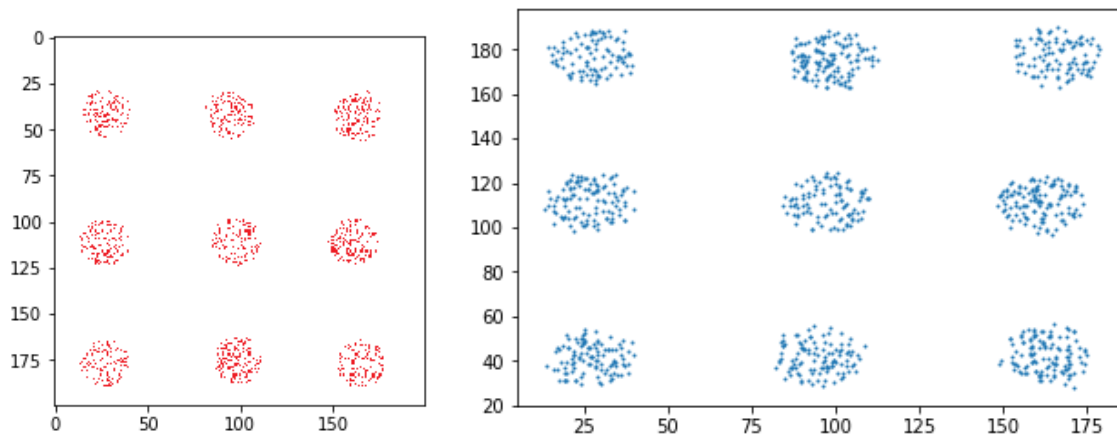
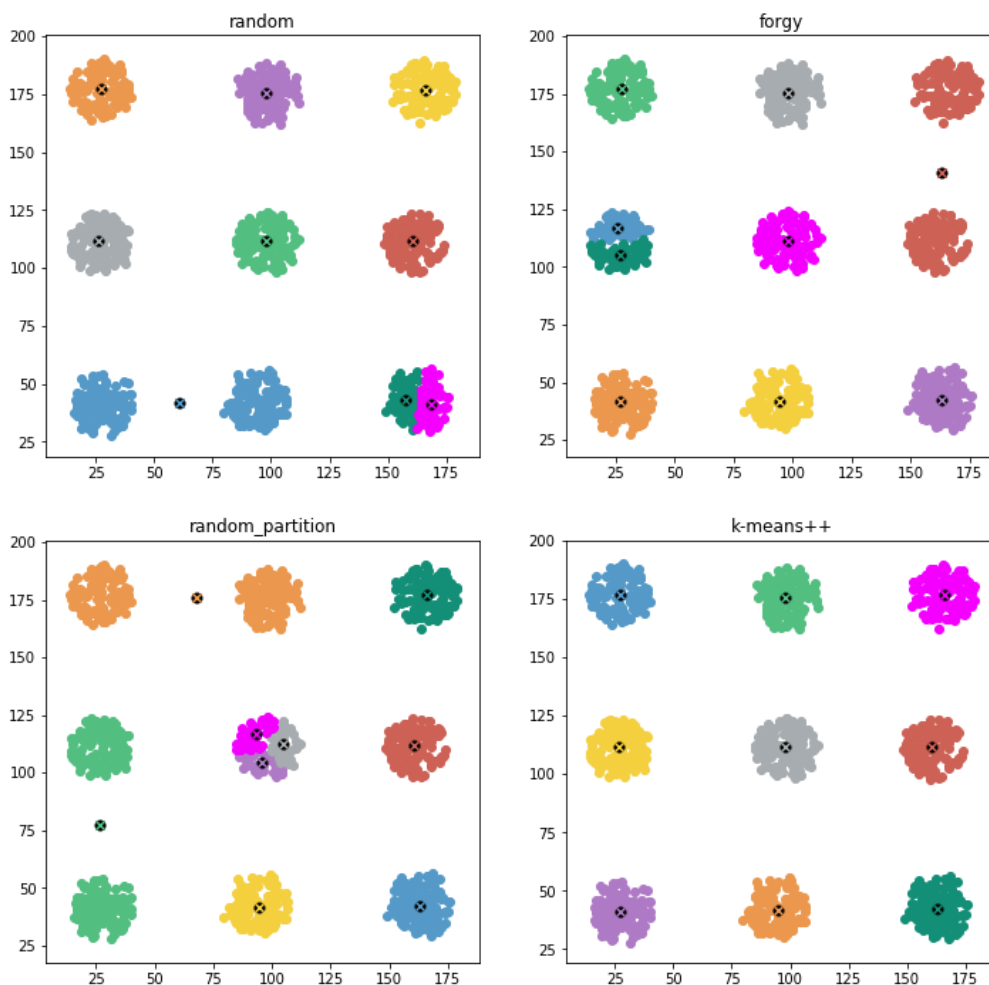


Raport – Kmeans initializations
Podstawy nauczania maszynowego
Wyk. Mateusz Woś

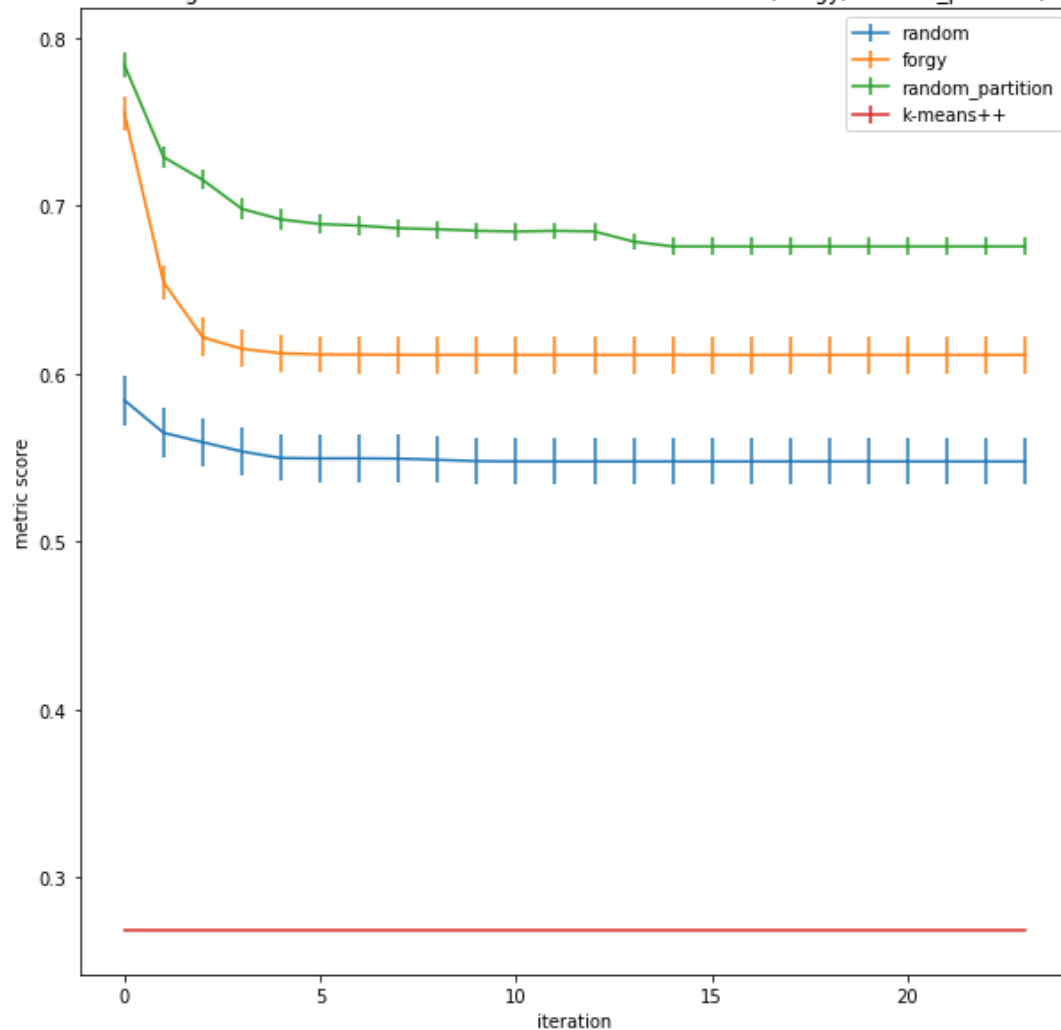
Dane do zadania ponownie wygenerowałem używając Painta. Zapisalem je do pliku aby w przyszłości łatwo z nich ponownie skorzystać.



Korzystając z gotowej implementacji k-means z scikit-learn zwizualizowałem klasteryzację dla każdej z metod inicjalizacji środków klastrów.



K-means score using Davies Bouldin score. Initialization methods: random, forgy, random_partition, k-means++

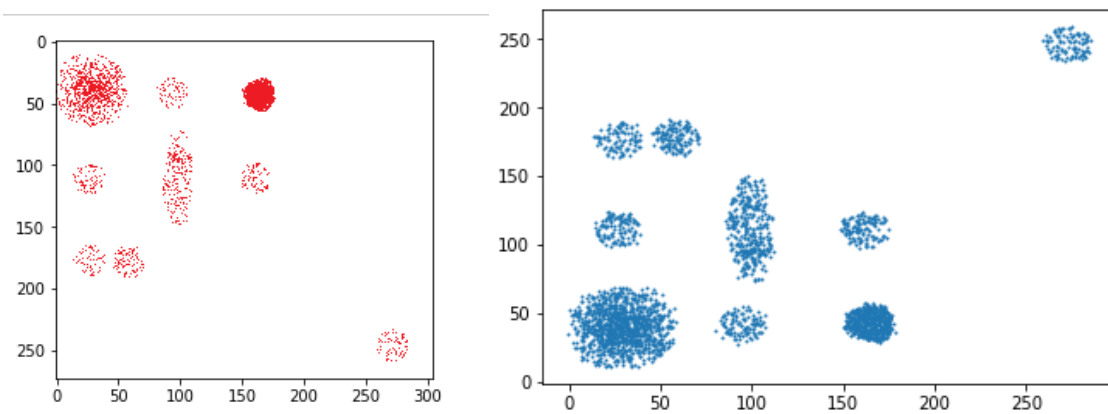


Korzystałem z metryki Daviesa Bouldina (mniejszy score -> lepiej)

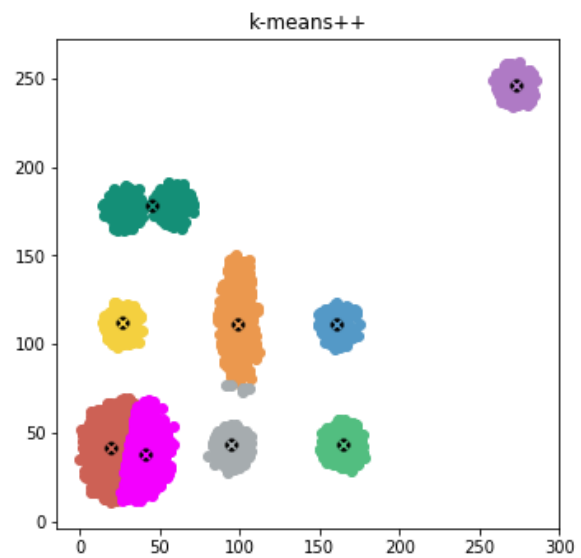
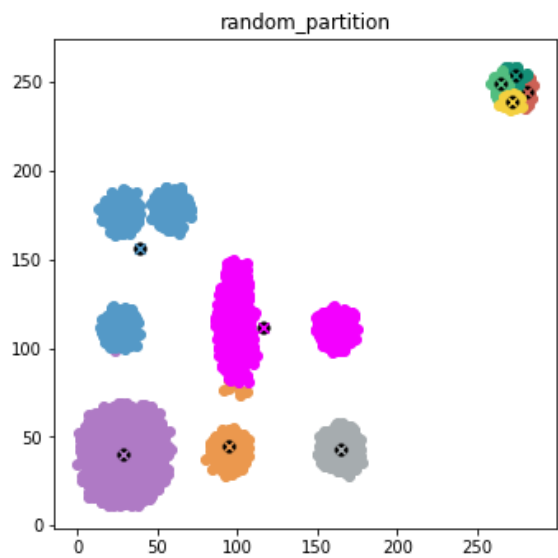
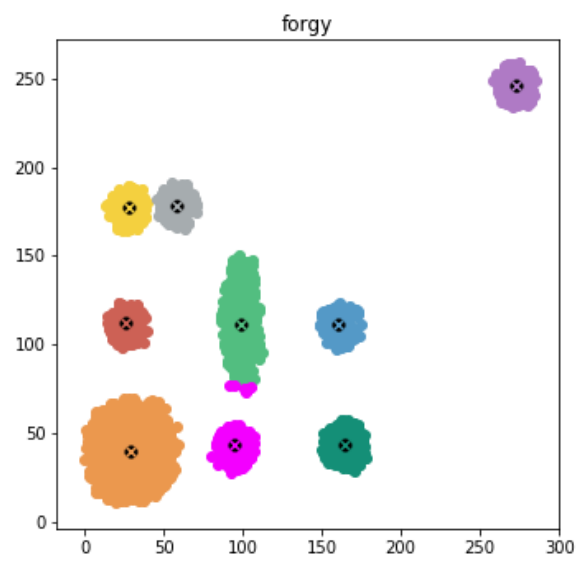
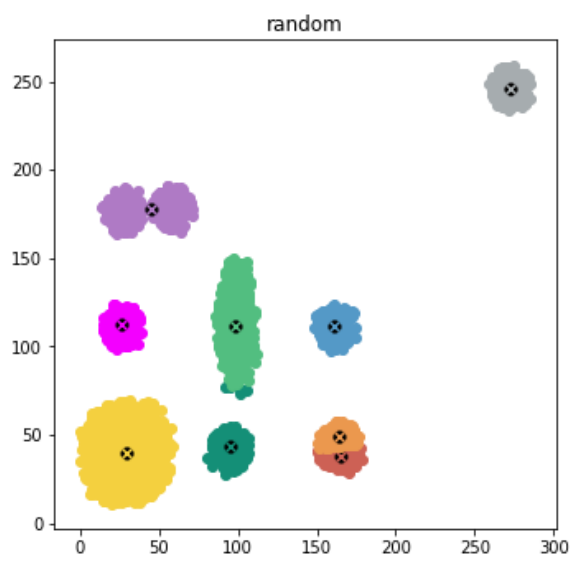
Random partition wypadł w tym przypadku najgorzej, lecz jak to bywa z losowością, czasami ta metoda bywała równie dobra jak czy kmeans++.

Kmeans++ jak widać powyżej jest bezkonkurencyjny. Praktycznie dla każdej ilości iteracji osiągał swój najlepszy wynik. Dla danego datasetu przy około 10 próbach zawsze idealnie wykonywał klasteryzację.

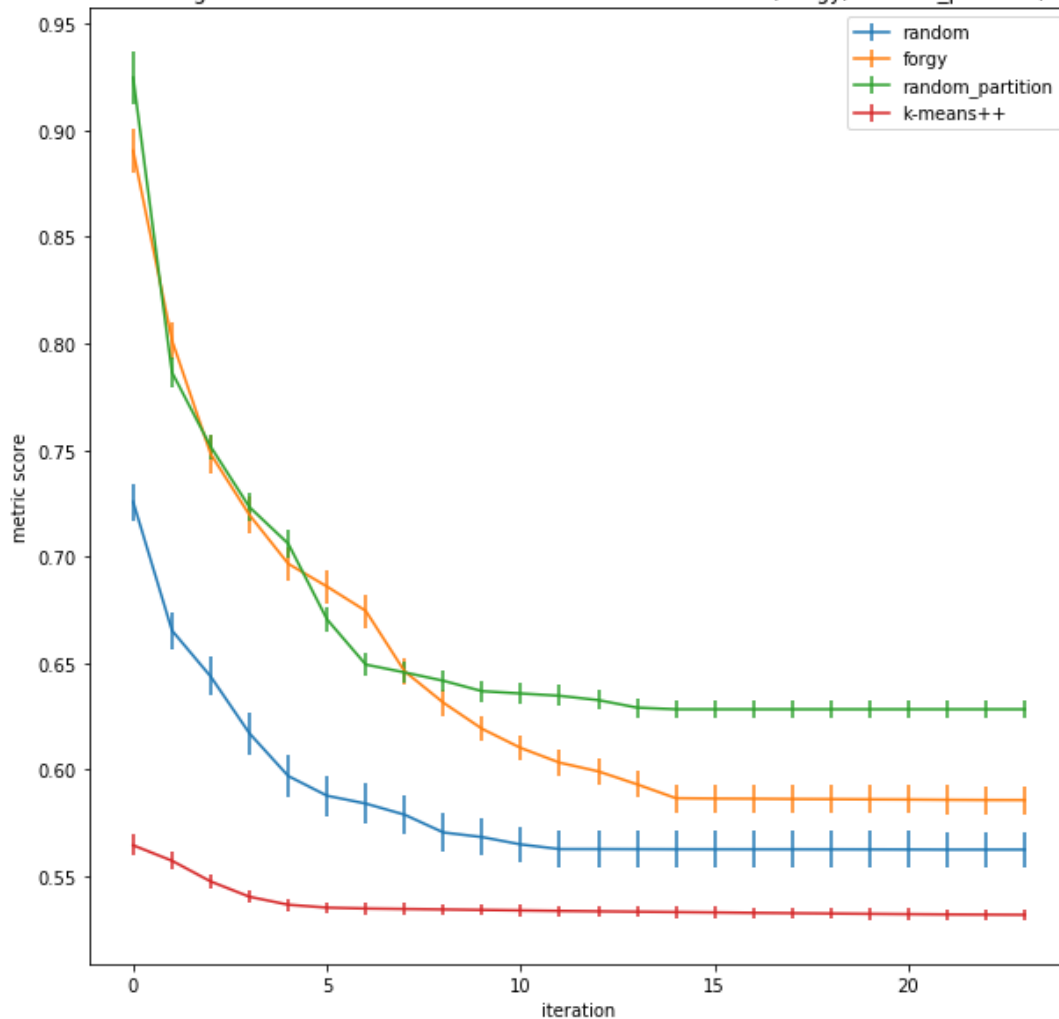
Wykresy dość szybko osiągają swój najlepszy wynik, powyżej około 5 iteracji każda z metod już nie stawała się lepsza. Dataset z którego korzystaliśmy był dość łatwy do klasteryzacji.



K-means clusterization depending on initialization method

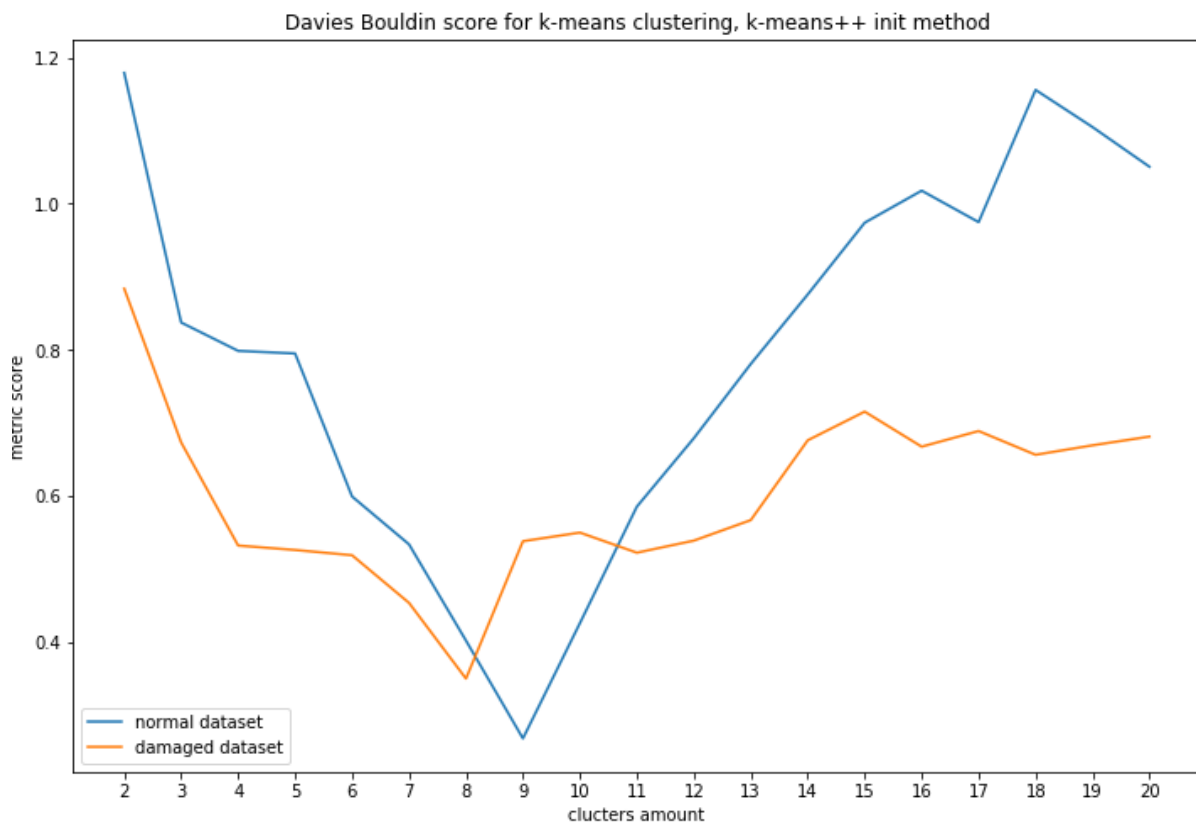


K-means score using Davies Bouldin score. Initialization methods: random, forgy, random_partition, k-means++



Dla zniekształconego datasetu sytuacja jest już ciekawsza. Widać tutaj, że metody nie spłaszczają się od razu. Metody random i random partition dopiero uzyskują swoje najlepsze możliwe wyniki przy około 20 iteracji na uruchomienie. Jak i poprzednio kmeans++ jest bezkonkurencyjny.

Dodatkowo ciekawe zależności można zaobserwować na wykresie z klastrami. Dobrze można zauważyć tutaj słabości algorytmu kmeans np. zbyt duży obszar i źle wylosowane centroidy -> kmeans uznaje taki klaster jako 2 oddzielne klastry; źle wylosowane punkty (patrz: random partition) i daleko oddalony kluster -> reszta klustrow nie brana pod uwagę bo do nich jest za daleko od początkowych centroidów



„Czy na ich podstawie można stwierdzić, że optymalne k to 9 (bo tyle mamy klastrów)?”

Tak, na powyższym wykresie widać, że najlepszy wynik dla normalnego dataset (w sumie dla uszkodzonego prawie też) jest dla $k=9$. Osobiście spodziewałem się większego, optymalnego k dla datasetu uszkodzonego.

Wnioski:

- Algorytm kmeans jest bardzo prosty i mało złożony, dzięki czemu nadaj się idealnie do klasteryzacji dużych datasetów.
- Niestety nie radzi sobie za bardzo z datasetami, w których nie ma wyraźnie wyróżnionych klastrów lub występują liczne anomalie.
- Musimy wiedzieć z góry ile klastrów spodziewamy się w danym datasetcie.
- Kmeans w zależności od metody inicjalizacji początkowych punktów, może dawać dość nieoczekiwane, niepoprawne wyniki. (dobrze to widać w metodzie random partition)