

## Sieci neuronowe i sztuczna inteligencja – laboratorium 5

Monika Błyszcz, 236623

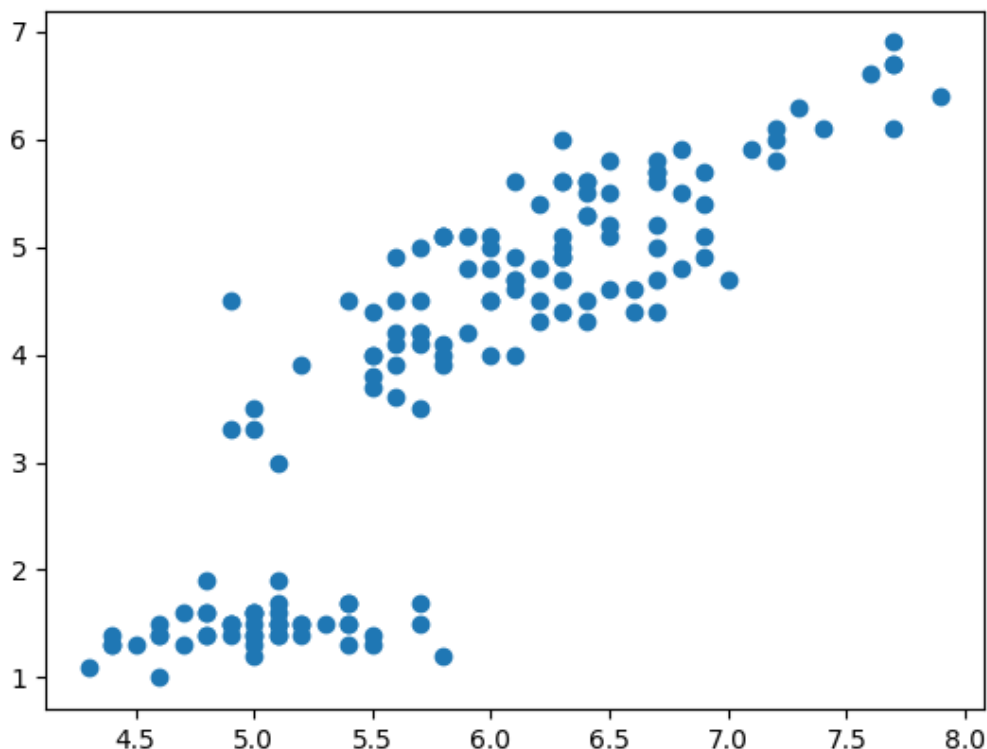
W zadaniu zaimplementowano 3 różne testy z kategorii korelacji. Wybrano korelację Pearsona, Spearmana i Kendalla.

- Współczynnik korelacji Pearsona można wykorzystać do podsumowania siły liniowej zależności między dwiema próbkami danych. Oblicza się go jako kowariancję dwóch zmiennych podzieloną przez iloczyn odchylenia standardowego każdej próbki danych. Jest to normalizacja kowariancji między dwiema zmiennymi, aby uzyskać interpretowalny wynik.
- Współczynnik korelacji Spearmana można wykorzystać do podsumowania siły między dwiema próbkami danych. Zamiast obliczać współczynnik przy użyciu kowariancji i odchyłeń standardowych dla samych próbek, statystyki te są obliczane na podstawie względnej rangi wartości na każdej próbce. Jest to powszechne podejście stosowane w statystyce nieparametrycznej, np. metodach statystycznych, w których nie zakładamy rozkładu danych, takich jak gaussowskie. W tym przypadku nie zakłada się liniowej tylko monotoniczną zależność między zmiennymi
- Korelacja rang Kendalla (współczynnik korelacji Kendalla) polega na tym, że oblicza znormalizowany wynik dla liczby pasujących lub zgodnych rankingów między dwiema próbkami. Test przyjmuje dwie próbki danych jako argumenty i zwraca współczynnik korelacji oraz wartość p. Jako test hipotezy statystycznej metoda zakłada ( $H_0$ ), że nie ma związku między dwiema próbkami

Poniżej przedstawiono zaimplementowany kod. Jako bazę danych wykorzystano bazę kwiatów iris\_datatset. Szukano korelacji między Sepal Length (długość działki kielicha) a Petal Length (długość płatka).

```
1 import pandas
2 from scipy.stats import pearsonr
3 from scipy.stats import spearmanr
4 from scipy.stats import kendalltau
5 from matplotlib import pyplot
6
7 #Load dataset
8 url="https://raw.githubusercontent.com/jbrownlee/Datasets/master/iris.csv"
9 names=["sepal-length", "sepal-width", "petal-length", "petal-width", "class"]
10 dataset = pandas.read_csv(url, names=names)
11 #data
12 data1 = dataset.loc[:, 'sepal-length']
13 data2 = dataset.loc[:, 'petal-length']
14
15 #show data
16 pyplot.scatter(data1, data2)
17 pyplot.show()
18
19 # Pearson's Correlation test
20 print("Pearson's correlation")
21 stat, p = pearsonr(data1, data2)
22 print('stat=%.3f, p=%.3f' % (stat, p))
23 if p > 0.05:
24     print('Probably independent')
25 else:
26     print('Probably dependent')
27
28 #Spearman's Rank Correlation Test
29 print(' ')
30 print("Spearman's correlation")
31 stat, p = spearmanr(data1, data2)
32 print('stat=%.3f, p=%.3f' % (stat, p))
33 if p > 0.05:
34     print('Probably independent')
35 else:
36     print('Probably dependent')
37
38 #Kendall's Rank Correlation Test
39 print(' ')
40 print("Kendall's correlation")
41 stat, p = kendalltau(data1, data2)
42 print('stat=%.3f, p=%.3f' % (stat, p))
43 if p > 0.05:
44     print('Probably independent')
45 else:
46     print('Probably dependent')
```

**Najpierw wykreślono zależność na wykresie:** widzimy, że prawie każda zmienna ma równomierny rozkład, a dodatni związek między zmiennymi jest widoczny przez ukośne grupowanie punktów od lewego dolnego do prawego górnego rogu wykresu. Zależność nie jest do końca liniowa, bo w dolnym lewym rogu występują odseparowane dane. Odrzucenie tych danych sprawi, że zależność będzie liniowa.



Otrzymano następujące wyniki:

```
C:\Users\Mo\AppData\Local\Microsoft\WindowsApps\python3.10.exe C:\Users\Mo\
Pearson's correlation
stat=0.872, p=0.000
Probably dependent
```

```
Spearman's correlation
stat=0.881, p=0.000
Probably dependent
```

```
Kendall's correlation
stat=0.718, p=0.000
Probably dependent
```

```
Process finished with exit code 0
```

Analizując powyższe wyniki widzimy, że występuje dość duża, pozytywna korelacja pomiędzy długością kielicha a długością płatka. Dla poszczególnych testów statystycznych mamy:

- **Korelacja Pearsona:** obie zmienne są dodatnio skorelowane i że korelacja wynosi 0,872. Sugeruje to wysoki poziom korelacji pomiędzy wartościami
- **Korelacja Spearmana:** test statystyczny wykazuje silną dodatnią korelację z wartością 0,881. Wartość p równa zero, co oznacza, że prawdopodobieństwo zaobserwowania danych, uwzględniając, że próbki są nieskorelowane, jest bardzo mało prawdopodobne. Pozwala to na odrzucenie hipotezy zerowej, że próbki są nieskorelowane.
- **Korelacja rang Kendalla:** wyniosła 0,718, co oznacza, że zależność jest wysoce skorelowana. Wartość p jest bliska zero (i wydrukowana jako zero), podobnie jak w korelacji Spearmana, co oznacza, że możemy odrzucić hipotezę zerową, że próbki są nieskorelowane.

Najwyższy współczynnik korelacji pokazała korelacja Spearmana, co jest zgodne z rzeczywistością, bo dane nie są w pełni liniowe.