

Politechnika Śląska w Gliwicach
Wydział Automatyki, Elektroniki i Informatyki



Programowanie Komputerów

Metoda grupowania danych K-średnich

autor	Michał Skorus
prowadzący	dr inż. Tomasz Moroń
rok akademicki	2018/2019
kierunek	teleinformatyka
rodzaj studiów	SSI
semestr	3
termin zajęć projektowych	czwartek, 13:45 – 15:15
grupa	2
termin oddania sprawozdania	2019-01-25
data oddania sprawozdania	2019-01-25

1. Treść zadania

Stworzyć program grupujący dane według algorytmu k-średnich inaczej zwany k-means. Grupowanie może być wykonywane na wielu wymiarach i wyniki mają być zapisywane do pliku tekstowego. Dla dwóch wymiarów ma wygenerować również bitmapę graficznego wyniku.

2. Analiza zadania

Zagadnienie przedstawia problem tworzenia zastosowania algorytmu k-średnich. Aby go w pełni zrozumieć należy znać zagadnienia tj. odległość euklidesowska oraz średnia arytmetyczna. Odległość euklidesowska między dwoma punktami wyraża się wzorem:

$$d(A, B) = \sqrt{\sum_{i=1}^n ((x_{iA} - x_{iB})^2)}, \text{ gdzie } n\text{-liczba wymiarów; } x\text{-współrzędna danego punktu.}$$

Średnia arytmetyczna natomiast wyraża się wzorem:

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \text{ gdzie } x\text{-element zestawu liczb; } n\text{-ilość elementów.}$$

W celu ułatwienia zrozumienia, utrzymania większej uwagi i ułatwienia wyobrażenia sobie działania projektu, będę wymiennie nazywał dane owieczkami, a centroidy pasterzami.

2.1. Struktury danych

W programie wykorzystano kontener pary list zawierających obiekty danych wejściowych oraz obiekty wg których dane będą pogrupowane. Wykorzystywane kontenery pochodzą z biblioteki Standard Template Library (STL). Wybrałem strukturę list, gdyż algorytm nie sortuje żadnych elementów, ani nie potrzebuje szukać konkretnych elementów, tylko wykonuje się dla całego kontenera.

2.2. Algorytmy

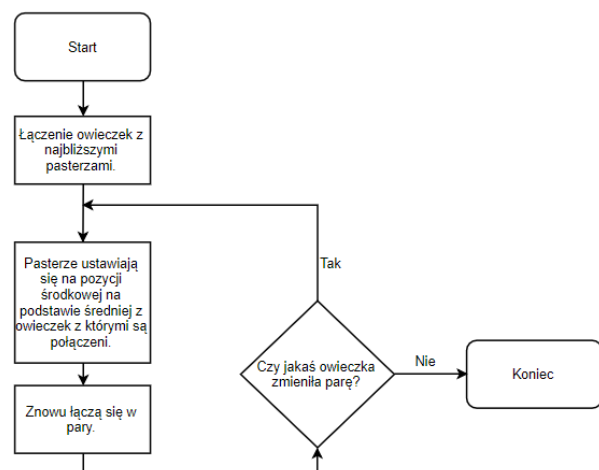
Algorytm k-średnich dzieli się na 3 etapy.

I ETAP – połączenie w pary danych z najbliższymi centroidami pod względem odległości euklidesowskiej.

II ETAP – centroidy ustawiają się na pozycji środkowej obliczonej ze średniej arytmetycznej danych z którymi są połączone.

III ETAP – powtarzanie poprzednich dwóch etapów do momentu w którym żaden element danych nie zmieni swojej pary.

Aby zoobrazować algorytm można wykorzystać schemat blokowy znajdujący się obok.



3. Specyfikacja zewnętrzna

Program uruchamiany jest z linii poleceń na dwa sposoby. Albo należy przekazać do programu nazwę pliku wejściowego po odpowiednim przełączniku "-i", np.

projekt.exe -i dane.txt

albo należy przekazać do programu przełącznik "-w" (without) np.

projekt.exe -w

co spowoduje wygenerowanie domyślnych ilości losowych danych.

Uruchomienie programu bez żadnego parametru lub z parametrem "-h"

program.exe

program.exe -h

powoduje wyświetlenie pomocy. Program można wywołać również z innymi przełącznikami takimi jak:

"-d" po którym należy podać wymiar. Jeżeli nie będzie podanego pliku wejściowego to w takim wymiarze będzie zrealizowana symulacja. Jeżeli będzie podany plik wejściowy program rozpozna wymiar. Minimalny wymiar to 1 a maksymalny to 499.

"-m" po którym należy podać ilość danych do pogrupowania. Jeżeli w pliku wejściowym znajdują się jedynie centroidy/pasterze to zostanie wygenerowane tyle danych ile się poda. Minimalna ilość danych to 1, a maksymalna to 49999.

"-k" po którym należy podać na ile grup mają zostać pogrupowane dane. Jeżeli w pliku wejściowym znajdują się jedynie dane/owce to zostanie wygenerowane tyle centroidów ile się poda. Minimalna ilość to 1, a maksymalna to 10000.

"-n" który spowoduje, że jeżeli program będzie realizowany w dwóch wymiarach to w wyniku przedstawionym za pomocą bitmapy danymi będą okręgi zamiast owiec.

Przykładowe sposoby użycia przełączników:

projekt.exe -d 25 -w -k 200 -m 10000

projekt.exe -n -k 30 -i owceplik.txt

Wszystkie parametry można podać w dowolnej kolejności natomiast muszą poprzedzać je przełączniki, aby miało to logiczny sens. Podanie większej ilości tych samych przełączników nie spowoduje zepsucia programu, program użyje argumentu podanego przy ostatnim przełączniku. Uruchomienie programu z wartościami wykraczającymi poza skalę wyświetli komunikaty:

Wymiar wykracza poza przedział [1;499]

Ilość podziału na grupy wykracza poza przedział [1;10000]

Ilość danych wykracza poza przedział [1;49999]

Wpisanie przełącznika -d, -m, -k i nie podanie po nim liczby zakończy program i wyświetli komunikat:

Nieprawidłowy parametr po <przełącznik>

Podanie nieprawidłowej nazwy pliku powoduje wyświetlenie odpowiedniego komunikatu:

Nie udało się otworzyć pliku

I wygenerowanie domyślnych ilości danych.

Plik wejściowy rozpoznaje obiekty dzięki nazwie. Jeżeli pierwszy wyraz linijki to:

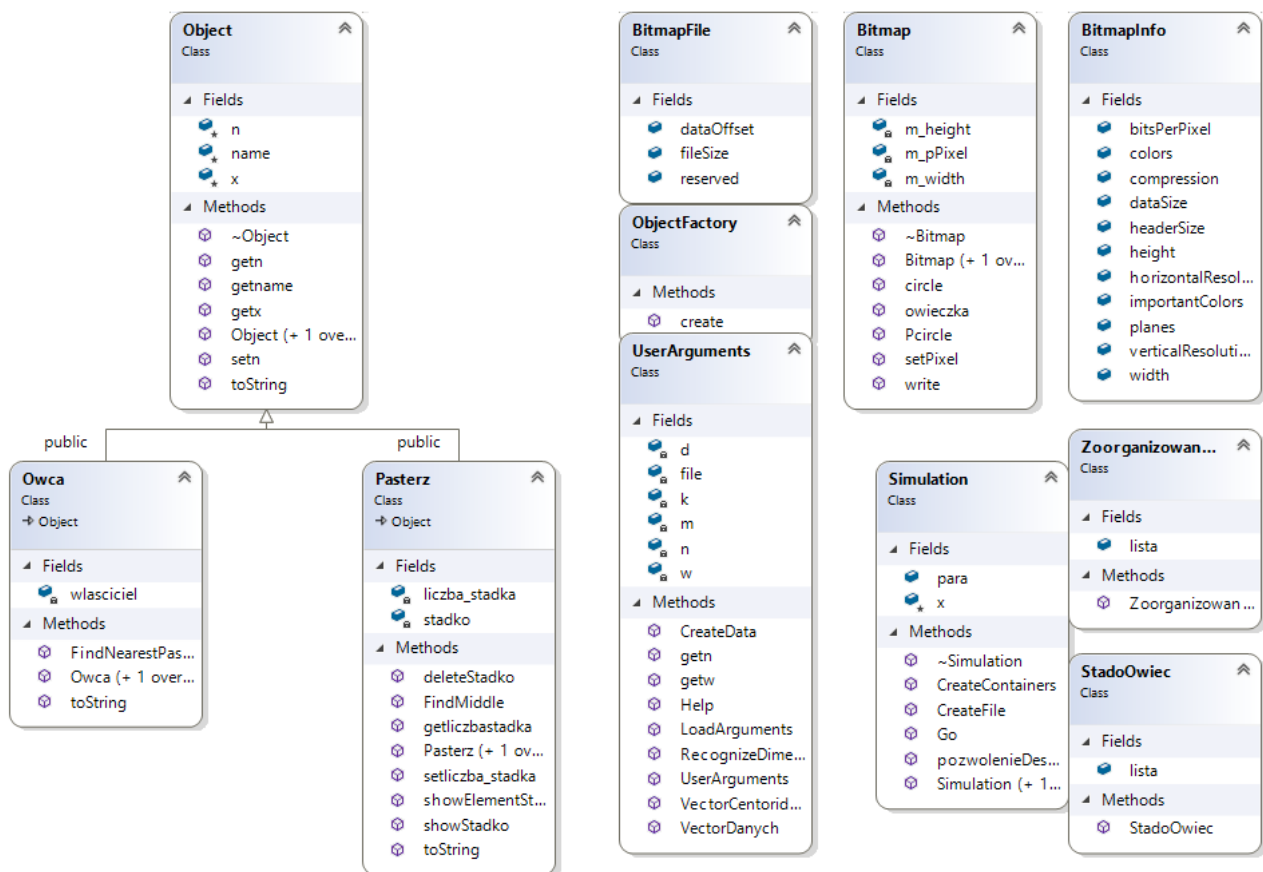
pasterz, centroid, centrum, pastor, Pasterz, PASTERZ, Centroid, CENTROID, Pastor, DobryPasterz, WladcaOwiec, KrolOwiec, owczarz, pastuszek

wtedy program uznaje linijkę jako centroidę/pasterza. Każdy inny pierwszy wyraz lub jego brak spowoduje uznanie linijki za owcę, pod warunkiem, że zgadza się z wymiarem. Program jest odporny na puste linijki, puste linijki ze znakami białymi, większą ilość spacji/tabulacji między danymi czy pusty plik wejściowy. Jeśli dane w pliku wykrócą poza przedział $[-10;10]$ lub zostanie podana mniejsza liczba współrzędnych niż wymiar, w którym realizowany jest program to zostanie wyświetlony komunikat o nie utworzeniu obiektu i podanie numeru linii, w której wystąpił problem, ale nie zatrzyma to działania programu.

4. Specyfikacja wewnętrzna

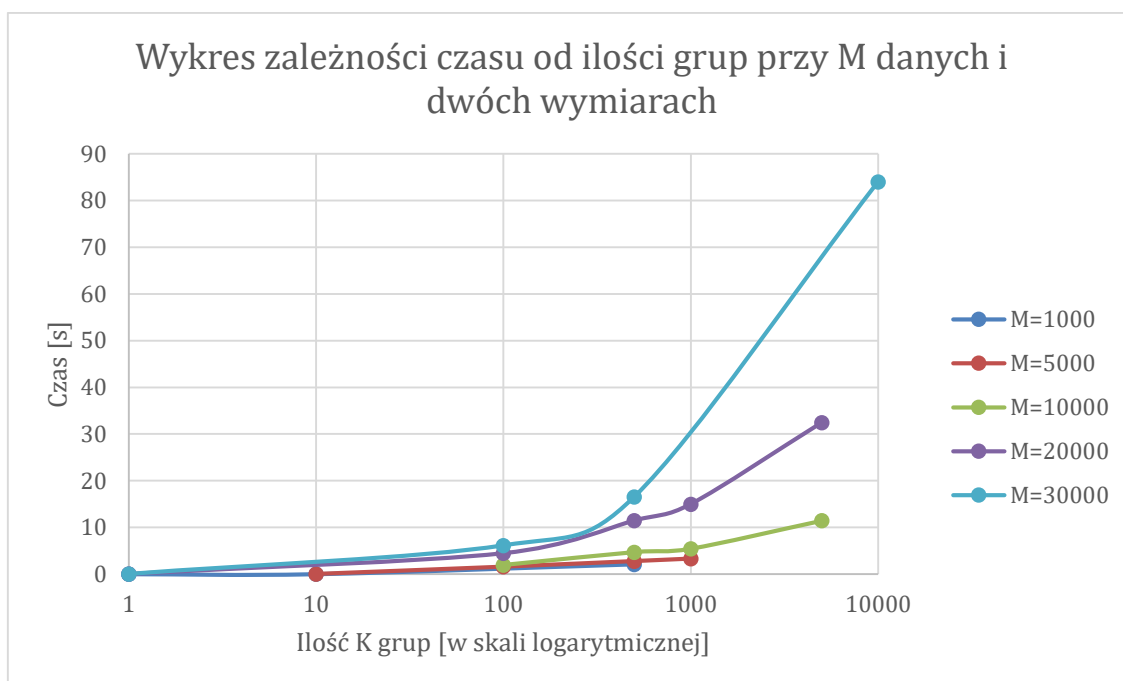
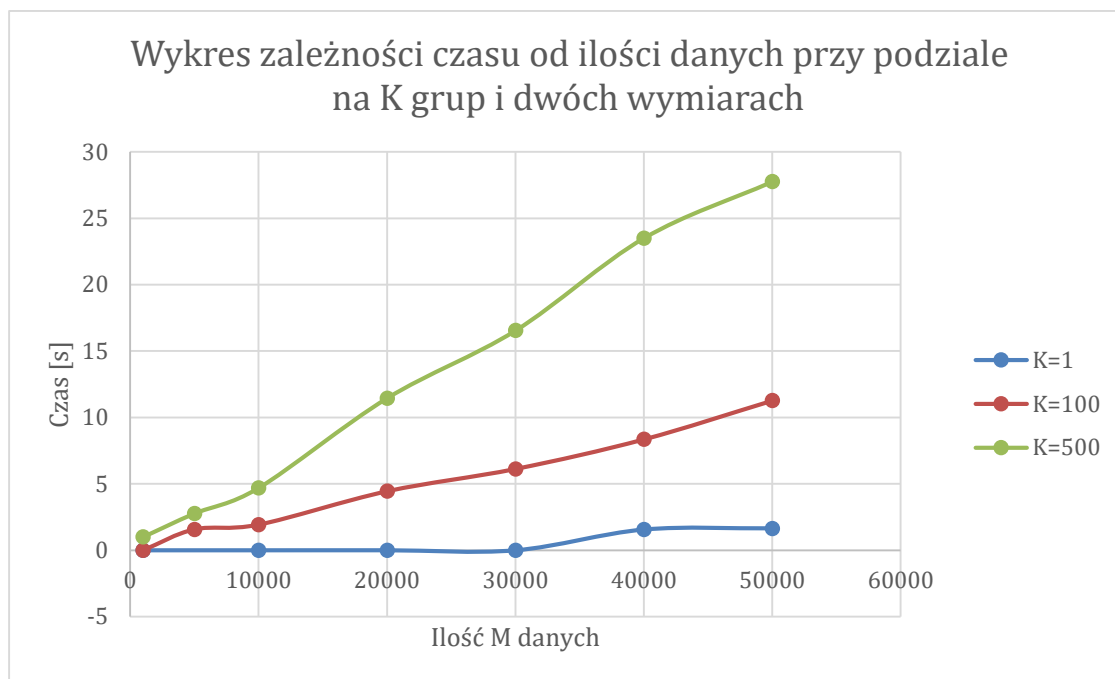
Program został zrealizowany zgodnie z paradygmatem obiektowym. Dokumentacja wewnętrzna wygenerowana Doxygenem jest dołączona w folderze sprawozdania w formacie pdf.

Diagram klas znajdujący się również w plikach na githubie prezentuje się w następujący sposób:



5. Efekty działania programu

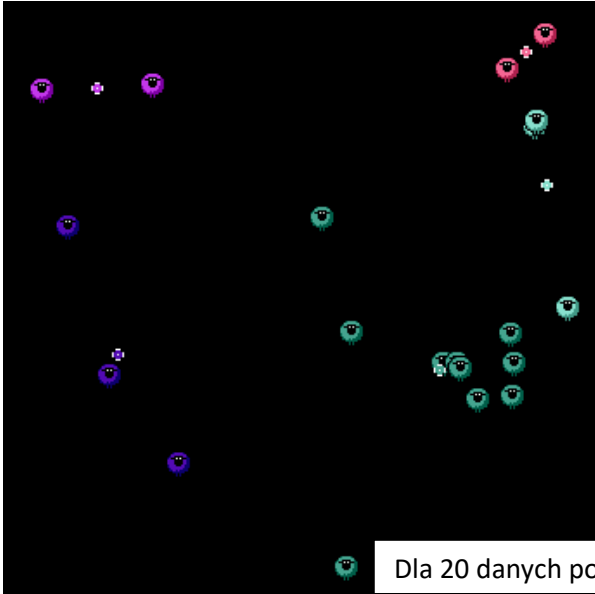
Program działa w optymalnym czasie co przedstawiają wykresy dla poszczególnych zestawów danych.



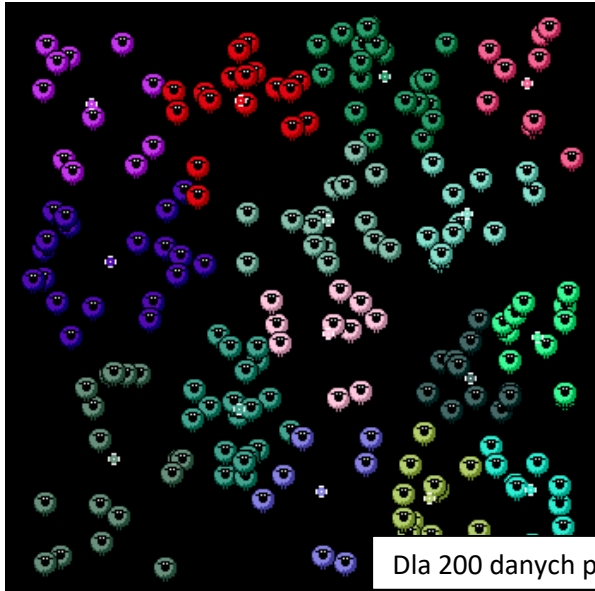
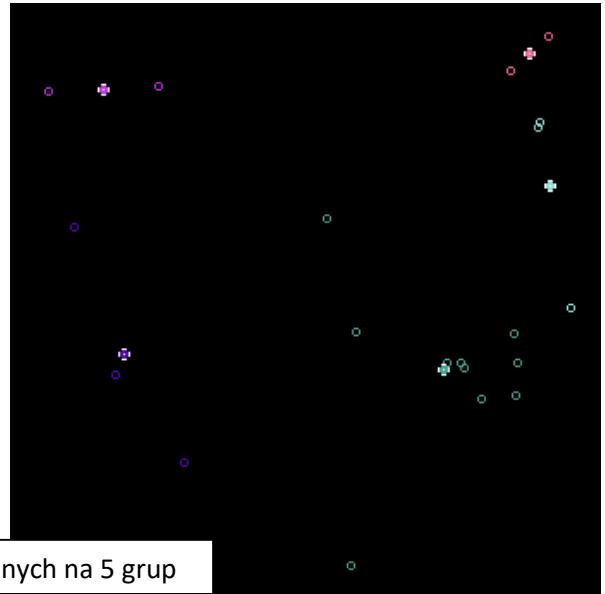
Platformą pomiarową był laptop z procesorem Intel Core i7 7700HQ, kartą graficzną GTX1050, posiadający 8GB RAMU

Bitmapa z wynikami wygląda dosyć oryginalnie z uwagi na grafikę. Każdej grupie losowany jest osobny kolor oraz wyświetlany jest w jej centrum centroid.

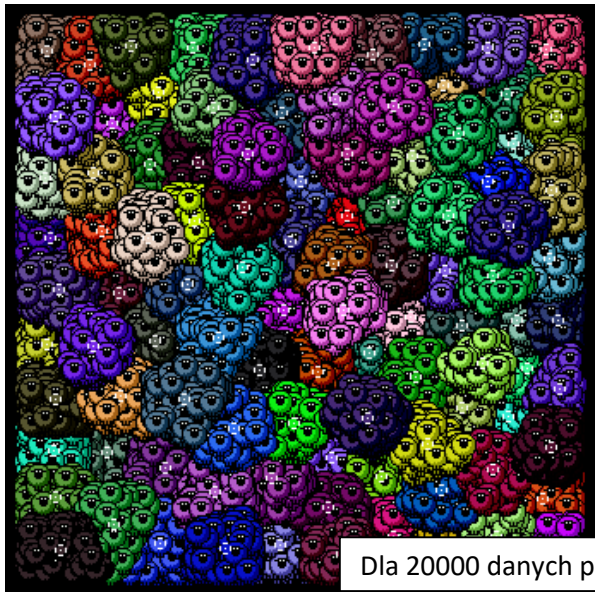
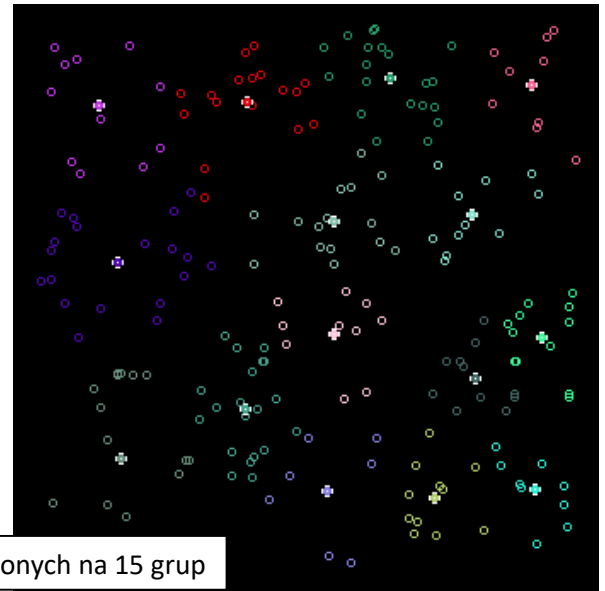
Przykładowe wyniki programu w postaci bitmapy z lewej bez przełącznika -n, z prawej z przełącznikiem:



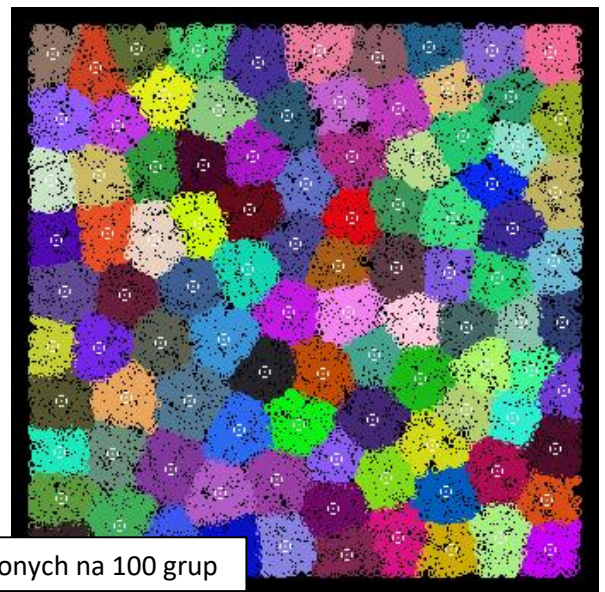
Dla 20 danych podzielonych na 5 grup



Dla 200 danych podzielonych na 15 grup



Dla 20000 danych podzielonych na 100 grup



Wyniki przedstawiane są zgodnie z układem kartezjańskim, czyli w lewym dolnym rogu są wartości $(-10;-10)$, w prawym górnym $(10;10)$, a w środku $(0;0)$.

Wyniki dla każdego wymiaru prezentują się w postaci pliku tekstowego o nazwie "wyniki.txt". Oto przykładowe pliki wynikowe: (Plik tekstowy wyświetla dane {dane są oddzielone enterami, natomiast współrzędne konkretnego wymiaru spacjami}, konsola pozycję centroidów)

```
wynik.txt — Notatnik
Plik Edycja Format Widok Pom
Grupa 1
-9.35 7.52 8.44 6.2
-8.4 2.6 7.48 8.26
Grupa 2
1.86 -1.24 8.85 9.54
5.15 -2.32 7.68 -3.53
6.41 -3.64 7.72 -2.32
7.6 -1.28 1.66 -9.74
5.68 -2.36 -4.41 -5.99
0.77 2.9 -5.34 7.7
9.67 -0.37 8.52 6.38
5.79 -2.51 -6.92 -2.76
```

```
C:\Users\Majkel\Desktop\Proj
Pasterz zawany pasterz
Wymiar = 4
-8.88 5.06 7.96 7.23
Liczba stadka: 2
Pasterz zawany pasterz
Wymiar = 4
5.37 -1.35 2.22 -0.09
Liczba stadka: 8
```

Dla 10 danych podzielonych na 2 grupy w 4 wymiarach.

```
wynik.txt — Notatnik
Plik Edycja Format Widok Pomoc
Grupa 1
-9.35 7.52 8.44 6.2 1.86
-3.53 6.41 -3.64 7.72 -2.32
-5.34 7.7 9.67 -0.37 8.52
-4.01 9.48 -7.59 -0.83 -1.02
-6.93 -2.84 9.79 0.18 -5.36
Grupa 2
-1.24 8.85 9.54 -8.4 2.6
6.38 5.79 -2.51 -6.92 -2.76
2.62 -5.14 -4.26 -5.94 -7.53
Grupa 3
7.48 8.26 5.15 -2.32 7.68
7.6 -1.28 1.66 -9.74 5.68
-2.36 -4.41 -5.99 0.77 2.9
-0.2 4.03 -6.15 4.42 7.87
```

```
C:\Users\Majkel\Desktop\Projekt\x64\Ro
Pasterz zawany pasterz
Wymiar = 5
-5.83 5.65 3.33 2.58 0.34
Liczba stadka: 5
Pasterz zawany pasterz
Wymiar = 5
2.59 3.17 0.92 -7.09 -2.56
Liczba stadka: 3
Pasterz zawany pasterz
Wymiar = 5
3.13 1.65 -1.33 -1.72 6.03
Liczba stadka: 4
```

Dla 12 danych podzielonych na 3 grupy w 5 wymiarach.

6. Testowanie

Program został przetestowany na różnego rodzaju plikach. Pliki niepoprawne (puste, z większą ilością znaków białych, z pustymi linijkami, z danymi wykraczającymi poza skalę) powodują zgłoszenie błędu lub wygenerowanie własnych danych wejściowych. Złe parametry wywołania wyświetlają błąd. Błąd wyświetla się również gdy plik nie zostanie otwarty. Program został sprawdzony pod kątem wycieków pamięci. Program jest zabezpieczony pod szerokim kątem.

7. Wnioski

Problem zaimplementowania algorytmu k-średnich choć z pozoru wydawał się prosty to w zastosowaniu przyspożył mi sporo zagwostek. Spowodowane było to głównie tym, że biblioteka STL i paradygmat obiektowy to dla mnie coś nowego dlatego potrzebowałem dłuższej chwili czasu na rozwiązanie pewnych problemów. Największym problemem okazało się dodanie do listy obiektu, który był w stworzony fabryką generyczną. Element ten miał typ klasy bazowej i należało użyć konwersji `dynamic_cast` aby problem znikł. Mimo wielu potyczek z debuggerem myślę, że w ogromnym stopniu ulepszyłem swoje umiejętności programowania oraz zasiliłem swoją wiedzę wzorcem obiektowym. Taki model tworzenia aplikacji pozwolił mi spojrzeć na programowanie z zupełnie innej strony, ułatwiającej tworzenie większych projektów. Udało mi się zrealizować wszystkie założenia początkowe co tylko dowodzi osiągnięcia przeze mnie nowych kwalifikacji.