

Zadanie 2 Pracownia Specjalistyczna - Eksploracja danych

Krzysztof Funkowski

```
In [1]: from math import *\nimport numpy as np\nimport matplotlib.pyplot as plt\nimport random as rnd
```

Zadania

- Budowanie macierzy kowariancji Σ_k dla poszczególnych zbiorów uczących $C_k = \{x_j\}$ w danych Irys ($k = 1, 2, 3$).
- Implementacja iteracyjnej procedury rozwiązywania zagadnienia własnego.
- Wyznaczenie wartości własnych λ_i i wektorów własnych k_i dla wybranej macierzy kowariancji Σ_k .
- Zwiualizować wektory cech x_i na płaszczyznę rozpiętą przez dwa wektory własne k_i

Funkcja służąca do obliczania kowariancji dla dwóch wektorów

```
In [2]: def cov(x, y, x_avg, y_avg):\n    elements_count = len(x) # ilość elementów w tablicy\n    elements_sum = 0\n\n    for x_i, y_i in zip(x, y):\n        elements_sum += elements_sum + ((x - x_avg)*(y - y_avg)) # sumowanie (x - średnia_arytmetyczna_x) * (y - średnia_arytmetyczna_y)\n\n    return elements_sum / elements_count # suma dzielona przez liczbę elementów
```

Funkcja obliczająca macierz kowariancji dla zbioru danych

```
In [3]: def cov_for_dataset(cov_dataset):\n    avgx = []\n    matrix_shape = cov_dataset.shape[1] # wymiary macierzy czyli liczba_kolumn_bioru_danych * liczba_kolumn_bioru_danych\n    for i in range(matrix_shape):\n        avgx.append(cov_dataset[:,i].mean()) # obliczanie średnich dla każdego parametru\n\n    matrix = []\n    for i in range(matrix_shape):\n        matrix_row = []\n        for j in range(matrix_shape):\n            matrix_shape_i += j\n            res = cov(cov_dataset[:,i], cov_dataset[:,j], avgx[i], avgx[j]) # wyznaczenie kowariancji dla każdej pary parametrów\n            matrix_row.append(res) # budowanie wiersza macierzy\n        matrix.append(matrix_row) # dodawanie wiersza macierzy do tablicy bazowej\n\n    return np.matrix(matrix) # przekształcanie typu tablicy na macierz
```

Funkcja do obliczania długości wektora własnego

```
In [4]: def eigenvector_len(v1):\n    dims = v1.shape[0] # ilość elementów wektora\n    length = 0\n\n    for i in range(dims):\n        val = v1[i].ravel().tolist()[0][0]\n        length = length + pow(val,2) # sumowanie elementów podniesionych do kwadratu\n\n    length = sqrt(length) # pierwiastkowanie wyniku sumy elementów wektora podniesionych do kwadratów\n    return length # zwracanie długości
```

Funkcja normalizująca, zwracająca wartość własną i wektor własny

```
In [5]: def eigenvector_normalization(v1):\n    eigenvalue = eigenvector_len(v1) # wyliczenie wartości własnej\n    eigenvector = v1 / eigenvalue # wyliczenie wektora własnego\n\n    return eigenvalue, eigenvector # zwrócenie wartości własnej i wektora własnego
```

Pomocnicza funkcja zmieniająca tablicę wektorów własnych w odpowiednią formę macierzy

```
In [6]: def rearrange_matrix(eigenvectors):\n    res = []\n\n    for i in range(len(eigenvectors)):\n        row = eigenvectors[i].ravel().tolist()[0]\n        res.append(row)\n\n    return np.matrix(res).T
```

Funkcja odpowiadająca za obliczanie wartości własnych i wektorów własnych

```
In [8]: def eigenvalues_eigenvectors(matrix, num_iterations):\n    n = matrix.shape[0] # jeden z wymiarów macierzy, ponieważ jest kwadratowa to wystarczy tylko jeden\n    eigenvalues = [] # lista na wartości własne\n    eigenvectors = [] # lista na wektory własne\n\n    for i in range(n):\n        # losowe dane wektora na początku\n        eigenvector = np.random.rand(n,1)\n\n        # algorytm iteracyjny wykonywany num_iterations razy\n        for _ in range(num_iterations):\n            eigenvector = np.dot(matrix, eigenvector) # mnożenie macierzy kowariancji przez wektor\n            eigenvalue, eigenvector = eigenvector_normalization(eigenvector) # normalizacja wektora\n\n        # usunięcie z macierzy znalezionej wartości własnej i wektora własnego\n        # matrix = wartość_własna * iloczyn_zewnętrzny wektora własnego\n        matrix = matrix - eigenvalue * np.multiply(eigenvector, eigenvector.T).T\n\n        '''if i % 2 == 0:\n            matrix = matrix - eigenvalue * np.multiply(eigenvector, eigenvector.T).T\n            eigenvector = eigenvector * -1\n            matrix = matrix - eigenvalue * np.multiply(eigenvector, eigenvector.T).T'''\n\n    eigenvalues.append(eigenvalue) # dodanie znalezionej wartości własnej do listy\n    eigenvectors.append(eigenvector) # dodanie znalezionej wektora własnego do listy\n\n    return eigenvalues, rearrange_matrix(eigenvectors) # zwrócenie wartości własnych i wektorów własnych
```

Funkcja do rysowania projekcji

```
In [9]: def draw_plot(datasets):\n    colors = ['red', 'blue', 'green']\n\n    print(len(datasets))\n\n    # dla każdego zbioru danych z danej kategorii\n    for i in range(len(datasets)):\n        # przesłanie do wyplotu punktach\n        for j in range(datasets[i][1].shape[0]):\n            if j==0:\n                # zaznaczanie punktu dla pierwszego (SepalLengthCm) i drugiego (SepalWidthCm) parametru\n                plt.scatter(datasets[i][1][j][0], datasets[i][1][j][1], color=colors[i], label=datasets[i][0])\n            else:\n                # zaznaczanie punktu dla pierwszego (SepalLengthCm) i drugiego (SepalWidthCm) parametru\n                plt.scatter(datasets[i][1][j][0], datasets[i][1][j][1], color=colors[i])\n\n    plt.grid()\n    plt.legend()\n    plt.show()
```

Wczytanie zbioru Irisa i konwersja na dataframe'a

```
In [8]: df = pd.read_csv('Iris.csv', sep=',', index_col=0)\n        dataset = df.to_numpy()\n        df
```

	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm	Species	
	Id					
	1	5.1	3.5	1.4	0.2	Iris-setosa
	2	4.9	3.0	1.4	0.2	Iris-setosa
	3	4.7	3.2	1.3	0.2	Iris-setosa
	4	4.6	3.1	1.5	0.2	Iris-setosa
	5	5.0	3.6	1.4	0.2	Iris-setosa

	146	6.7	3.0	5.2	2.3	Iris-virginica
	147	6.3	2.5	5.0	1.9	Iris-virginica
	148	6.5	3.0	5.2	2.0	Iris-virginica
	149	6.2	3.4	5.4	2.3	Iris-virginica
	150	5.9	3.0	5.1	1.8	Iris-virginica
	150 rows x 5 columns					

Wyodrębnienie parametrów ze zbioru

```
In [9]: cov_dataset = np.array(dataset[:,[0,1,2,3]], dtype="double")\n        cov_dataset
```

	array([[5.1, 3.5, 1.4, 0.2],\n [4.9, 3. , 1.4, 0.2],\n [4.7, 3.2, 1.5, 0.2],\n [4.6, 3.1, 1.5, 0.2],\n [5. , 3.6, 1.4, 0.2],\n [5.4, 3.9, 1.7, 0.4],\n [4.6, 3.4, 1.4, 0.3],\n [5. , 3.4, 1.5, 0.2],\n [4.4, 2.9, 1.4, 0.2],\n [4.9, 3.1, 1.5, 0.1],\n [5.4, 3.7, 1.5, 0.2],\n [4.9, 3.4, 1.6, 0.2],\n [4.8, 3. , 1.4, 0.1],\n [4.3, 3. , 1.1, 0.1],\n [5.8, 4. , 1.2, 0.2],\n [5.7, 4.4, 1.5, 0.4],\n [5.4, 3.9, 1.3, 0.4],\n [5.1, 3.5, 1.4, 0.3],\n [5.7, 3.8, 1.7, 0.3],\n [5.1, 3.8, 1.5, 0.3],\n [5.4, 3.4, 1.7, 0.2],\n [5.1, 3.7, 1.5, 0.4],\n [4.6, 3.6, 1. , 0.2],\n [5.1, 3.3, 1.7, 0.5],\n [4.8, 3.4, 1.9, 0.2],\n [5. , 3. , 1.6, 0.2],\n [5. , 3.4, 1.6, 0.4],\n [5.2, 3.5, 1.5, 0.2],\n [5.2, 3.4, 1.4, 0.2],\n [4.7, 3.2, 1.6, 0.2],\n [4.8, 3.1, 1.6, 0.2],\n [5.4, 3.4, 1.5, 0.4],\n [5.2, 4.1, 1.5, 0.1],\n [5.4, 4.2, 1.4, 0.2],\n [4.9, 3.1, 1.5, 0.1],\n [5. , 3.2, 1.2, 0.2],\n [5.5, 3.5, 1.3, 0.2],\n [4.9, 3.1, 1.5, 0.1],\n [4.4, 3. , 1.3, 0.2],\n [5.1, 3.4, 1.5, 0.2],\n [5. , 3.5, 1.2, 0.3],\n [4.5, 2.3, 1.3, 0.3],\n [4.4, 3.2, 1.3, 0.2],\n [5. , 3.5, 1.6, 0.6],\n [5.1, 3.8, 1.5, 0.4],\n [4.8, 3. , 1.4, 0.3],\n [5.1, 3.8, 1.6, 0.2],\n [4.9, 3.2, 1.4, 0.2],\n [5.3, 3.7, 1.5, 0.2],\n [5. , 3.3, 1.4, 0.2],\n [7. , 3.2, 4.7, 1.4],\n [6.4, 3.2, 4.5, 1.5],\n [6.9, 3.1, 4.9, 1.5],\n [5.5, 2.3, 4. , 1.3],\n [6.5, 2.8, 4.6, 1.5],\n [5.7, 2.8, 4.5, 1.3],\n [6.3, 3.3, 4.7, 1.6],\n [4.9, 2.4, 3.3, 1. ,],\n [6.5, 2.9, 4.6, 1.3],\n [5.2, 2.7, 3.9, 1.4],\n [5. , 2. , 3.5, 1. ,],\n [5.8, 3. , 4.2, 1.5],\n [6. , 2.2, 4. , 1. ,],\n [6.1, 2.9, 4.7, 1.4],\n [5.6, 2.9, 3.6, 1.3],\n [6.7, 3.1, 4.4, 1.4],\n [5.8, 3.6, 4.5, 1.5],\n [5.9, 2.7, 4.1, 1. ,],\n [6.1, 2.8, 4.7, 1.2],\n [5.6, 2.5, 3.9, 1.1],\n [5.9, 3.2, 4.8, 1.8],\n [6.1, 2.8, 4. , 1.3],\n [6.3, 2.5, 4.9, 1.5],\n [6.1, 2.8, 4.7, 1.2],\n [6.4, 2.9, 4.3, 1.3],\n [6.6, 3. , 4.4, 1.4],\n [6.8, 3.1, 4.6, 1.4],\n [6.7, 3. , 4.5, 1.7],\n [5.7, 2.6, 3.5, 1. ,],\n [5.5, 2.4, 3.8, 1.1],\n [5.8, 2.7, 3.9, 1.2],\n [6. , 2.7, 5.1, 1.6],\n [5.4, 2. , 4.5, 1.5],\n [6. , 3.4, 4.5, 1.6],\n [6.7, 3.1, 4.7, 1.5],\n [6.3, 2.3, 4.4, 1.3],\n [5.8, 2. , 4.1, 1.3],\n [5.5, 2.5, 4. , 1.3],\n [5.5, 2.6, 4.4, 1.2],\n [6.1, 3. , 4.6, 1.4],\n [5.8, 2.6, 4. , 1.2],\n [5. , 2.3, 3.8, 1. ,],\n [5.6, 2.7, 4.2, 1.3],\n [5.7, 3. , 4.2, 1.2],\n [5.7, 2.9, 4.2, 1.3],\n [6.2, 2.9, 4.3, 1.3],\n [5.1, 2.5, 3. , 1.1],\n [5.7, 2.8, 4.1, 1.3],\n [6.3, 3.3, 6. , 2.5],\n [5.8, 2.7, 5.1, 1.9],\n [7.1, 3. , 5.9, 2.1],\n [6.3, 2.9, 5.6, 1.8],\n [6.5, 3. , 5.8, 2.2],\n [7.6, 3. , 6.6, 2.1],\n [4.9, 2.5, 4.5, 1.7],\n [7.3, 2.9, 6.3, 1.8],\n [6.7, 2.5, 5.8, 1.8],\n [7.2, 3.6, 6.1, 2.5],\n [6.5, 3.2, 5.1, 2. ,],\n [6.4, 2.7, 5.1, 1.9],\n [6.8, 3. , 5.5, 2.1],\n [5.7, 2.5, 5. , 2. ,],\n [5.8, 2.8, 6.1, 2.4],\n [6.4, 3.2, 5.3, 2.3],\n [6.5, 3. , 5.5, 1.8],\n [7.7, 3.8, 6.7, 2.2],\n [7.7, 2.6, 6.9, 2.3],\n [6. , 2.2, 5. , 1.5],\n [6.9, 3.2, 5.7, 2.3],\n [5.6, 2.8, 4.9, 2. ,],\n [7.7, 2.8, 6.7, 2. ,],\n [6.3, 2.7, 4.9, 1.8],\n [6.7, 3.3, 5.7, 2.1],\n [7.2, 3.2, 6. , 1.8],\n [6.2, 2.8, 4.8, 1.8],\n [6.1, 3. , 4.9, 1.8],\n [6.4, 2.8, 5.6, 2.1],\n [7.2, 3. , 6.5, 1.6],\n [7.4, 2.8, 6.1, 1.9],\n [7.9, 3.8, 6.4, 2. ,],\n [6.4, 2.8, 5.6, 2.2],\n [6.3, 2.8, 5.1, 1.5],\n [6.1, 2.6, 5.6, 1.4],\n [7.7, 3. , 6.1, 2.3],\n [6.3, 3.4, 5.6, 2.4],\n [6.4, 3.1, 5.5, 1.8],\n [6. , 3. , 4.8, 1.8],\n [6.5, 3.1, 5.4, 2.1],\n [6.7, 3.1, 5.6, 2.4],\n [6.9, 3.1, 5.2, 2.3],\n [5.8, 2.7, 5.1, 1.9],\n [6.8, 3.2, 5.9, 2.3],\n [6.7, 3. , 5.2, 2.3],\n [6.3, 2.5, 5. , 1.9],\n [6.5, 3. , 5.2, 2.1],\n [6.2, 3.4, 5.4, 2.3],\n [5.9, 3. , 5.1, 1.8]])
--	--

Macierz kowariancji dla zbioru Irisa

```
In [10]: cov_matrix = cov_for_dataset(cov_dataset)\n        cov_matrix
```

	matrix([[0.68112222, -0.03906607, 1.26519111, 0.51345778],\n [-0.03906607, 0.18675807, -0.319588 , -0.11719467],\n [1.26519111, -0.319588 , 0.89242489, 1.28774489],\n [0.51345778, -0.11719467, 1.28774489, 0.57855156]])
--	--

Występujące gatunki w zbiorze danych

```
In [11]: species = np.array(list(set(dataset[:,4])))\n        species
```

	array(['Iris-setosa', 'Iris-versicolor', 'Iris-virginica'], dtype='<U15')
--	---

Wartości własne i wektory własne na podstawie macierzy kowariancji

```
In [12]: num_iterations = 10000\n        eigenvalues, eigenvectors = eigenvalues_eigenvectors(cov_matrix, num_iterations)\n\n        print(f"\"Wartość własna:\\n {eigenvalues}\"")\n        print(f"\"Wektor własny:\\n {eigenvectors}\"")\n\n        Wartości własne: [ 4.19667316316379795, 0.24962861448331192, 0.87898041537352695, 0.8235251482784947]\n\n        Wektor własny: [[ 0.6862863 , 0.6953989 , 0.5899728 , 0.31725455]\n        [-0.08226889, -0.72971237 0.59641869 -0.3240435]\n        [-0.85657211 -0.1757674 0.67252488 -0.47971899]\n        [-0.35884393 -0.87479647 0.54966991 0.75112054]]
```

Wartości własne i wektory własne wygenerowane przez gotowe narzędzie z biblioteki NumPy w celu sprawdzenia poprawności

```
In [13]: # porównanie wyników metody zaimplementowanej w domu z gotową z biblioteki NumPy\n        eigenvalues2, eigenvectors2 = np.linalg.eig(cov_matrix)\n        eigenvalues2, eigenvectors2
```

	array([[4.19667316, 0.24962861, 0.87898042, 0.82352514],\n [-0.0822689 , -0.72971237, 0.59641869, -0.3240435],\n [-0.85657211, -0.1757674 , 0.67252488, -0.47971899],\n [-0.35884393, 0.87479647, 0.54966991, 0.75112054]])
--	---

Dla każdego gatunku wyznaczyć macierz kowariancji, następnie wyznaczenie na jej podstawie wartości własnych i wektorów własnych. Kolejno projekcja pierwszych dwóch parametrów (SepalLengthCm, SepalWidthCm) zbioru danych dla gatunków rozpiętych na dwa pierwsze wektory własne "badanego" w danym momencie badany.

```
In [20]: # dla każdego gatunku\n        for s in species:\n            current_datasets = [] # lista w której będą zbiory do rysowania\n            current_species_dataset = df[df['species'].str.contains(s)] to_numpy() # wybieranie danych tylko dla badanego gatunku\n            current_species_dataset = current_species_dataset[:, [0, 1, 2, 3]] # wybranie parametrów ze zbioru danych\n\n            num_iterations = 10000 # ustalenie liczby iteracji\n            current_cov_matrix = cov_for_dataset(current_species_dataset) # macierz kowariancji dla głównego w tym momencie gatunku\n            wartości_własne_i_wektory_własne_dla_badanego_gatunku\n            current_eigenvalues, current_eigenvectors = eigenvalues_eigenvectors(current_cov_matrix, num_iterations)\n            print(f"\"{s}\\n macierz cov:\\n{current_cov_matrix}\\n\\n wartości własne: {current_eigenvalues}\\n\\n wektory własne: {current_eigenvectors[:, :2]}\\n\\n\"\n            current_datasets.append(s, projection) # dodanie do listy do rysowania\n\n            other_species = species[species!=s] # wyznaczenie zbioru gatunków które nie są aktualnym badanym\n            # dla każdego gatunku który nie jest aktualnie badany / nie jest głównym w danej iteracji\n            for os in other_species:\n                other_species_dataset = df[df['species'].str.contains(os)] to_numpy() # wybieranie danych tylko dla niezłównego gatunku\n                other_species_dataset = other_species_dataset[:, [0, 1, 2, 3]] # wybranie parametrów ze zbioru danych\n\n                # projekcja zbioru danych dla gatunku innego niż główny w danej iteracji, rozpiętego na dwa pierwsze wektory własne\n                other_projection = other_species_dataset * current_eigenvectors[:, :2]\n                current_datasets.append(os, other_projection) # dodanie do listy do rysowania\n\n        draw_plot(current_datasets) # rysowanie projekcji
```

	Iris-setosa:\n macierz cov:\n [[0.121764 0.098292 0.015816 0.010336]\n [0.088292 0.142276 0.011448 0.011380]\n [0.015816 0.011448 0.029584 0.005584]\n [0.010336 0.011288 0.005584 0.011264]]\n\n wartości własne: [0.2337492578217778, 0.035534325661655904, 0.026312239196171977, 0.089212186380394142]\n\n wektory własne: [[0.6662863 , 0.69592481 -0.43478033 -0.86773998]\n [0.73535931 -0.61961464 -0.26584074 -0.86517629]\n [0.69478706 0.49163428 0.83463886 -0.22954169]\n [0.07042406 0.08502085 0.28913895 0.97163499]]
--	---

	dwa pierwsze wektory własne:\n [[0.6662863 , 0.69592481]\n [-0.73535931 -0.61961464]\n [0.69478706 0.49163428]\n [0.07042406 0.08502085]]
--	--



	Iris-versicolor:\n macierz cov:\n [[0.2611804 0.08348 0.17924 0.054664]\n [0.08348 0.0855 0.081 0.04806]\n [0.17924 0.081 0.2164 0.07164]\n [0.054664 0.04806 0.07164 0.038324]]\n\n wartości własne: [0.4781164652566476, 0.07093641392669654, 0.05368056334125243, 0.009594557476003636]\n\n wektory własne: [[0.68627376 -0.66908906 0.26588336 0.1822796]\n [0.39534703 0.56746531 0.72617766 -0.22891939]\n [0.62366312 0.34332698 0.83463886 -0.3196679]\n [0.21498369 0.33539513 -0.06366081 0.91564087]]
--	--

	dwa pierwsze wektory własne:\n [[0.68627376 -0.66908906]\n [0.39534703 0.56746531]\n [0.62366312 0.34332698]\n [0.21498369 0.33539513]]
--	---



	Iris-virginica:\n macierz cov:\n [[0.386256 0.091888 0.297224 0.048112]\n [0.091888 0.101924 0.069952 0.046676]\n [0.297224 0.069952 0.295496 0.047048]\n [0.048112 0.046676 0.047848 0.073924]]\n\n wartości własne: [0.68113497414689959, 0.1844280142375677, 0.05124951922482228, 0.033588053789052519]\n\n wektory własne: [[0.74101679 -0.16525895 -0.53445017 0.37341165]\n [0.2832772 0.74864279 -0.3253749 -0.54068495]\n [0.62789179 -0.16942778 0.65152357 -0.39059336]\n [0.12377451 0.61928804 -0.4206953 0.64587225]]
--	---

	dwa pierwsze wektory własne:\n [[0.74101679 -0.16525895]\n [0.2832772 0.74864279]\n [0.62789179 -0.16942778]\n [0.12377451 0.61928804]]
--	---

