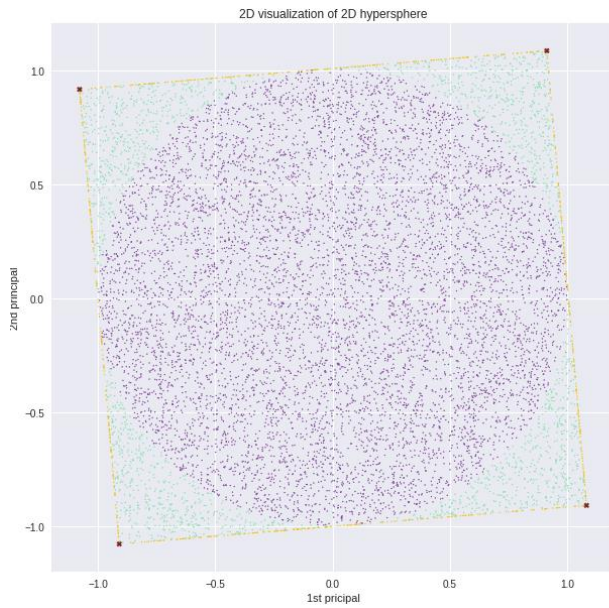


Raport – Condensed Nearest Neighbours
Podstawy nauczania maszynowego
Wyk. Mateusz Woś

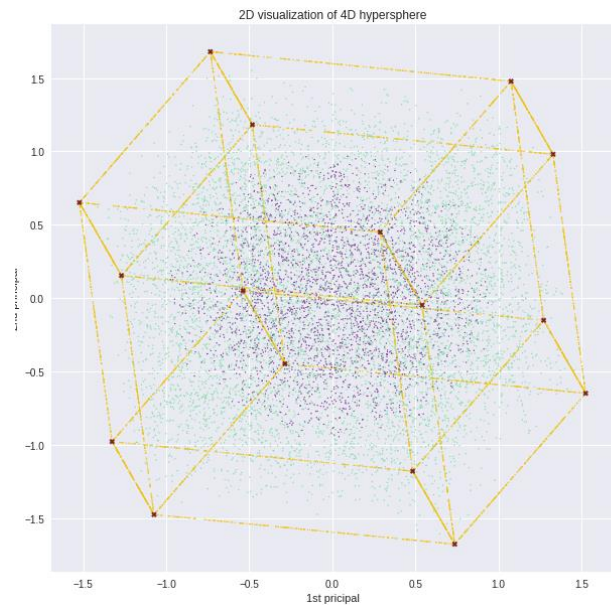
Zadanie a)

Wygenerowane dane składają się z 10000 punktów w środku hipersześcianu i dodatkowych 100 punktów na pojedynczą krawędź.

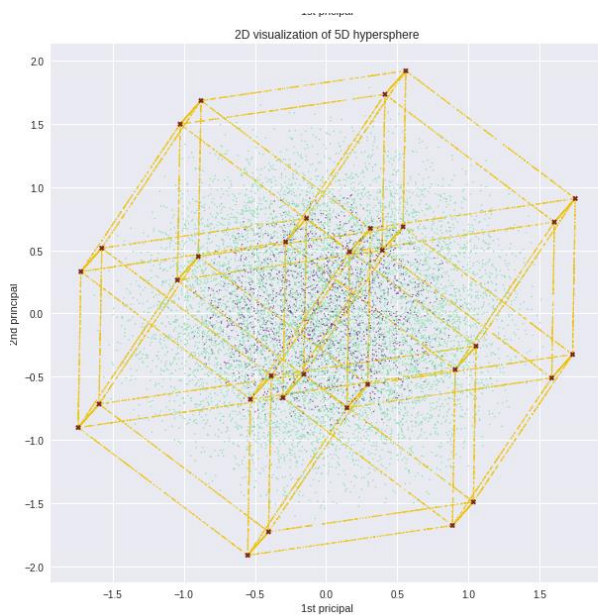
Data colors:{"corners" -> "dark red", "edges" -> "yellow", "inside" -> "purple", "outside" -> "green"}



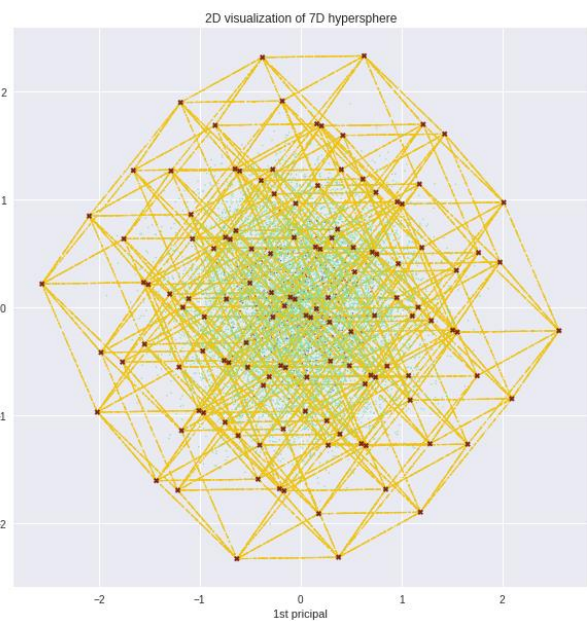
2D -> 2D



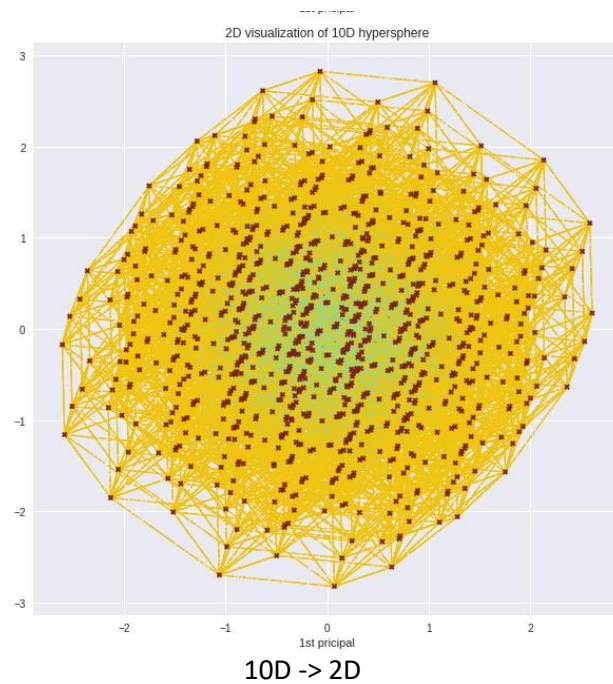
4D -> 2D



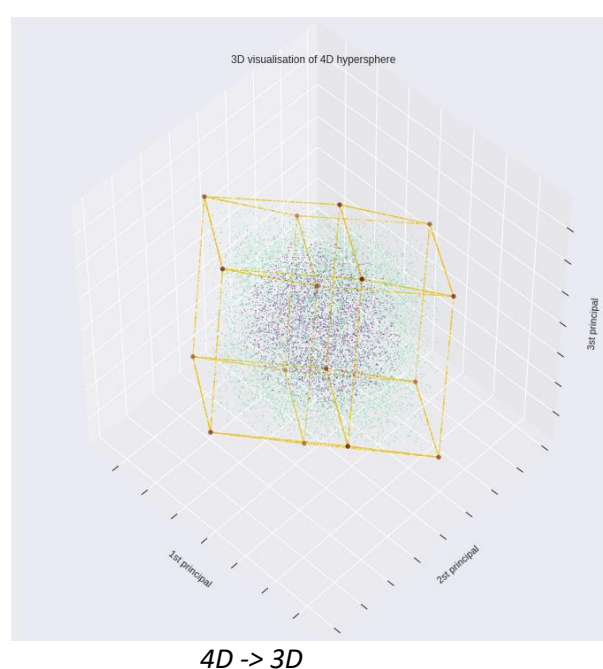
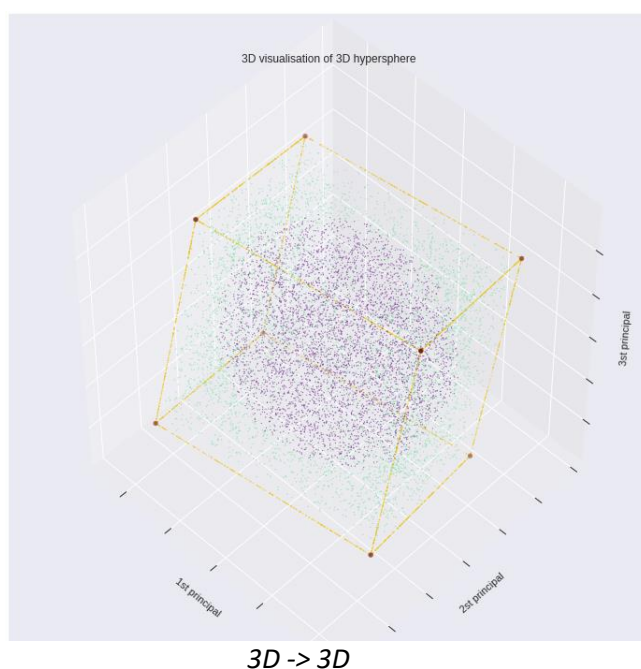
5D -> 2D

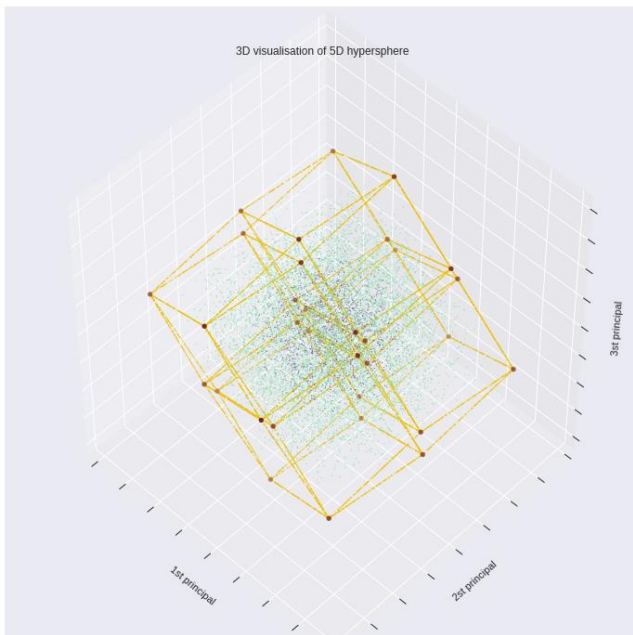


7D -> 2D

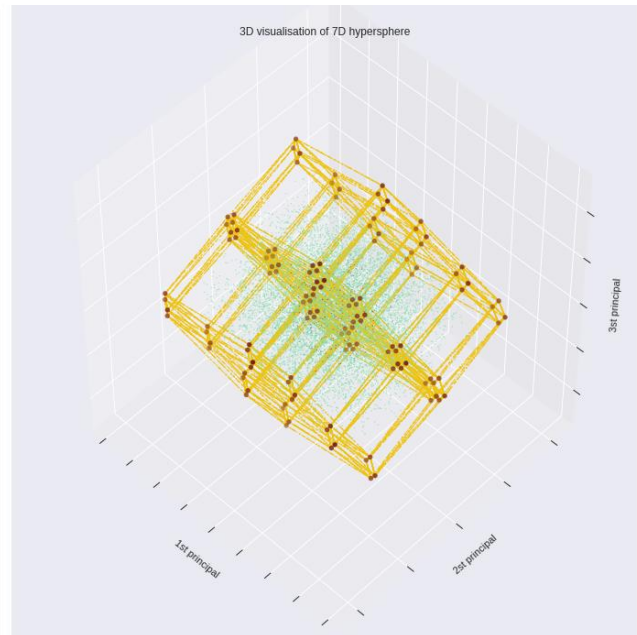


Niestety ze względu na niesamowicie długie czasy transformacji datasetu przez PCA z 13D do 2D byłem zmuszony na delikatną zmianę zadania i wizualizację rzutu 10D hipersześcianu na 2D płaszczyznę. Tak samo postąpiłem w wizualizacji rzutów dla 3D – zastąpiłem ostatni podpunkt wersją dla 10D.

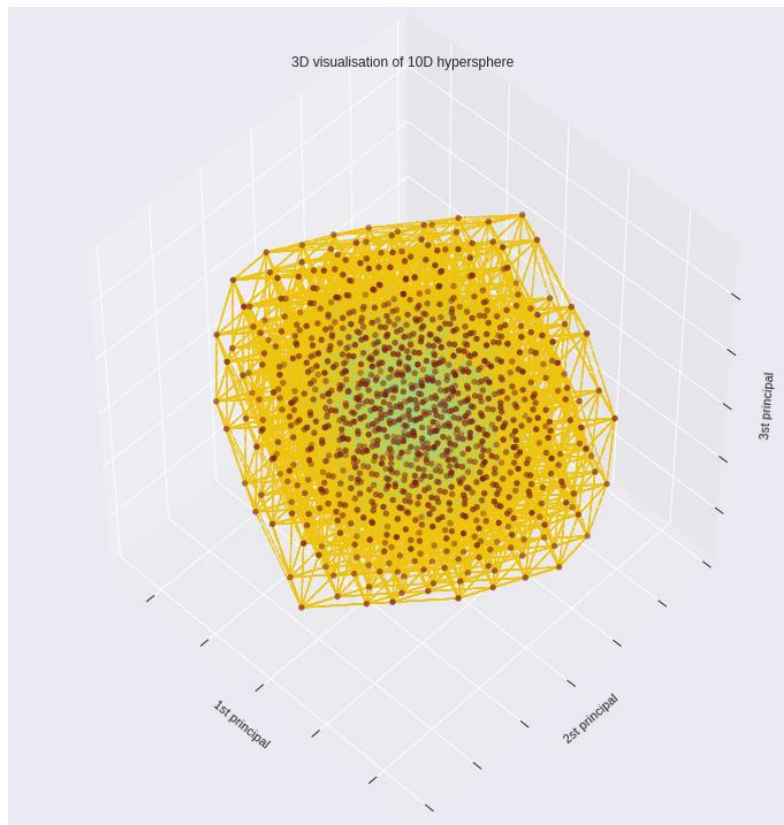




5D -> 3D

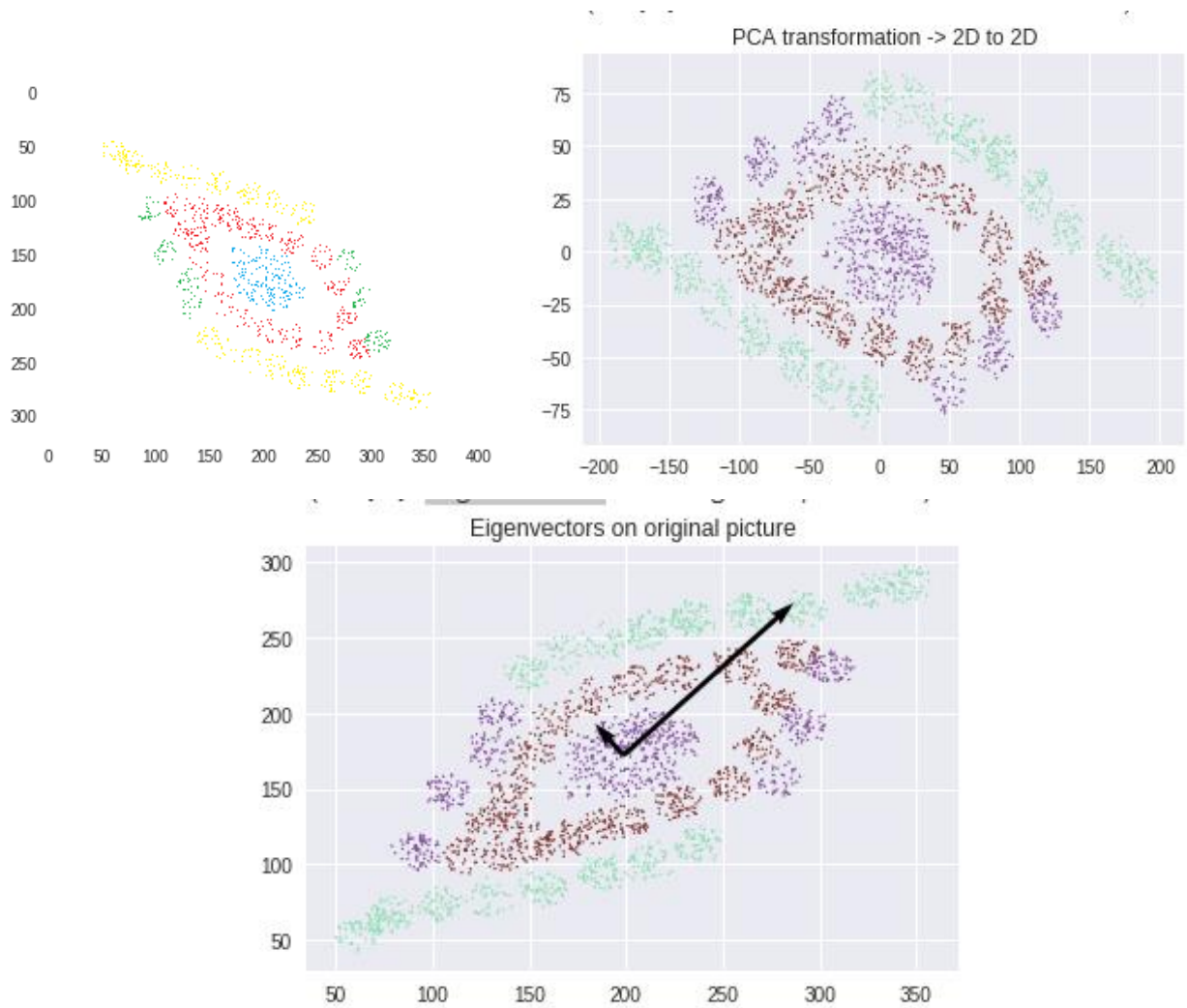


7D -> 3D

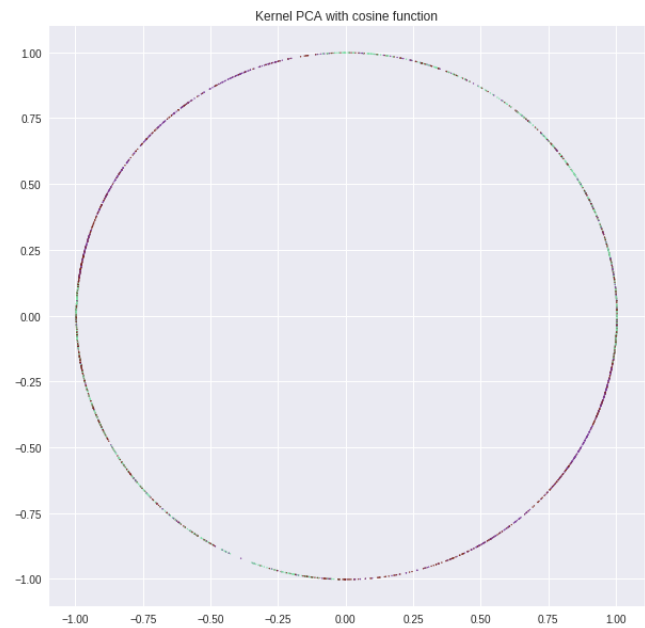
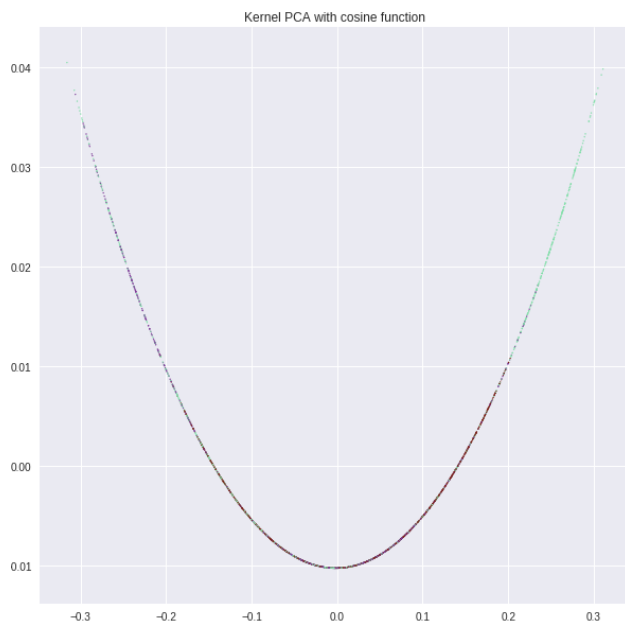


10D -> 3D

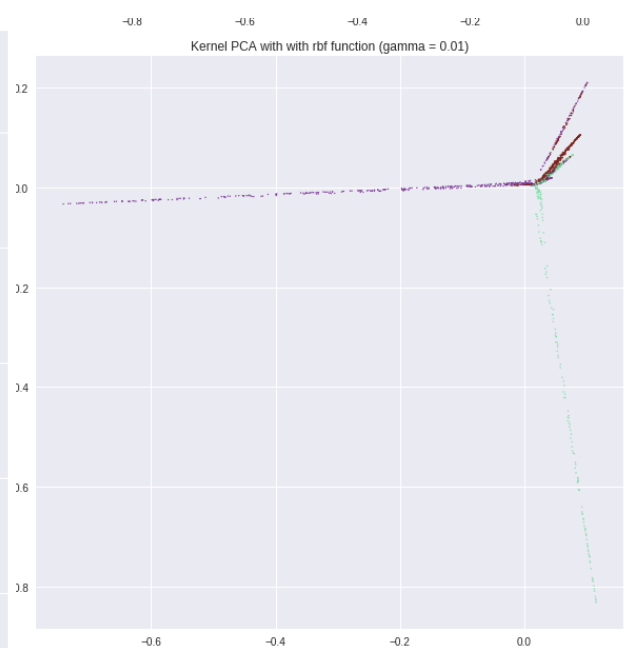
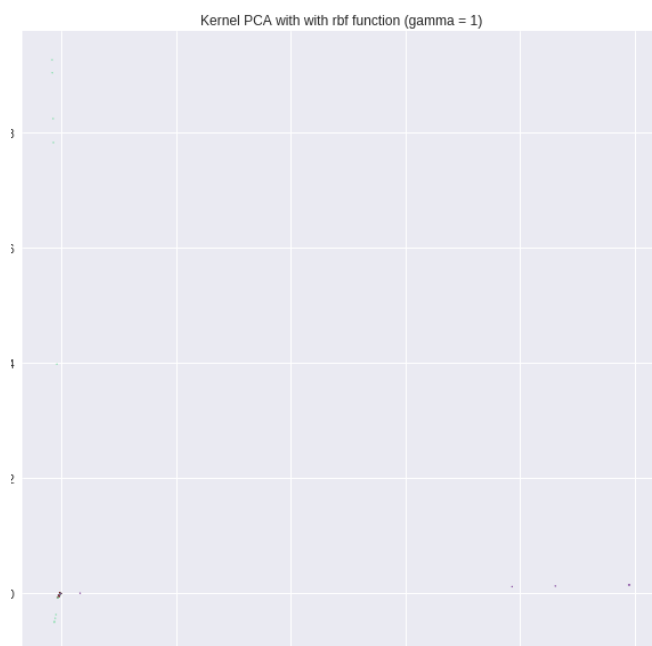
Zadanie b)

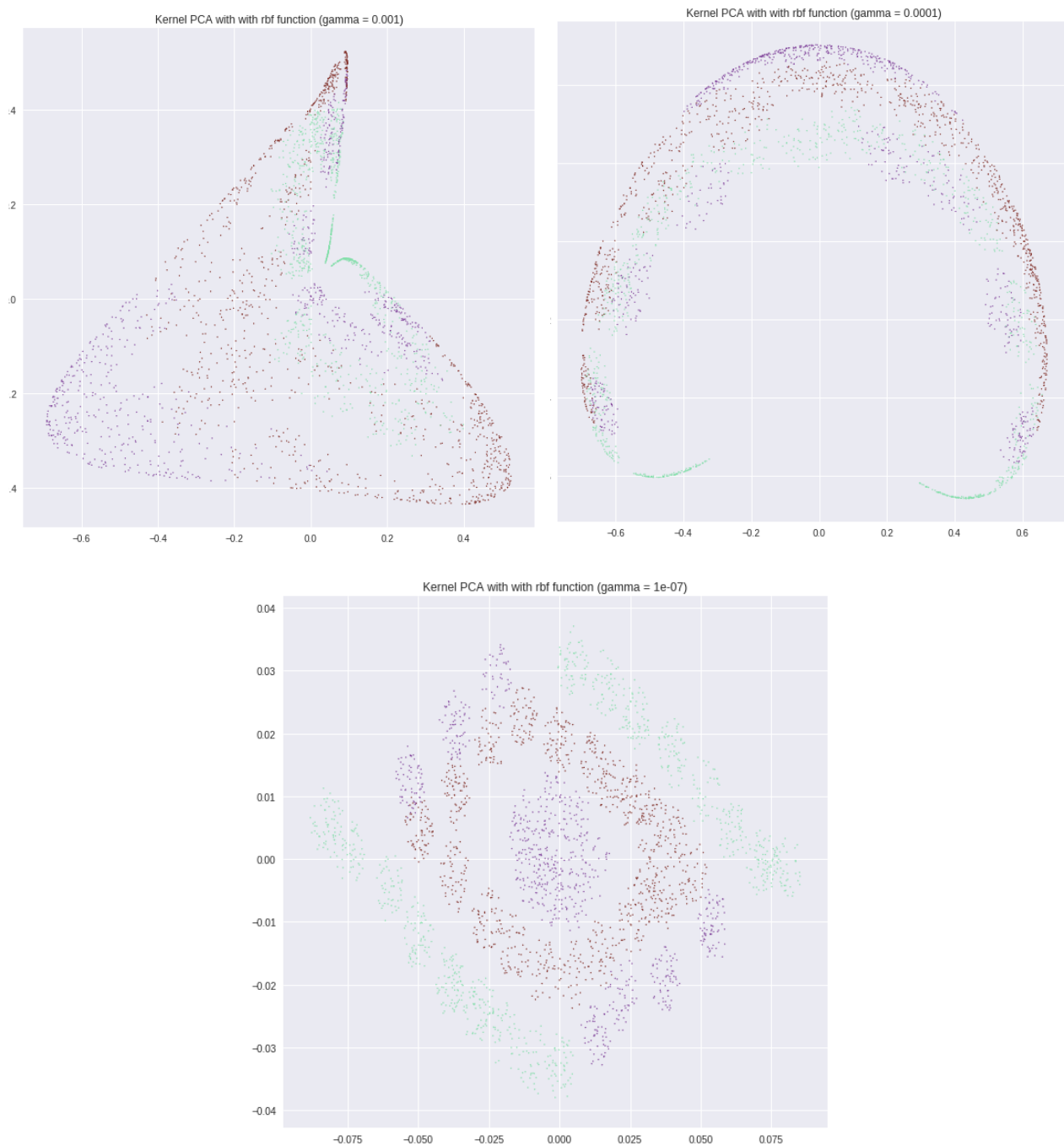


Po przejściu z 2D do 2D widać dokładnie, że dane są scentralizowane wokół punktu (0, 0)



Powyżej widać Cosine karnel PCA przed wyśrodkowaniem danych i po.





Powyżej widać Radial Basis Function kernel PCA z różnymi parametrami gamma

Wnioski:

- Jak można zauważyć na wykresach (jak i zadania z 1 cwiczen) ukazujących rzut hipersześcianów na płaszczyzny 2 i 3 wymiarowe prawdopodobieństwa pojawienia się punktów w środku hiperkuli maleje znacząco ze wzrostem wymiarów.
- Wraz ze wzrostem wymiarów rzuty stawały się coraz bardziej nieczytelne. Rzut stawał się praktycznie zbiorem samych krawędzi.
- Ograniczenie się do 2/3 głównych składowych pozwala nam na usunięcie mniej znaczących składowych, dzięki czemu pozbywamy się cech które mogłyby naruszyć poprawność działania naszego modelu, a także dzięki zmniejszeniu ilości wymiaru samo przetwarzanie byłoby szybsze.

- Przy zastosowaniu PCA z 2D do 2D widzimy, że baza została zmieniona. Obraz został skoncentrowany wokół punktu (0, 0)
- Przy użyciu PCA z kernel cosine można zauważyć, iż rozstawienie punktów jest zależne od kąta pod którym one leżą od punktu (0, 0) w pierwotnym datasetcie. PCA z kernel cosine dla datasetu skoncentrowanego w (0, 0) wygląda jak okrąg, przed wysrodkowaniem jak hiperbola.
- PCA z kernel RBF dla bardzo małych wartości γ działa podobnie do PCA.
- PCA z kernel RBF dla współczynnika $\gamma = 1$ dla tych danych na wykresie widać bardzo niewiele, wraz z maleniem współczynnika widać separację kolorów. Czym większy współczynnik tym słabsza separacja.