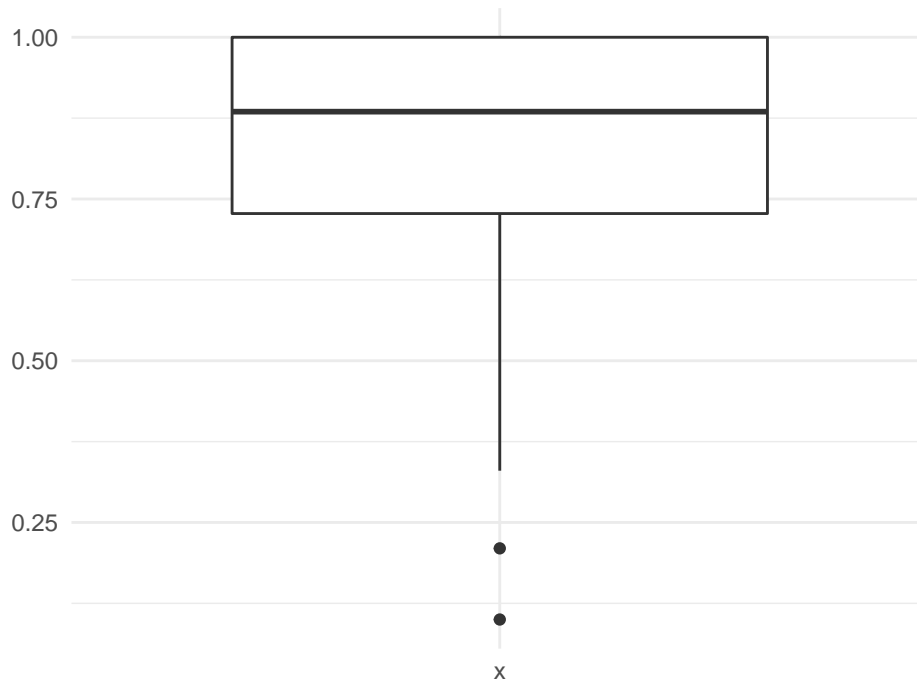
Udało się osiągnąć takie same wyniki dla thresholdGini równego 0.3 i 0.5. Widać również że ten algorytm radzi sobie bardzo dobrze z różnymi rodzajami zbiorów danych. Zgodnie z tym co jest napisane w opisie algorytmu na CRANie sprawuje się on lepiej niż algorytm ward.D z biblioteki stats, który z kolei był najlepszy ze wszystkich dostępnych w niej algorytmów. Różnica w wynikach jest dość znacząca bo wynosi prawie 0.1, co jest spora różnica w jakości działania. Jednocześnie długość działania algorytmu hclust2 nie jest większa, więc śmiało można powiedzieć że jest on lepszy niż hclust z biblioteki stats. Sprawdźmy czy może tym razem standaryzacja danych wejściowych pomoże poprawić wyniki.



Tym razem standaryzacja nie przyniosła praktycznie żadnych zmian w wynikach, więc można powoli dojść do wniosku, na podstawie działań tych algorytmów które przedstawiłem do tej pory, iż nie ma ona zbyt wielkiego wpływu na wyniki, bądź nawet czasem powoduje gorsze sprawowanie się algorytmu.

Poraz kolejny sprawdźmy rozrzut wyników stworzonych przez algorytm.

Rand index on best one



Tutaj widać już znaczną różnicę względem poprzednich algorytmów. Rozstrzał wyników jest dużo mniejszy, a i same wyniki są znacznie wyższe. Śmiało można stwierdzić że ten algorytm prezentuje się zdecydowanie najlepiej. Generuje najlepsze wyniki i robi to w sposób najbardziej stabilny, co powoduje, że ciężko mieć do niego jakiegokolwiek zastrzeżenia.

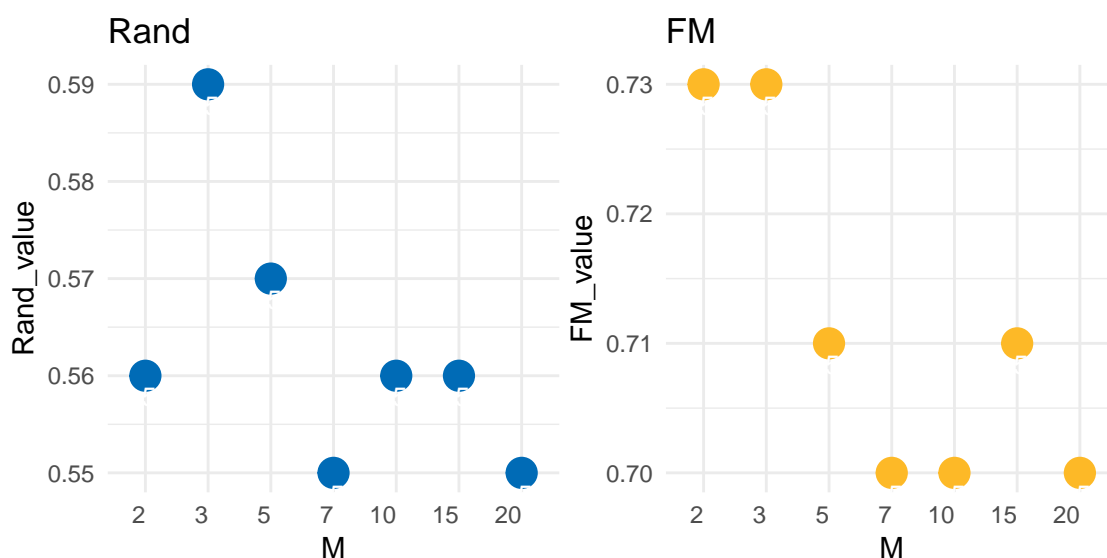
Sprawdźmy jeszcze dokładny rozkład wyników w poszczególnych kubelkach.

Przedziały	Wartosci
(0,0.25]	2
(0.26,0.5]	5
(0.51,0.75]	5
(0.76,1]	34

Algorytm uzyskał oszałamiające 34 na 46 wyników w najwyższych kubelku, co oznacza prawie 75% zbiorów udało mu się sklasyfikować co najmniej bardzo dobrze. Jest to zdecydowanie najlepszy wynik spośród wszystkich z porównywanych przeze mnie algorytmów.

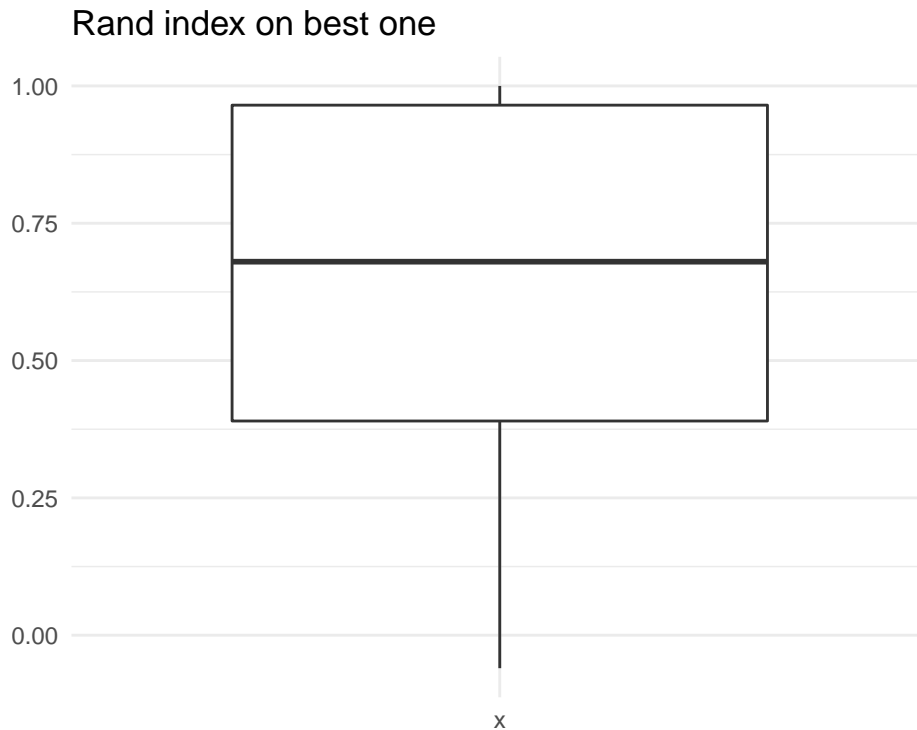
Cmeans

Jako swój algorytm zdecydowałem się użyć algorytmu cmeans z pakietu e1071. Przeprowadza on tzw. fuzzy clustering, który charakteryzuje się tym iż tutaj przykładowe punkty mogą należeć do kilku clusterów, co w normalnym procesie clusteringu nie ma miejsca. Algorytm ten zwraca również niejako prawdopodobieństwo tego że jakiś punkt może należeć do któregoś clusteru, co również wyróżnia go od innych algorytmów. Cmeans udostępnia parametr M, który określa “zamglenie” zbioru. Sterując tym parametrem można zmieniać jego zachowanie, a co za tym idzie uzyskiwać inne wyniki. Z racji tego że poprzednie próby pokazały że standaryzacja nie daje praktycznie żadnych pozytywnych efektów w działaniu algorytmów zdecydowałem, iż teraz nie będę przedstawiał wyników po standaryzacji danych. Zdecydowałem się również na troszeczkę inne przedstawienie danych osiągniętych przez ten algorytm, gdyż myślę że dzięki temu lepiej widoczne będzie jego zachowanie.



Widać zarówno po zachowaniu współczynnika FM jak i Randa, jak wraz ze wzrostem parametru zamglenia wyniki delikatnie spadają, jednakże nie są to znaczące wielkości. Algorytm dość małe wahania wyników w zależności od parametru zamglenia.

Sprawdźmy poraż kolejny jak przedstawia się rozkład wyników na poszczególnych zbiorach.



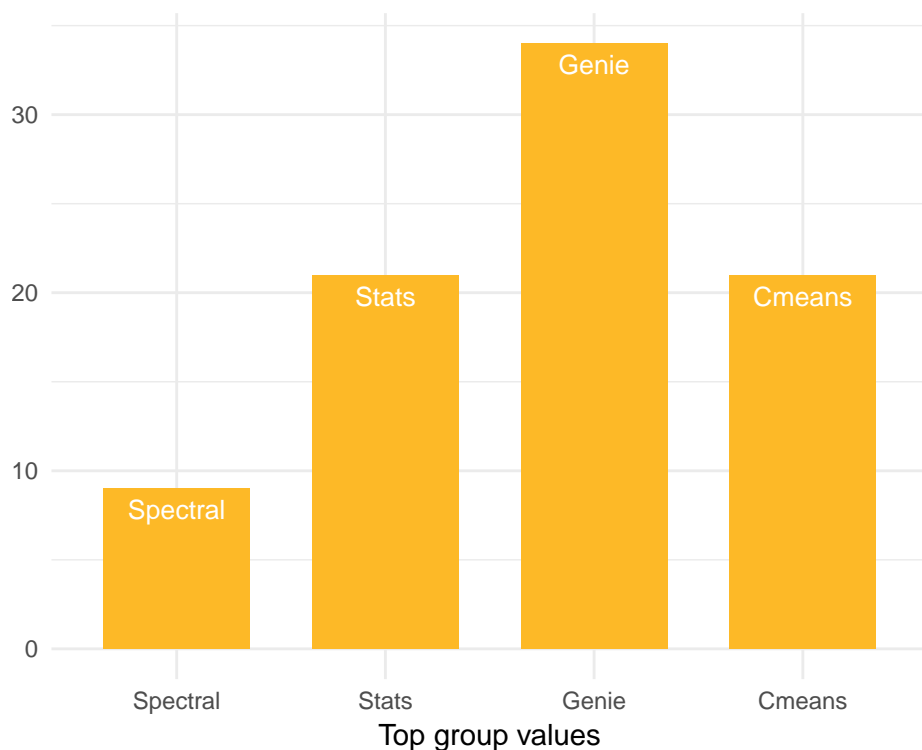
Algorytm ten prezentuje dość spory rozrzut, śmiałym stwierdzić że jeden z większych spośród testowanych przeze mnie. Jednocześnie jego średnia działania dość dobra. Sprawdźmy do jakich kubelków wpadają poszczególne zbiory, co myślę że daje jeden z bardziej miarodajnych obrazów jakości działania algorytmu./

Przedziały	Wartosci
0	1
(0,0.25]	8
(0.26,0.5]	7
(0.51,0.75]	9
(0.76,1]	21

Cmeans spisał się gorzej od algorytmów z pakietu stats, ale jednak lepiej od mojego algorytmu spektralnego. Warto zauważyć że cmeans działa bardzo szybko, co z pewnością można odnotować na jego plus. Jednakże dokładność jego wyliczeń pozostawia wiele do życzenia.

Wnioski

Sprawdźmy na jakiej ilości zbiorów ze wszystkich 46 testowanych przeze mnie każdy z algorytmów zdołał zdobyć wynik powyżej 0.75, co można będzie uznać za najlepszej jakości miernik sprawności działania tych wszystkich algorytmów.



Tak jak pisałem wcześniej genie zdecydowanie wygrywa swoją dokładnością oraz umiarkowaną predkością działania. Pozostałe algorytmy sprawują się dobrze, jednakże można mieć co do nich poważne zastrzeżenia co do sporego rozrzutu przez nie generowanego. Podsumowując, widać iż problem clusteringu nie jest problemem trywialnym. Algorytmy często mają różnego rodzaju problemy z odpowiednim pogrupowaniem danych. Jednakże, jak udało mi się pokazać w tej analizie istnieją takie algorytmy, które radzą sobie z tym problemem bardzo dobrze.