

Eksploracja danych

Piotr Lipiński

Lista zadań nr 2 – Minikurs Pythona z NumPy

Zadanie 0. (rozgrzewkowe, 1 punkt, zadanie należy przesłać emailiem)

Utwórz skrypt definiujący poniższe zmienne:

$a = [1, 2, 3, 4, \dots, 100]$ (wektor złożony z liczb całkowitych od 1 do 100)

$b = [1, 3, 5, 7, \dots, 99]$ (wektor złożony z liczb całkowitych nieparzystych od 1 do 99)

$c = [-1.00 * \pi, -0.99 * \pi, \dots, -0.01 * \pi, 0, 0.01 * \pi, \dots, 0.99 * \pi, 1.00 * \pi]$

$d = [-1.00 * \pi, -0.99 * \pi, \dots, -0.01 * \pi, 0.01 * \pi, \dots, 0.99 * \pi, 1.00 * \pi]$

$e = [e_1, e_2, \dots, e_{100}]$, gdzie $e_i = \sin(i)$, jeśli $\sin(i) > 0$, lub $e_i = 0$ w przeciwnym przypadku

A = macierz rozmiaru 10×10 zawierająca liczby całkowite od 1 do 100: w pierwszym wierszu od lewej 1, 2, ..., w drugim wierszu od lewej 11, 12, ..., itd. (wskazówka: użyć polecenia `reshape`)

B = macierz trójdzielna rozmiaru 100×100 mająca na głównej przekątnej liczby całkowite od 1 do 100, a poniżej i powyżej głównej przekątnej liczby od 99 do 1

C = macierz trójkątna górna wypełniona jedynekami (łącznie z główną przekątną)

D = macierz rozmiaru 2×100 , w której pierwszy wiersz zawiera elementy $d_{1i} = 1 + 2 + \dots + i$, a drugi wiersz zawiera elementy $d_{2i} = i!$

E = macierz rozmiaru 100×100 mająca 1 w pozycji (i, j) , jeśli i dzieli j , lub 0 w przeciwnym przypadku.

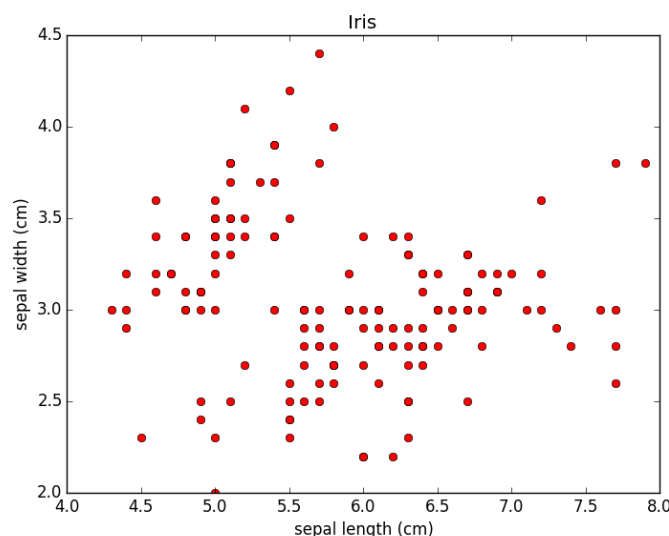
Zadanie 1. (1 punkt)

a) Wczytaj dane IRIS. Można to szybko zrobić korzystając z biblioteki SciKit za pomocą polecenia:

```
from sklearn import datasets
iris = datasets.load_iris()
```

Zobacz co zawiera `iris.data`, `iris.target`, `iris.feature_names`, `iris.target_names`.

b) Przedstaw wczytane dane na wykresie w poniższy sposób (zwróć uwagę na kolory i typ znaczników, opisy osi i tytuł wykresu):



c) Zmień zakres osi: oś X powinna pokazywać wartości od 3 do 9, a oś Y od 1 do 5.

d) Zmień podziałki na osiach, tak aby zaznaczone były tylko liczby całkowite.

e) Każdy gatunek irysa zaznacz innym kolorem.

f) Zapisz rysunek do pliku `zadanie1.png`.

Zadanie 2. (1 punkt)

- a) Zrób rysunek podobny do tego z poprzedniego zadania, ale umieść na nim tylko irysy gatunku *setosa* i *versicolor* (nie rysuj irysów gatunku *versicolor*).
- b) Dodaj do rysunku prostą o równaniu $y = 2x - 8$.
- c) Irysy gatunku *setosa* znajdujące się pod narysowaną linią zaznacz na czerwono, pozostałe na zielono. Irysy gatunku *virginica* znajdujące się nad narysowaną linią zaznacz na czerwono, pozostałe na zielono.
- d) Zapisz rysunek do pliku zadanie2a.png.
- e) Spróbuj zmienić równanie prostej z punktu b) tak, aby zmniejszyć liczbę czerwonych punktów.
- f) Zapisz rysunek do pliku zadanie2b.png.

Zadanie 3. (1 punkt)

- a) Zrób rysunek przedstawiający 10 punktów o następujących współrzędnych (1, 10), (2, 10), (3, 11), (4, 12), (5, 18), (6, 18), (7, 19), (8, 26), (9, 19), (10, 26).
- b) Dodaj do rysunku prostą o równaniu $y = 2x + 5$.
- c) Zapisz rysunek do pliku zadanie3a.png.
- d) Dla każdego punktu danych policz jego odległość od wyznaczonej prostej.
- e) Spróbuj zmienić równanie prostej z punktu b) tak, aby zmniejszyć sumę odległości punktów danych od prostej.
- f) Zapisz rysunek do pliku zadanie3b.png.

Zadanie 4. (1 punkt)

Napisz program generujący zestaw dwuwymiarowych danych losowych złożony z K chmur punktów, taki że:

- a) każda chmura punktów składała się z 1000 punktów o współrzędnych (x, y) , gdzie x pochodzi z rozkładu normalnego $N(a_i, 1)$, y pochodzi z rozkładu normalnego $N(b_i, 1)$, zaś (a_i, b_i) to centrum i -tej chmury punktów,
 - b) centra chmur punktów (a_i, b_i) tworzą wielokąt foremny o boku o zadanej długości d .
- Uruchom program dla $K = 7$, $K = 11$ i $K = 23$ oraz $d = 5$, $d = 10$ i $d = 15$. Zrób rysunki przedstawiające wyniki.

Zadanie 5. (1 punkt)

Dla danych z poprzedniego zadania policz odległość każdego punktu danych od każdego centrum chmury i na sporządzonych wykresach zaznacz kolorem czerwonym te punkty danych, które znajdują się bliżej centrum innej chmury niż chmury, z której pochodzą, a kolorem zielonym pozostałe punkty danych. Jak zależy frakcja punktów czerwonych od długości d ? Jaka powinna być wartość d , żeby punkty czerwone stanowiły około 10% wszystkich punktów danych?

UWAGA: Proszę nie korzystać z żadnych funkcji wbudowanych ani bibliotecznych liczących odległości, iloczyny skalarne, itp. Proszę sprawdzić działanie swoich funkcji na przykładowych danych (najlepiej dość dużych rozmiarów). Proszę spróbować ocenić efektywność swoich obliczeń.