

## Zadanie kolokwium 2

Wojciech Ganobis 310519, Bartosz Troszka 309912

05/05/20

Celem zadania jest obliczenie współczynnika korelacji między zachorowaniami, na podstawie dostarczonych danych. Najpierw omówmy metodę zrobienia tego zadania.

Do wykonania zadania użyjemy współczynnika korelacji Pearsona. Wyraża się on wzorem:

$$\begin{aligned} r_{xy} &= \frac{\text{cov}(r, y)}{\sigma_X \sigma_Y} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}} = \\ &= \frac{n \cdot \sum x_i y_i - \sum x_i \cdot \sum y_i}{\sqrt{[n \cdot \sum x_i^2 - (\sum x_i)^2] \cdot [n \cdot \sum y_i^2 - (\sum y_i)^2]}}, \end{aligned}$$

gdzie  $r_{XY} \in [-1, 1]$ . Im większa wartość bezwzględna, tym silniejsza zależność między cechami. Jeśli  $r_{XY} == 1$  to zależność jest dokładnie liniowa. Natomiast jeśli  $r_{XY} == 0$  to nie ma liniowej zależności pomiędzy cechami. Zależności można opisać za pomocą tabelki:

Korelacje	Ujemne	Dodatnie
Slabe	-0.5 do 0.0	0.0 do 0.5
Silne	-1.0 do -0.5	0.5 do 1.0

Kolejną rzeczą ważną w zadaniu jest wybranie punktu początkowego. "Pierwsze" dane w każdym państwie są niedokładne oraz niezetelne. Zaczniemy więc obliczanie korelacji gdy liczba zakazanych będzie zwiększać się każdego dnia.

Przejdźmy teraz do obliczenia pierwszej korelacji.

- Pierwsze państwa między którymi będziemy szukać korelacji to Włochy i Hiszpania. Jako punkt początkowy dla Włoch wybierzemy 22 luty, ponieważ po tym dniu każdego następnego dnia liczba zarazonych przybywa. Dla Hiszpani natomiast wybierzemy 27 luty, z tego samego powodu. Z racji że dane miały być do 30 kwietnia punkt końcowy dla Hiszpani to 30 kwietnia, a dla Włoch 25 kwietnia, aby liczba dni była taka sama. Do obliczenia wykorzystamy wzór  $\frac{n \cdot \sum x_i y_i - \sum x_i \cdot \sum y_i}{\sqrt{[n \cdot \sum x_i^2 - (\sum x_i)^2] \cdot [n \cdot \sum y_i^2 - (\sum y_i)^2]}}$ .

Dla ułatwienia obliczeń zrobimy tabelkę:

x	y	$x_i y_i$	$x_i^2$	$y_i^2$
6	10	60	36	100
67	13	871	4489	169
48	7	336	2304	49
105	13	1365	11025	169
...	...	...	...	...

$\sum$ kolumna1	$\sum$ kolumna2	$\sum$ kolumna3	$\sum$ kolumna4	$\sum$ kolumna5
192991	225045	988355044	810008099	1286773571

Teraz można łatwo obliczyć współczynnik wstawiając do wzoru:

$$\frac{64 \cdot 988355044 - 192991 \cdot 225045}{\sqrt{[64 \cdot 810008099 - (192991)^2] \cdot [64 \cdot 1286773571 - (225045)^2]}} = 0.8827529322$$

Wynik ten oznacza że korelacja jest silna dodatnia.

- Obliczmy teraz korelację dla innych państw. Weźmy np. Polskę oraz Francję. Za punkt początkowy Polski weźmy 7 marzec, a końcowy 3 maja. Dla Francji będzie to 27 luty i 24 kwietnia. Metoda ta sama więc nie będę zapisywał obliczeń tylko sam wynik który wynosi:

$$0,7512537726$$

Oznacza to również że korelacja jest silna, jednak słabsza niż między Hiszpanią a Włochami.

- Obliczmy teraz korelację między naszymi sąsiadami mianowicie między Czechami a Słowacją. Za punkt początkowy Słowacji weźmy 5 marca, a końcowy 2 maja. W Czechach natomiast dzień później czyli 6 marca do 3 maja. Korelacja wynosi:

$$0,60372796$$

Korelacja jest już na skraju korelacji silnej i korelacji słabej.

- Porównamy teraz USA oraz Polskę, przypuszczam że korelacja będzie słaba, ponieważ po Polsce wirus nie rozprzestrzenił się tak szybko z powodu wprowadzenia ograniczeń oraz z powodu mniejszej ilości zaludnienia. Sprawdźmy więc matematycznie czy nasz przypuszczenia są poprawne. W USA za pierwszy dzień weźmiemy 22 lutego, a ostatni. W Polsce natomiast pierwszy będzie 7 marca, a ostatni 3 maja. Nasz wynik to:

$$0,6353313598$$

Widzimy więc że mimo różnej gęstości zaludnienia korelacja zachodzi całkiem niezłe. Jest nawet większa niż między Czechami a Słowacją, co jest zaskakujące.

- Sprawdźmy teraz korelację między państwami należącymi do innych kontynentów. Weźmy za przykład Brazylię oraz Azerbejdżan. W Azerbejdżanie wyniki będziemy brać z dni od 13 marca do 3 maja, natomiast z Brazylii od 5 maja do 25 kwietnia. Nasza korelacja to:

$$0,1799759062$$

Jak widać korelacja jest fatalna, oznacza to że dane nie są zbyt dobre, wynika to prawdopodobnie z małej ilości robionych testów.

- Sprawdźmy jeszcze korelację między USA a Azerbejdżanem. Azerbejdżan będzie zawierał dni od 13 marca do 3 maja, a USA od 22 lutego do 13 kwietnia. Wynik to:

$$-0,08316804048$$

Oznacza to że korelacja praktycznie nie istnieje, jest bardzo bliska 0.

Z powyższych obliczeń można wnioskować że zachodzi korelacja między zachorowaniami między większością państw. Jednak model nie jest uniwersalny co widać między korelacją Brazylii a Azerbejdżanem oraz Azerbejdżanu oraz USA.

Zadanie zostało wykonane na wersji punktowa 10+4 ponieważ obliczono korelację dla większej ilości państw.