

Klasteryzacja danych za pomocą sieci Kohonena

Aleksander Kłak, prowadzący Dr Marek Bazan

18 czerwca 2023

Spis treści

1	Cel ćwiczenia	1
2	Algorytm Kohonena	1
3	Miary podobieństwa	2
4	Wyniki	2
4.1	Dane Iris	2
4.1.1	Liczba klas: 2	2
4.1.2	Liczba klas: 3	7
4.2	Dane giełdowe dla trzech spółek	11
4.2.1	Spółka AOS	11
4.2.2	Spółka AMP	12
4.2.3	Spółka APD	12
5	Wnioski	12
6	Kod	13
6.1	Sposoby mierzenia odległości między danymi a reprezentantami	13
6.2	Algorytm Kohonena	14
6.3	Skrypt	15

1 Cel ćwiczenia

Zadaniem jest zaimplementowanie algorytmu Kohonena w celu grupowania danych (klasteryzacji) oraz sprawdzenia jego działania na dwóch zestawach danych: zestawie danych Iris oraz zestawie danych giełdowych

2 Algorytm Kohonena

Algorytm Kohonena jest jednym z algorytmów samoorganizujących map i służy do grupowania danych na podstawie ich podobieństwa. W kontekście tego zadania należy zastosować algorytm Kohonena, aby znaleźć klastry (grupy) danych w zbiorach danych Iris oraz danych giełdowych.

3 Miary podobieństwa

Iloczyn skalarny: Miara podobieństwa, która oblicza iloczyn skalarny dwóch wektorów. Iloczyn skalarny między dwoma wektorami oblicza się jako sumę iloczynów odpowiadających sobie elementów obu wektorów. Iloczyn skalarny jest większy, gdy wektory są bardziej skierowane w podobnych kierunkach, a mniejszy, gdy wektory są bardziej skierowane w przeciwnych kierunkach.

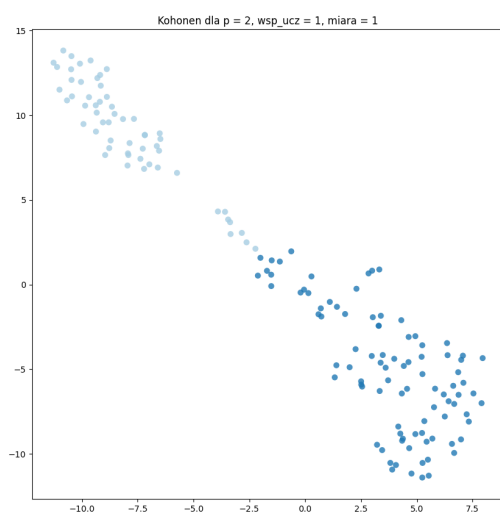
Odległość euklidesowa: Miara podobieństwa, która oblicza odległość pomiędzy dwoma punktami w przestrzeni. Odległość euklidesowa między dwoma wektorami oblicza się jako pierwiastek kwadratowy sumy kwadratów różnic między odpowiadającymi sobie elementami obu wektorów. Odległość euklidesowa jest niższa, gdy wektory są bardziej podobne i mają bliższe wartości elementów.

Odległość manhattańska oblicza sumę wartości bezwzględnych różnic między odpowiadającymi sobie elementami dwóch wektorów. Ta miara podobieństwa jest często używana w przypadkach, gdy istotne jest porównanie względnych różnic między wartościami elementów, a nie ich rzeczywiste wartości.

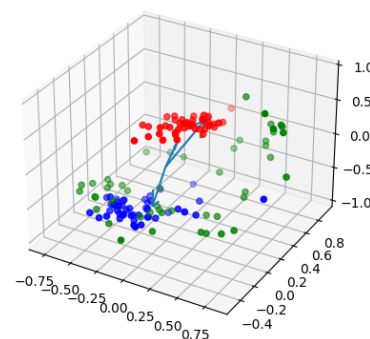
4 Wyniki

4.1 Dane Iris

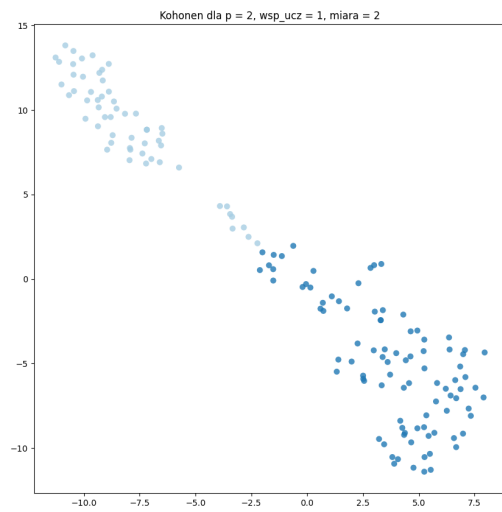
4.1.1 Liczba klas: 2



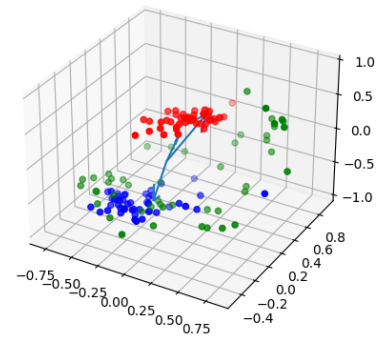
Wizualizacja3D t-SNE dla p = 2, wsp_ucz = 1, miara = 1



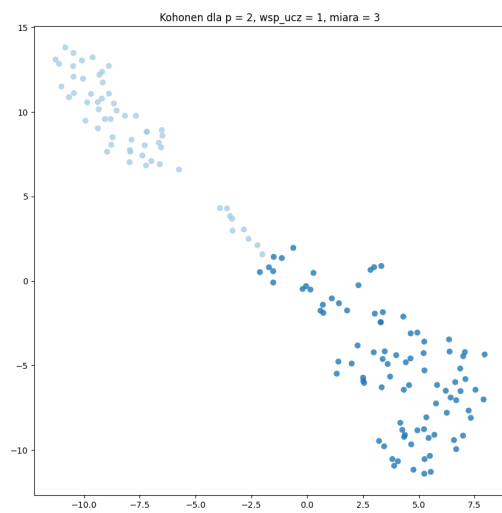
Rysunek 1: Współczynnik uczenia liniowo malejący (1), metoda miary podobieństwa iloczyn skalarny



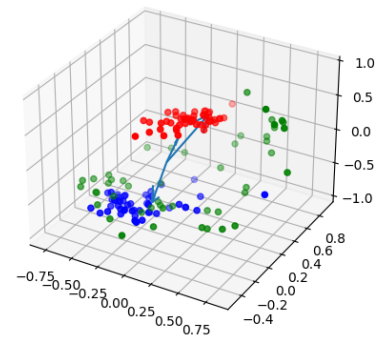
Wizualizacja3D t-SNE dla $p = 2$, $wsp_ucz = 1$, $miara = 2$



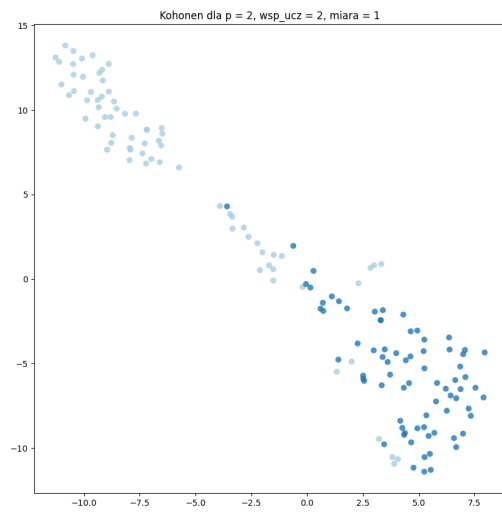
Rysunek 2: Współczynnik uczenia liniowo malejący (1), metoda miary podobieństwa odległość euklidesowa



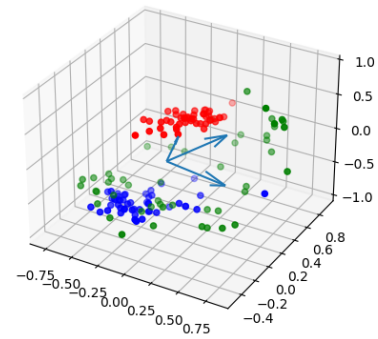
Wizualizacja3D t-SNE dla $p = 2$, $wsp_ucz = 1$, $miara = 3$



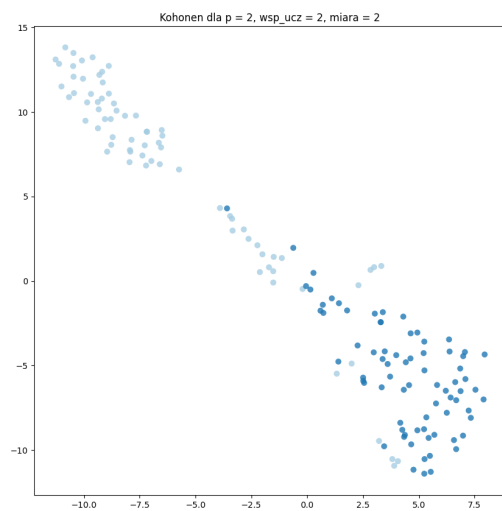
Rysunek 3: Współczynnik uczenia liniowo malejący (1), metoda miary podobieństwa odległość manhattańska



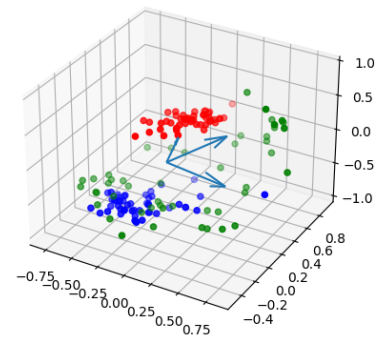
Wizualizacja3D t-SNE dla $p = 2$, $wsp_ucz = 2$, miara = 1



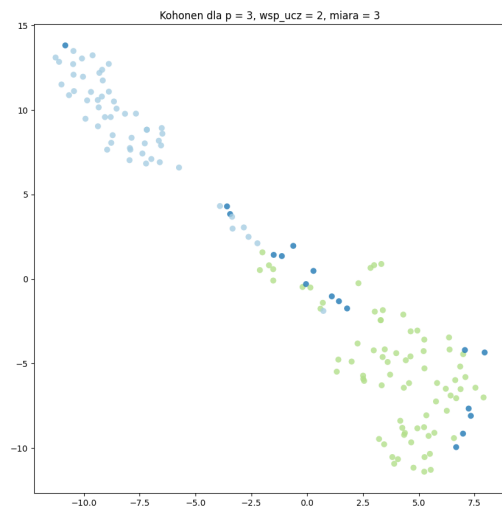
Rysunek 4: Współczynnik uczenia wykładniczo malejący (2), metoda miary podobieństwa iloczyn skalarny



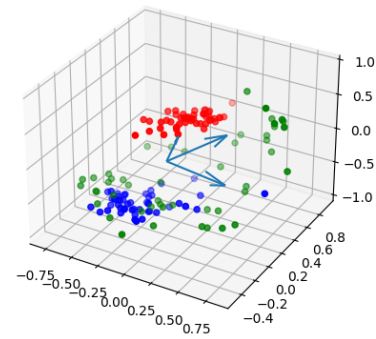
Wizualizacja3D t-SNE dla $p = 2$, $wsp_ucz = 2$, miara = 2



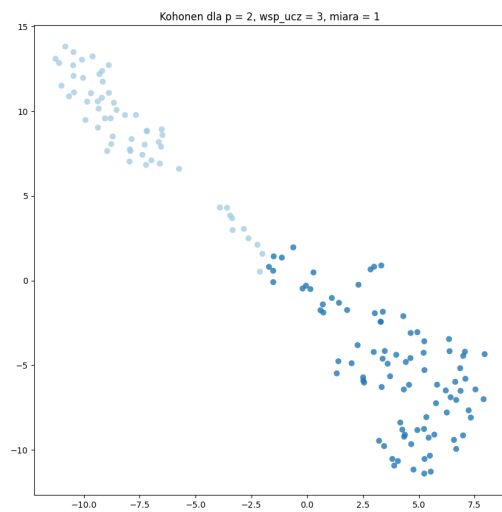
Rysunek 5: Współczynnik uczenia wykładniczo malejący (2), metoda miary podobieństwa odległość euklidesowa



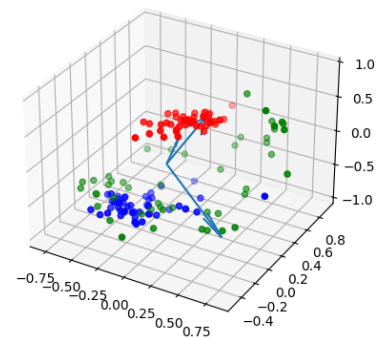
Wizualizacja3D t-SNE dla $p = 2$, $wsp_ucz = 2$, $miara = 3$



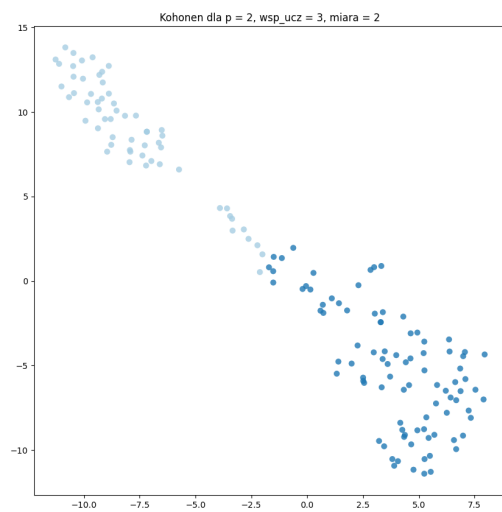
Rysunek 6: Współczynnik uczenia wykładniczo malejący (2), metoda miary podobieństwa odległość manhat-
tańska



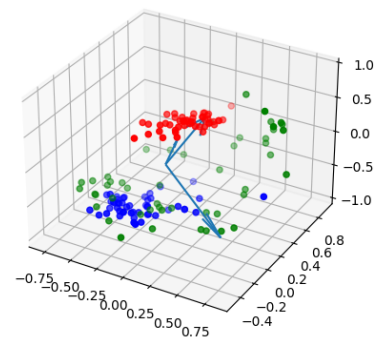
Wizualizacja3D t-SNE dla $p = 2$, $wsp_ucz = 3$, $miara = 1$



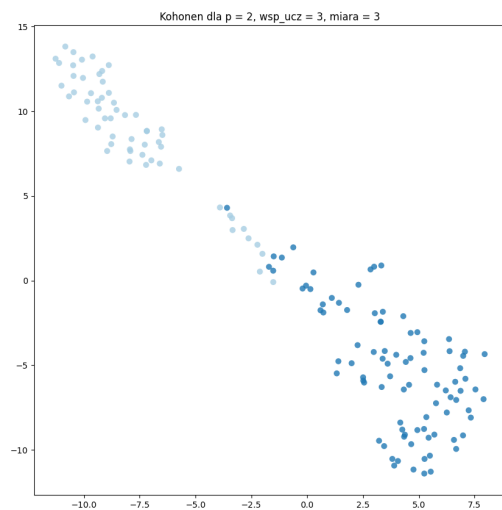
Rysunek 7: Współczynnik uczenia malejący hiperbolicznie (3), metoda miary podobieństwa iloczyn skalarny



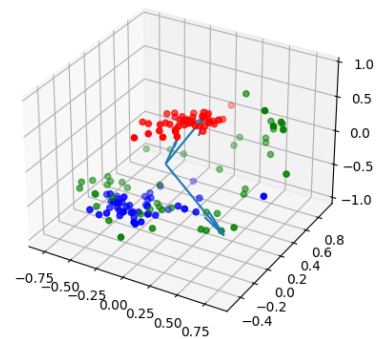
Wizualizacja3D t-SNE dla $p = 2$, $wsp_ucz = 3$, miara = 2



Rysunek 8: Współczynnik uczenia malejący hiperbolicznie (3), metoda miary podobieństwa odległość euklidesowa

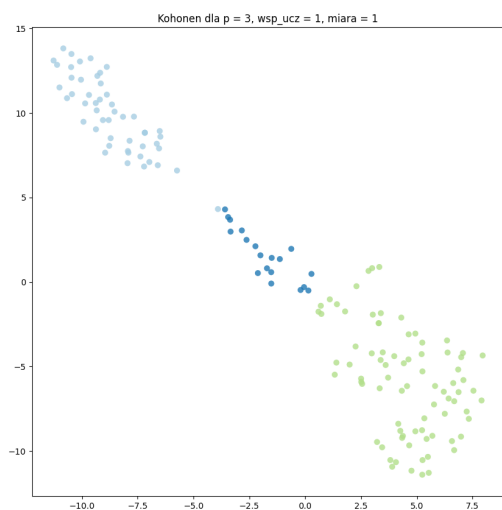


Wizualizacja3D t-SNE dla $p = 2$, $wsp_ucz = 3$, miara = 3

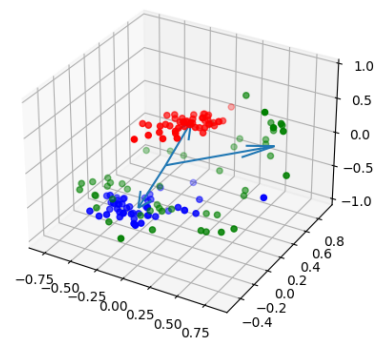


Rysunek 9: Współczynnik uczenia malejący hiperbolicznie (3), metoda miary podobieństwa odległość manhattańska

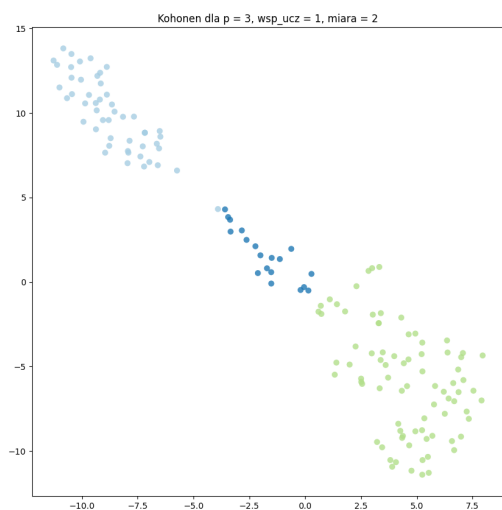
4.1.2 Liczba klas: 3



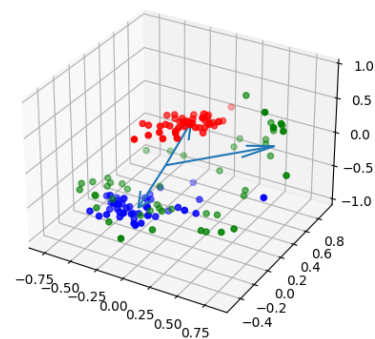
Wizualizacja3D t-SNE dla $p = 3$, $wsp_ucz = 1$, $miara = 1$



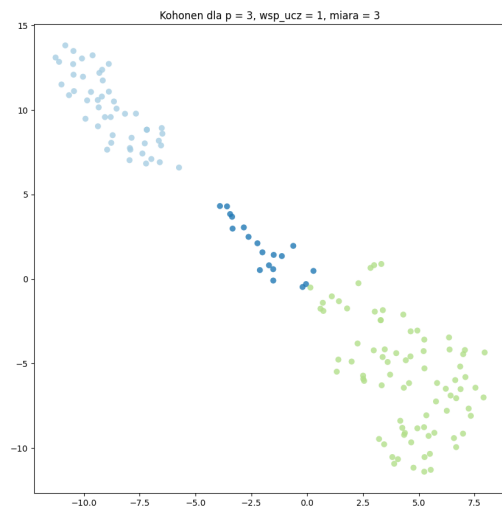
Rysunek 10: Współczynnik uczenia liniowo malejący (1), metoda miary podobieństwa iloczyn skalarny



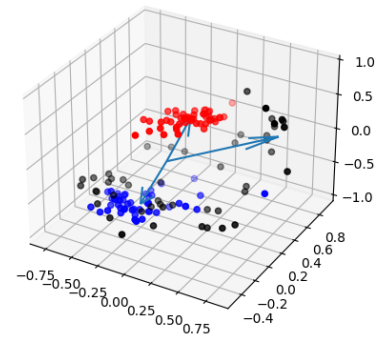
Wizualizacja3D t-SNE dla $p = 3$, $wsp_ucz = 1$, $miara = 2$



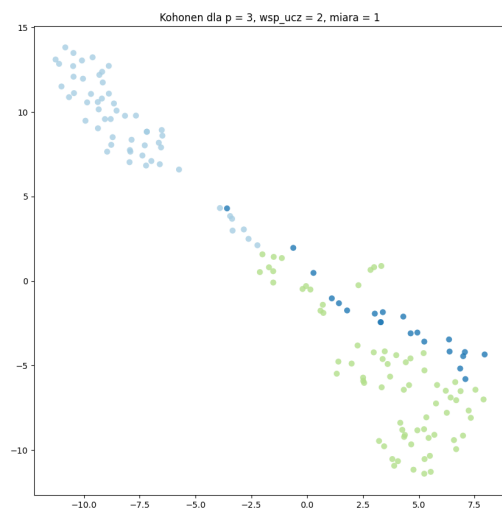
Rysunek 11: Współczynnik uczenia liniowo malejący (1), metoda miary podobieństwa odległość euklidesowa



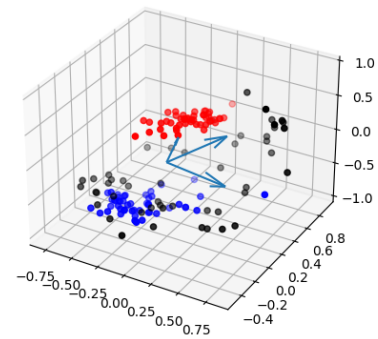
Wizualizacja3D t-SNE dla $p = 3$, $wsp_ucz = 1$, miara = 3



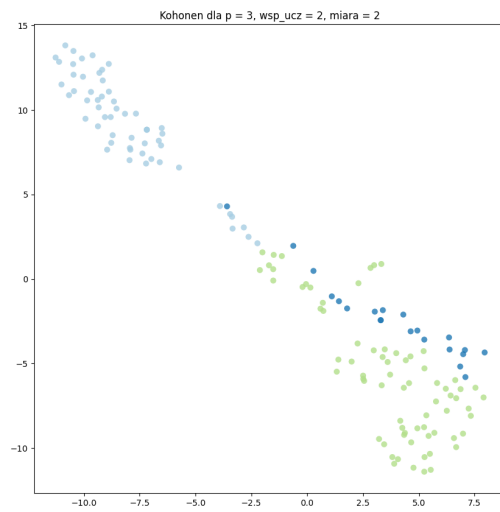
Rysunek 12: Współczynnik uczenia liniowo malejący (1), metoda miary podobieństwa odległość manhattańska



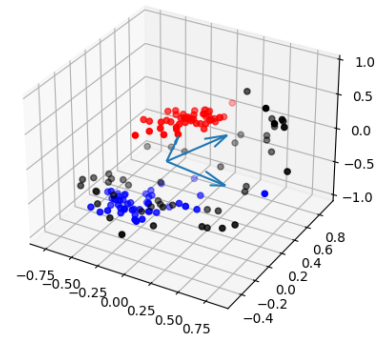
Wizualizacja3D t-SNE dla $p = 3$, $wsp_ucz = 2$, miara = 1



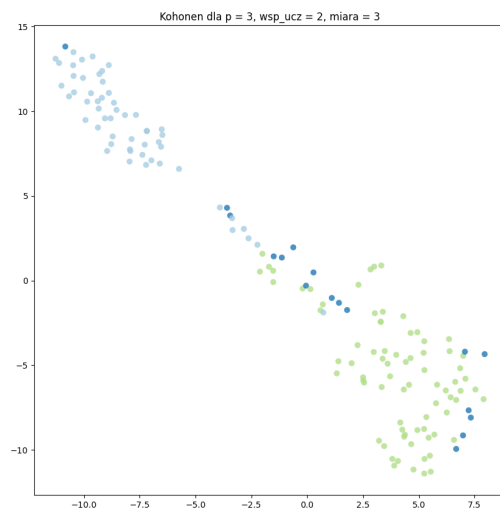
Rysunek 13: Współczynnik uczenia wykładniczo malejący (2), metoda miary podobieństwa iloczyn skalarny



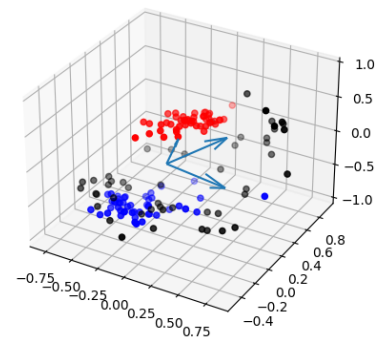
Wizualizacja3D t-SNE dla $p = 3$, $wsp_ucz = 2$, miara = 2



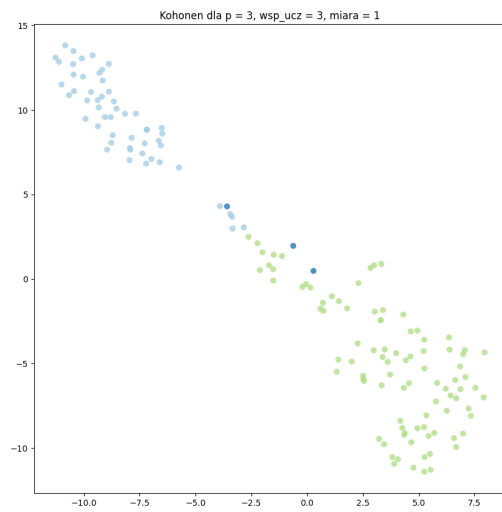
Rysunek 14: Współczynnik uczenia wykładniczo malejący (2), metoda miary podobieństwa odległość euklidesowa



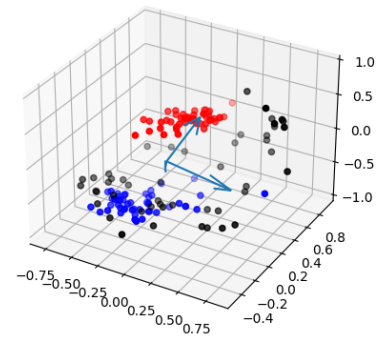
Wizualizacja3D t-SNE dla $p = 3$, $wsp_ucz = 2$, miara = 3



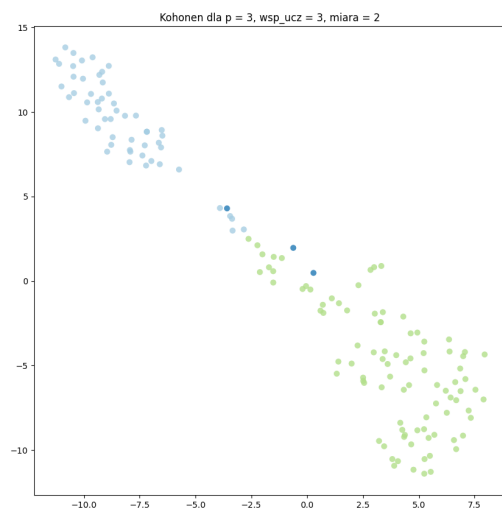
Rysunek 15: Współczynnik uczenia wykładniczo malejący (2), metoda miary podobieństwa odległość manhattanśka



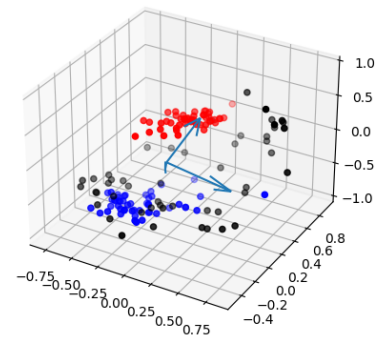
Wizualizacja3D t-SNE dla $p = 3$, $wsp_ucz = 3$, miara = 1



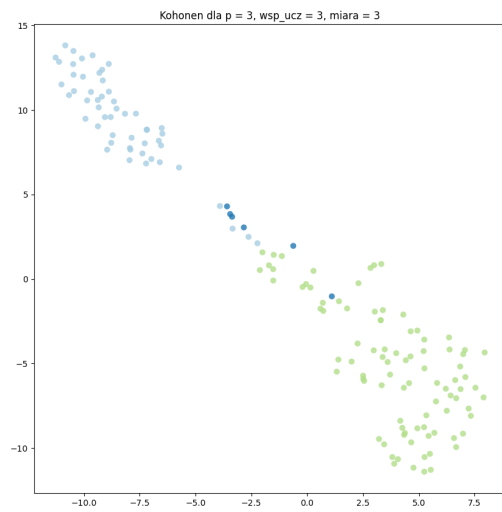
Rysunek 16: Współczynnik uczenia malejący hiperbolicznie (3), metoda miary podobieństwa iloczyn skalarny



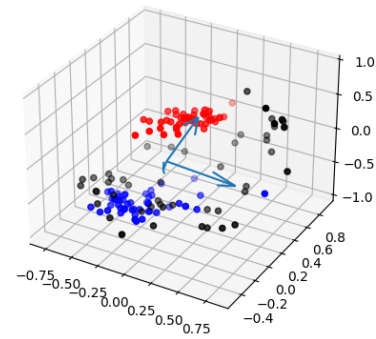
Wizualizacja3D t-SNE dla $p = 3$, $wsp_ucz = 3$, miara = 2



Rysunek 17: Współczynnik uczenia malejący hiperbolicznie (3), metoda miary podobieństwa odległość euklidesowa



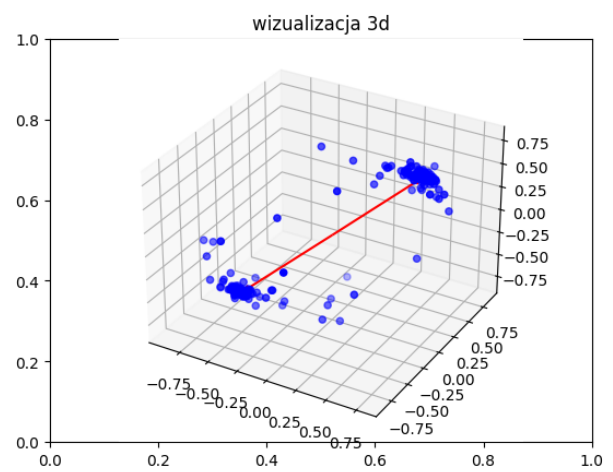
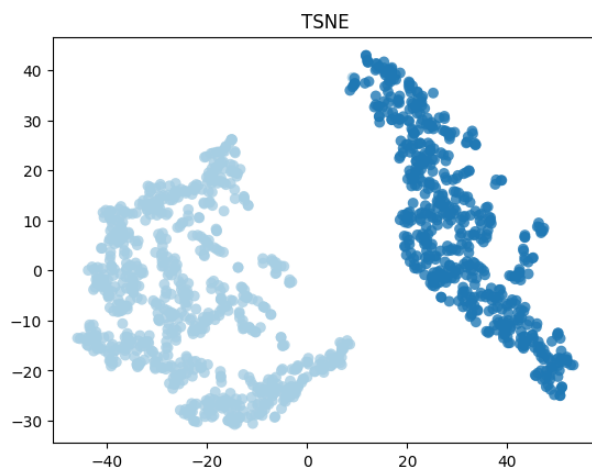
Wizualizacja3D t-SNE dla $p = 3$, $wsp_ucz = 3$, $miara = 3$



Rysunek 18: Współczynnik uczenia malejący hiperbolicznie (3), metoda miary podobieństwa odległość manhattańska

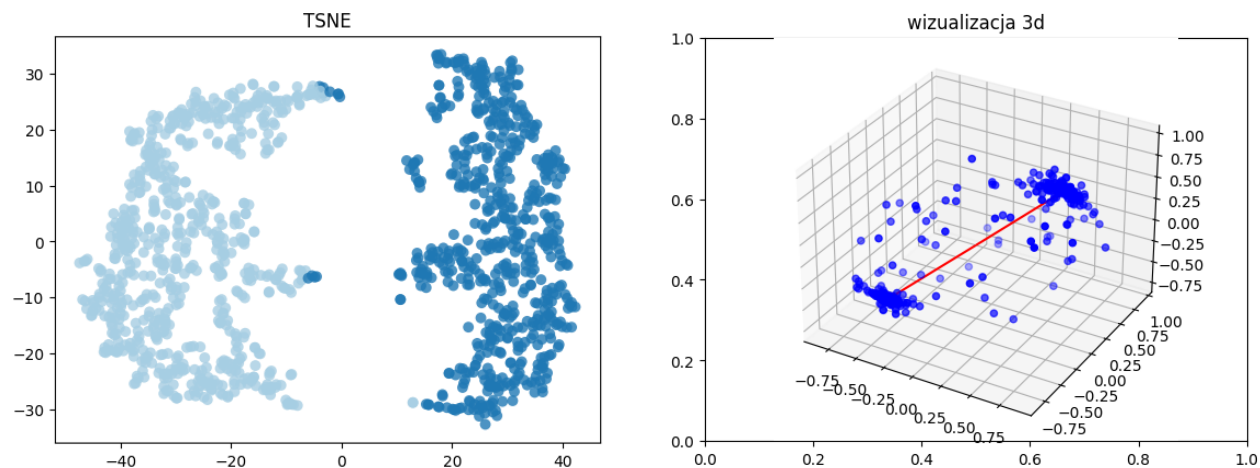
4.2 Dane giełdowe dla trzech spółek

4.2.1 Spółka AOS



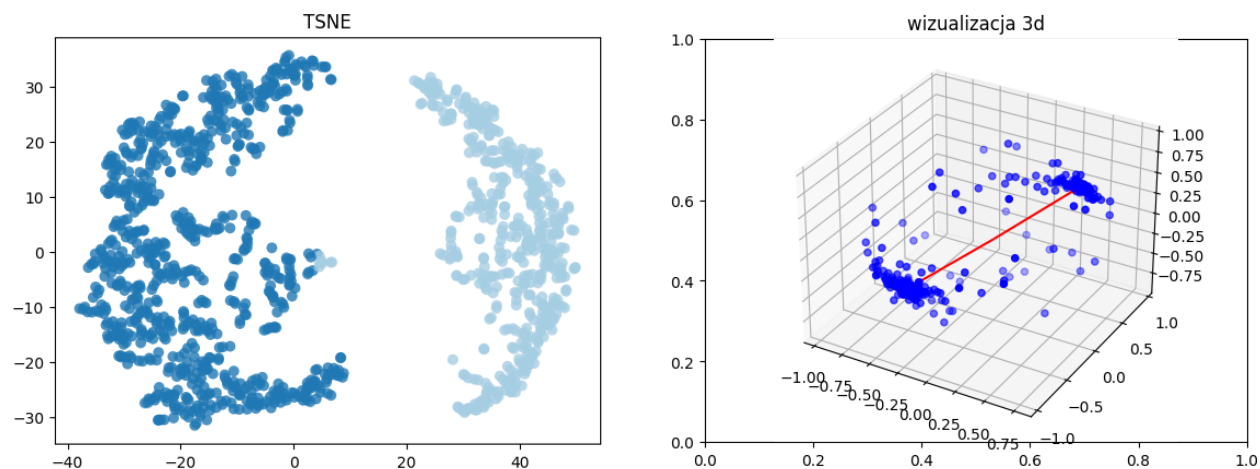
Rysunek 19: t-SNE, wizualizacja 3D, liniowo malejący, metoda manhattańska

4.2.2 Spółka AMP



Rysunek 20: t-SNE, wizualizacja 3D, liniowo malejący, metoda manhattańska

4.2.3 Spółka APD



Rysunek 21: t-SNE, wizualizacja 3D, liniowo malejący, metoda manhattańska

5 Wnioski

Dla zbioru irysów nie zaobserwowano różnicy między metodami iloczynu skalarnego i odległości euklidesowej. Zmiana współczynnika uczenia miała widoczny wpływ na umiejscowienie centrów danych.

6 Kod

6.1 Sposoby mierzenia odległości między danymi a reprezentantami

Listing 1: Metody mierzenia odległości

```
1 def calc_scalar_product(X, T_value, num_repres, repres):
2     measurements = [np.dot(rep, X[T_value % len(X)]) for rep in repres[:
3         num_repres]]
4     max_index = np.argmax(measurements)
5     return max_index
6
7 def calc_euclidean_distance(X, T_value, num_repres, repres):
8     measurements = [np.linalg.norm(rep - X[T_value % len(X)]) for rep in repres
9         [: num_repres]]
10    min_index = np.argmin(measurements)
11    return min_index
12
13 def calc_absolute_difference(X, T_value, num_repres, repres):
14     measurements = []
15
16     for rep in repr[: num_repres]:
17         absolute_diff = np.abs(rep - X[T_value % len(X)])
18         measurements.append(np.sqrt(np.sum(absolute_diff)))
19
20     min_index = np.argmin(measurements)
21     return min_index
22
23 measure_methods = [calc_scalar_product, calc_euclidean_distance,
24     calc_absolute_difference]
```

6.2 Algorytm Kohonena

Listing 2: Kohonen

```
1 def normalize_data(X):
2     data_mean = np.mean(X, axis=0)
3     X_normalized = (X - data_mean) / np.linalg.norm(X - data_mean, axis=1, keepdims=True)
4     return X_normalized
5
6 def initialize_vectors(X, num_vectors):
7     random_generator = np.random.RandomState(0)
8     vectors = random_generator.normal(loc=0.0, scale=0.01, size=(num_vectors, len(X[0])))
9     vectors /= np.linalg.norm(vectors, axis=1, keepdims=True)
10    return vectors
11
12 def update_learning_rate(learning_rate_option, initial_learning_rate, C1, C2,
13    current_iteration, total_iterations):
14    if learning_rate_option == 1:
15        return initial_learning_rate * (total_iterations - current_iteration) /
16            total_iterations
17    elif learning_rate_option == 2:
18        return initial_learning_rate * math.exp(-C1 * current_iteration)
19    elif learning_rate_option == 3:
20        return C1 / (C2 + current_iteration)
21    else:
22        raise ValueError("Invalid learning_rate_option")
23
24 def Kohonen(X, T, initial_learning_rate, num_vectors, learning_rate_option, measure_option,
25    C1, C2, normalize):
26    """
27    X - dane
28    T - liczba iteracji
29    init_lr - współczynnik uczenia
30    num_vecs - liczba reprezentant w
31    lr_option - 1 - malejacy, 2 - wykładniczy, 3 - hiperboliczne zmniejszanie
32    measure_opt - 1 - skalar, 2 - euklidesowa, 3 - manhattanska
33    C1 - parametr dla lr_option = 2
34    C2 - parametr dla lr_option = 3
35    norm - czy normalizowac dane
36    """
37    if normalize:
38        X = normalize_data(X)
39
40    vectors = initialize_vectors(X, num_vectors)
41
42    measurements_table = np.zeros(len(X))
43
44    for iteration in range(T):
45        selected_index = measure_option(X, iteration, num_vectors, vectors)
46        measurements_table[iteration % len(X)] = selected_index
47
48        vectors[selected_index] += initial_learning_rate * (X[iteration % len(X)] - vectors[
49            selected_index])
50        vectors[selected_index] /= np.linalg.norm(vectors[selected_index])
51
52        initial_learning_rate = update_learning_rate(learning_rate_option,
53            initial_learning_rate, C1, C2, iteration, T)
54
55    return vectors, measurements_table, X
```

6.3 Skrypt

Listing 3: Skrypt

```
1 # Definicja parametrów
2 num`iterations` = 10000, learning`rate`initial = 0.1, num`representatives` = 2
3 learning`rate`type = 1, measure`method` = measure`methods`[1]
4 param`C1` = 0.3, param`C2` = 0.4
5 normalize`data` = True
6 plot`title` = 'p=2 IRIS LR=1, MO=2'
7
8 # Wczytanie danych
9 dataset = datasets.load`iris`()
10 data = dataset.data
11
12 # Obliczenie Kohonena
13 kohonen`result`, predictions`, transformed`data` = Kohonen(
14     data,
15     num`iterations`,
16     learning`rate`initial,
17     num`representatives`,
18     learning`rate`type,
19     measure`method`,
20     param`C1`,
21     param`C2`,
22     normalize`data`
23 )
24
25 # Wykresy
26 predictions`np` = np.array(predictions)
27 plt.figure()
28 plt.plot(predictions`np`, 'o', markersize=2)
29 plt.vlines([50,100],0,3, colors='r', linestyle='dashed')
30 plt.title('Kohonen -'.format(plot`title`))
31 plt.xlabel('Dane', fontsize=15)
32 plt.ylabel('Przypisanie', fontsize=15)
33 plt.show()
34
35 # Wykresy t-SNE
36 tsne`result` = TSNE(random`state`=1).fit`transform`(transformed`data`)
37 plt.figure()
38 plt.scatter(tsne`result`[:,0], tsne`result`[:,1], lw=0, s=50, alpha=0.8, edgecolors='black',
39             c=np.array(sns.color`palette`("Paired"))[predictions`np`.astype(np.int)])
40 plt.title('TSNE -'.format(plot`title`))
41 plt.show()
42
43 # Wykresy 3D
44 center = [0,0,0]
45 fig = plt.figure()
46 plt.title('Wizualizacja 3D -'.format(plot`title`))
47 ax = fig.add_subplot(111,projection='3d')
48 u1, v1, w1 = zip(center, center)
49 u2, v2, w2 = zip(kohonen`result`[0, 0:3], kohonen`result`[1, 0:3])
50 ax.quiver(u1, v1, w1, u2, v2, w2, arrow`length`ratio=0.1, colors='black')
51 ax.scatter(transformed`data`[0:49,0], transformed`data`[0:49,1], transformed`data`[0:49,2], c='r')
52 ax.scatter(transformed`data`[50:99,0], transformed`data`[50:99,1], transformed`data`[50:99,2], c='g')
53 ax.scatter(transformed`data`[99:149,0], transformed`data`[99:149,1], transformed`data`[99:149,2], c='b')
54 plt.show()
```
