

PROJEKT JĘZYK R

zbiór danych: Video Game Sales [1]
autor: Marcin Belicki
numer indeksu: 273417

1. Opis zbioru danych

Zbiór danych zawiera listę gier ze sprzedażą przekraczającą 100 000 kopii. Dane zostały pozyskane ze strony vgchartz.com.

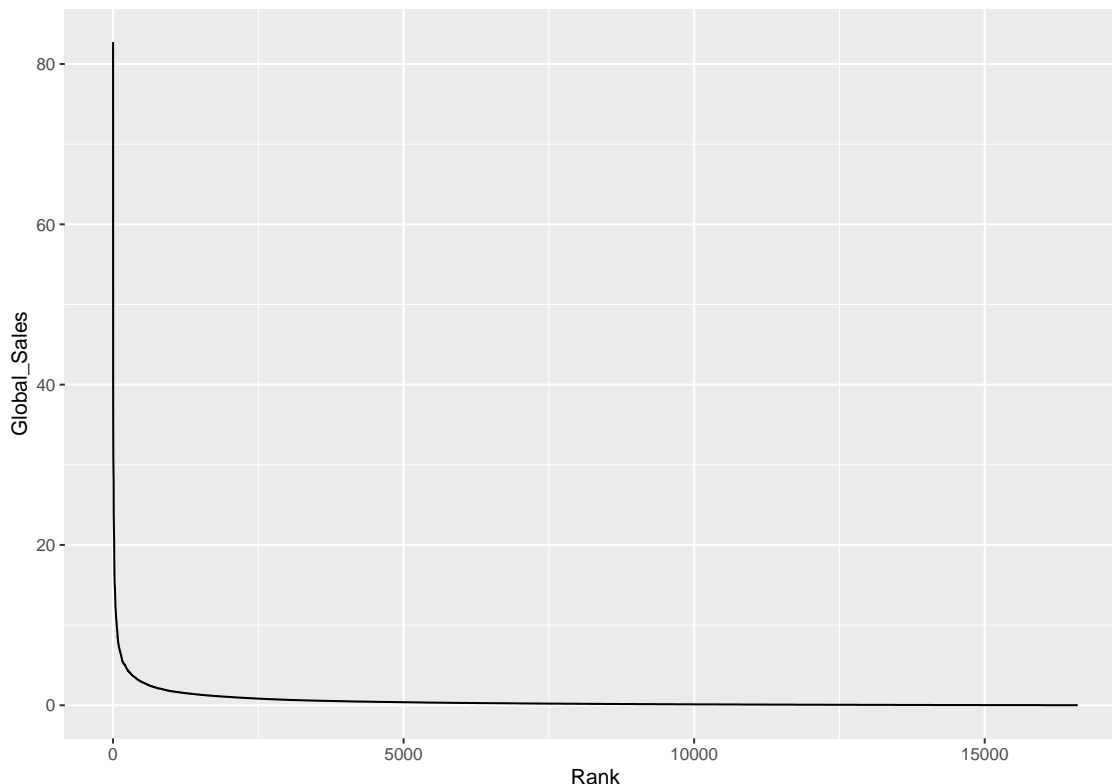
Opisy poszczególnych pól:

Rank	ranga gry pod względem sprzedaży ogółem
Name	nazwa gry
Platform	platforma, na którą gra została wydana (np. PC, PS4 itp.)
Year	rok wydania gry
Genre	gatunek gry
Publisher	wydawca gry
NA_Sales	liczba sprzedanych kopii w Ameryce Północnej (w milionach)
EU_Sales	liczba sprzedanych kopii w Europie (w milionach)
JP_Sales	liczba sprzedanych kopii w Japonii (w milionach)
Other_Sales	liczba sprzedanych kopii w reszcie świata (w milionach)
Global_Sales	liczba sprzedanych kopii na całym świecie (w milionach) - zmienna główna

Import danych został do środowiska R zostł wykonany za pomocą kodu:

```
1 data = read.csv("vgsales.csv", sep = ",")
2
3 list2env(data, .GlobalEnv)
```

2. Ilustracja graficzna zmiennej głównej



Rysunek 1: Wykres zależności zmiennej głównej od rangi

Wykres został wygenerowany za pomocą kodu:

```
1 ggplot(  
2   data = data,  
3   aes(  
4     y = Global_Sales,  
5     x = Rank  
6   )  
7 )+  
8   geom_line()
```

3. Obliczenie podstawowych statystyk opisowych zmiennej głównej

Obliczenia zostały wykonane za pomocą kodu:

```
1 description_statistics = function (data) {  
2  
3   number_of_records = length(data)  
4   number_of_records_inverted = 1/number_of_records  
5  
6   mean = mean(data)  
7  
8   variance = var(data) * (number_of_records - 1) * number_of_records_  
   inverted  
9  
10  standard_deviation = variance ^ .5  
11  inside = (data - mean)/standard_deviation  
12  
13  assymetry = sum(inside^3)*number_of_records_inverted  
14  curtosis = sum(inside^4)*number_of_records_inverted - 3  
15  
16  c(min, first_quartil, median, third_quartil, max) %<-% fivenum (data)  
17  
18  interquartile_range = third_quartil - first_quartil  
19  
20  quartil_deviation = interquartile_range * .5  
21  
22  mean_quartil = (first_quartil + third_quartil) * .5  
23  
24  quartil_assymetry_coefficient = (mean_quartil - median) / quartil_  
   deviation  
25  
26  quartil_variance_coefficient = quartil_deviation / median * 100  
27  
28  list(  
29    number_of_records = number_of_records,  
30    mean = mean,  
31    standard_deviation = standard_deviation,  
32  
33    assymetry = assymetry,  
34    curtosis = curtosis,  
35  
36    min = min,  
37    first_quartil = first_quartil,  
38    median = median,  
39    third_quartil = third_quartil,  
40    max = max,  
41  
42    interquartile_range = interquartile_range,  
43    quartil_deviation = quartil_deviation,  
44    mean_quartil = mean_quartil,
```

```

45     quartil_assymetry_coefficient = quartil_assymetry_coefficient,
46     quartil_variance_coefficient = quartil_variance_coefficient
47 )
48 }
49 }
50
51 Global_Sales_stats = description_statistics(Global_Sales)
52
53 Global_Sales_stats

```

Uzyskane wyniki

Liczba prób

$$n(\text{Global_Sales}) = 16598$$

Średnia

$$\text{mean}(\text{Global_Sales}) = 0.5374407$$

Odchylenie standardowe

$$\sigma(\text{Global_Sales}) = 1.554981$$

Współczynnik asymetrii

$$A(\text{Global_Sales}) = 17.39907$$

Kurtoza

$$K(\text{Global_Sales}) = 603.7501$$

Wartość minimalna

$$\min(\text{Global_Sales}) = 0.01$$

Pierwszy kwartył

$$Q_1(\text{Global_Sales}) = 0.06$$

Mediana

$$M_e(\text{Global_Sales}) = 0.17$$

Trzeci kwartył

$$Q_3(\text{Global_Sales}) = 0.47$$

Wartość maksymalna

$$\max(\text{Global_Sales}) = 82.74$$

Rozstęp między kwartyłowy

$$IQR(\text{Global_Sales}) = 0.41$$

Odchylenie ćwiartkowe

$$Q(\text{Global_Sales}) = 0.205$$

Kwartył średni

$$\bar{Q}(\text{Global_Sales}) = 0.265$$

Kwartylny współczynnik asymetrii

$$A_k(\text{Global_Sales}) = 0.4634146$$

Kwartylny współczynnik zmienności

$$V_k(\text{Global_Sales}) = 120.5882$$

4. Dobór zmiennych (przynajmniej trzech), które mogą mieć wpływ na zmienną główną

Do analizy zostały wybrane 3 następujące zmienne.

```
1 NA_Sales
2 EU_Sales
3 JP_Sales
```

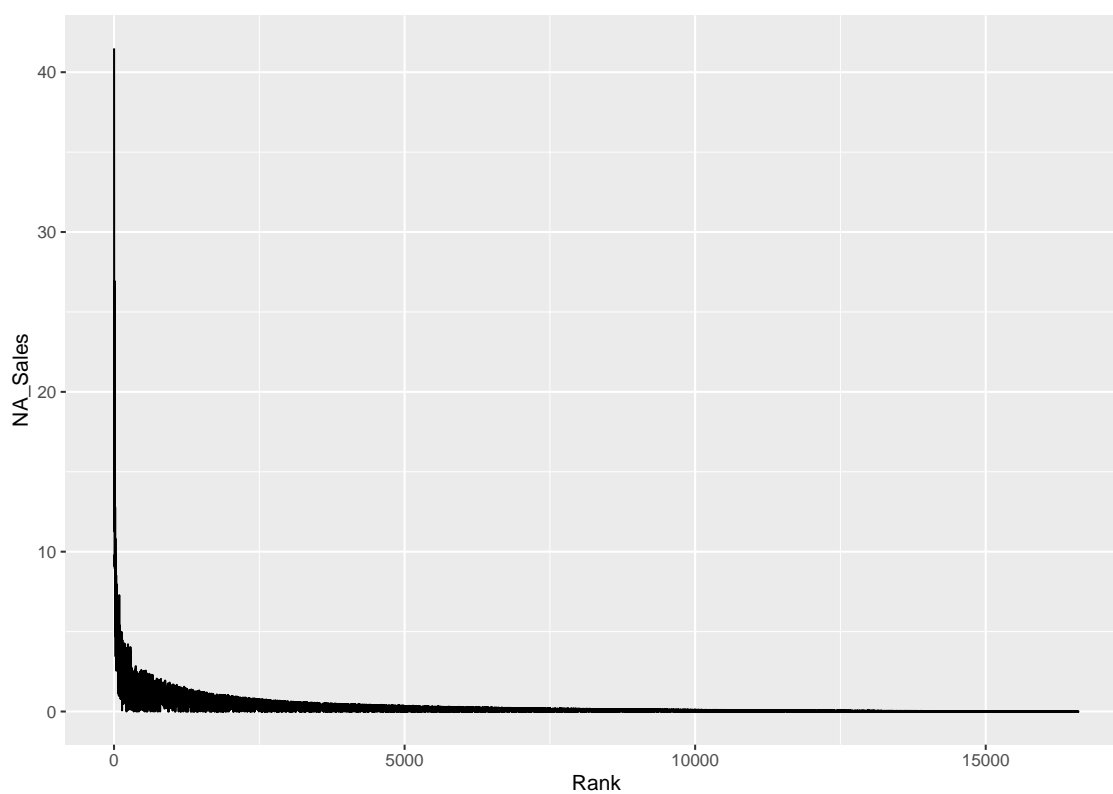
Te zmienne mogą pozwolić ustalić który region ma najważniejszy wpływ na globalną sprzedaż.

5. Graficzna prezentacja wybranych w poprzednim kroku zmiennych

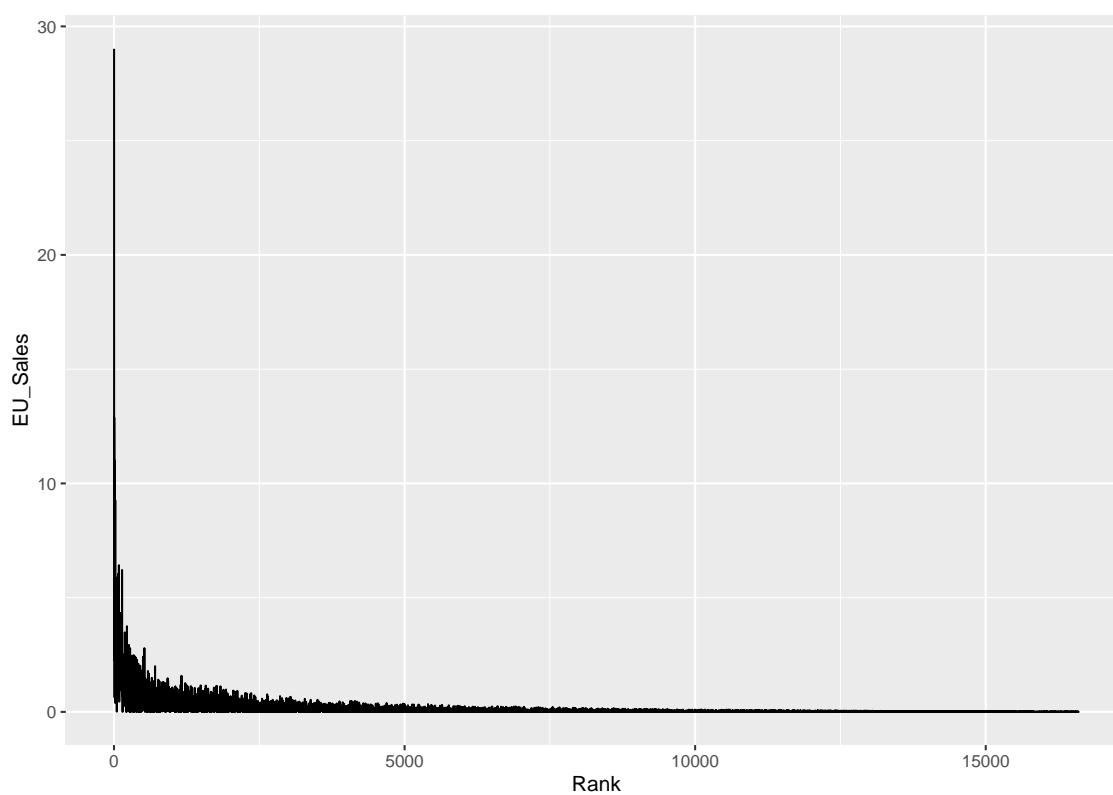
Na poniższych wykresach zostały przedstawione zmienne (wybrane w poprzednim punkcie) w zależności od rangi wiodącej.

Wykresy zostały uzyskane za pomocą poniższego kodu:

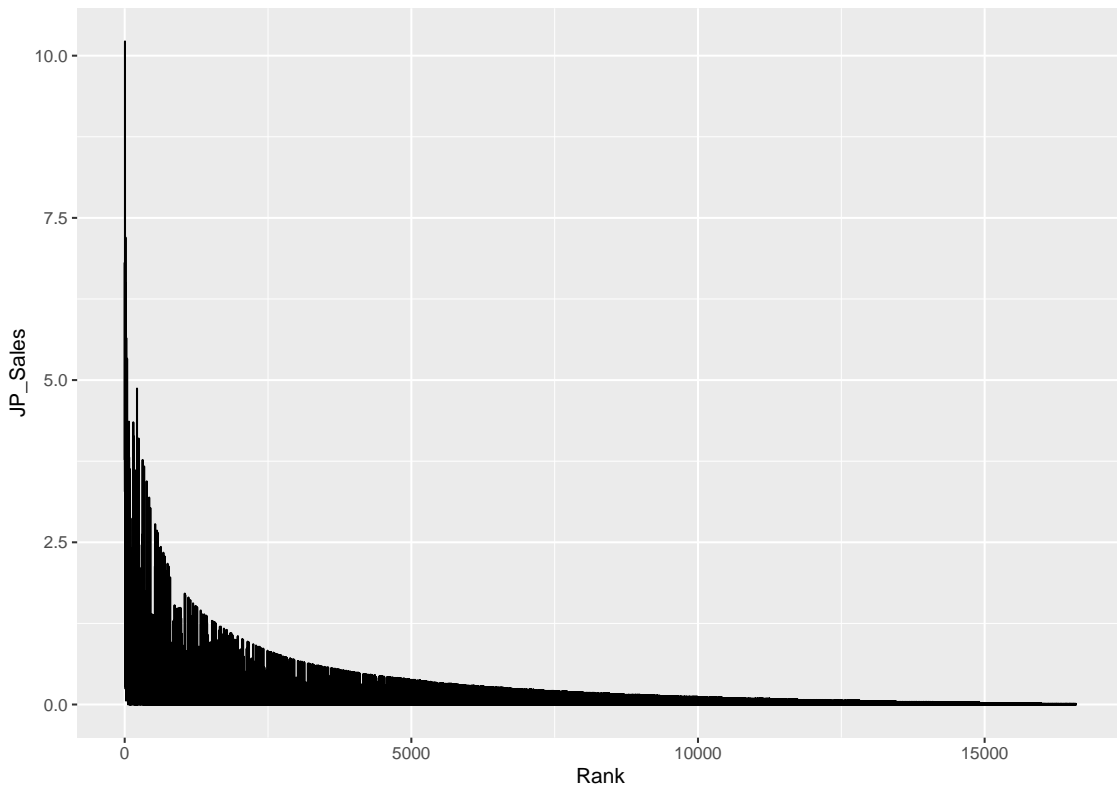
```
1 ggplot(  
2   data = data,  
3   aes(  
4     y = NA_Sales,  
5     x = Rank  
6   )  
7 )+  
8   geom_line()  
9  
10 ggplot(  
11   data = data,  
12   aes(  
13     y = EU_Sales,  
14     x = Rank  
15   )  
16 )+  
17   geom_line()  
18  
19 ggplot(  
20   data = data,  
21   aes(  
22     y = JP_Sales,  
23     x = Rank  
24   )  
25 )+  
26   geom_line()
```



Rysunek 2: Wykres zależności zmiennej `NA_Sales` od rangi



Rysunek 3: Wykres zależności zmiennej `EU_Sales` od rangi



Rysunek 4: Wykres zależności zmiennej JP_Sales od rangi

6. Charakterystyka powyższych zmiennych (statystyki opisowe lub rozkłady liczebności, w zależności od klasy zmiennych)

Wyniki przedstawione w tabeli zostały uzyskane za pomocą poniższego kodu:

```

1 NA_Sales_stats = description_statistics(NA_Sales)
2
3 EU_Sales_stats = description_statistics(EU_Sales)
4
5 JP_Sales_stats = description_statistics(JP_Sales)
6
7 NA_Sales_stats
8
9 EU_Sales_stats
10
11 JP_Sales_stats

```

data	NA_Sales	EU_Sales	JP_Sales
$n(\text{data})$	16598	16598	16598
$\text{mean}(\text{data})$	0.2646674	0.146652	0.07778166
$\sigma(\text{data})$	0.8166584	0.505336	0.3092813
$A(\text{data})$	18.79793	18.87383	11.20545
$K(\text{data})$	648.9344	755.7997	194.1751
$\min(\text{data})$	0	0	0
$Q_1(\text{data})$	0	0	0
$M_e(\text{dat})$	0.08	0.02	0
$Q_3(\text{data})$	0.24	0.11	0.04
$\max(\text{data})$	41.49	29.02	10.22
$IQR(\text{data})$	0.24	0.11	0.04
$Q(\text{data})$	0.12	0.055	0.02
$Q(\text{data})$	0.12	0.055	0.02
$A_k(\text{data})$	0.3333333	0.6363636	1
$V_k(\text{data})$	150	275	Inf

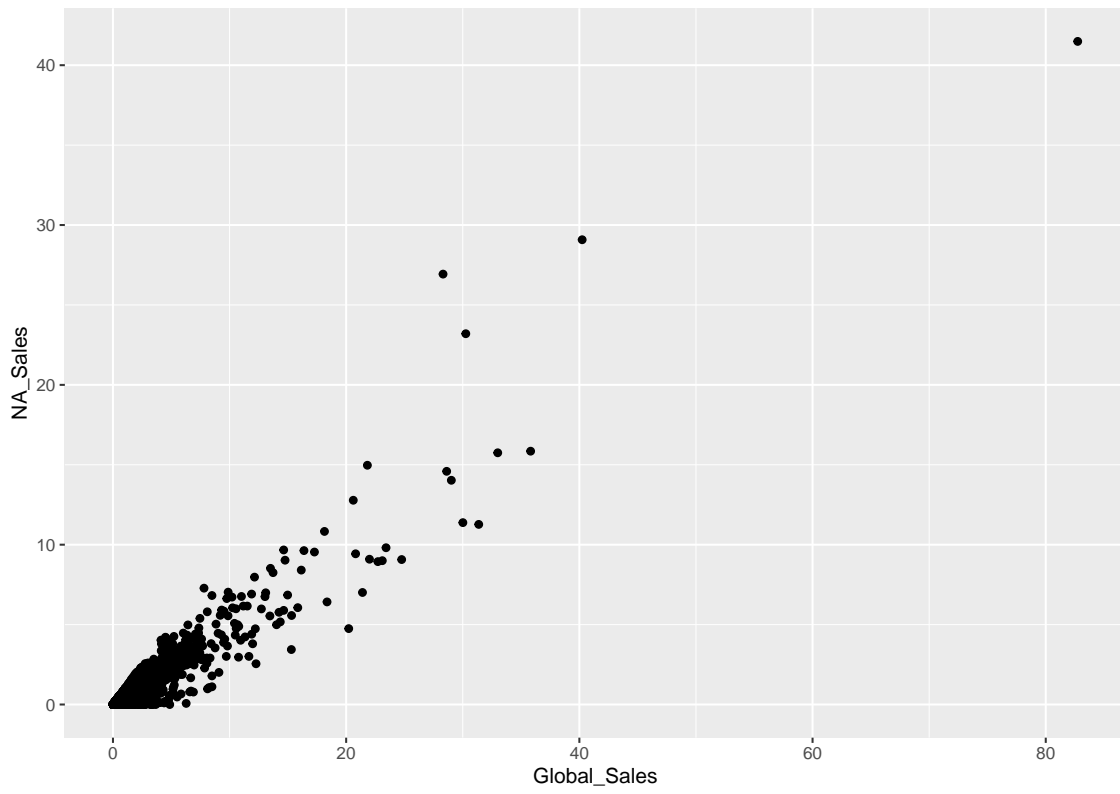
7. Graficzna prezentacja zależności

Wykresy zostały uzyskane za pomocą poniższego kodu:

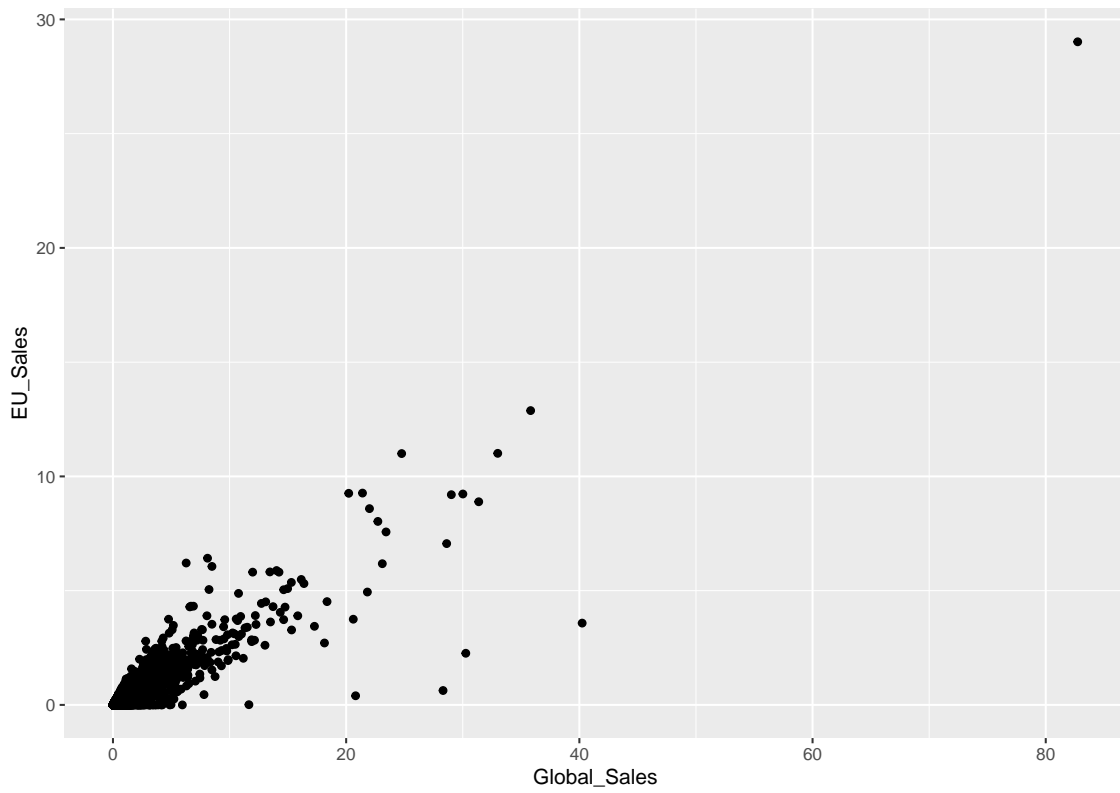
```

1 ggplot(
2   data = data,
3   aes(
4     y = NA_Sales,
5     x = Global_Sales
6   )
7 )+
8   geom_point()
9
10 ggplot(
11   data = data,
12   aes(
13     y = EU_Sales,
14     x = Global_Sales
15   )
16 )+
17   geom_point()
18
19 ggplot(
20   data = data,
21   aes(
22     y = JP_Sales,
23     x = Global_Sales
24   )
25 )+
26   geom_point()

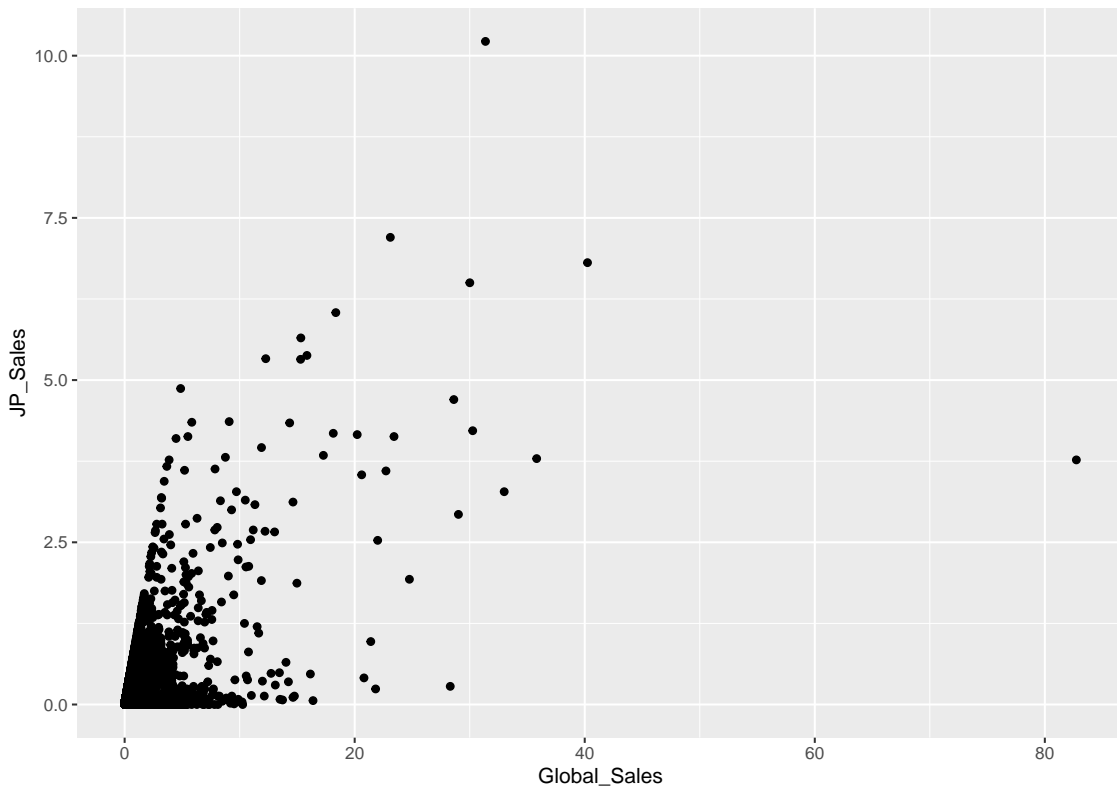
```

Rysunek 5: Wykres zależności zmiennej `NA_Sales` od `Global_Sales`



Rysunek 6: Wykres zależności zmiennej `EU_Sales` od `Global_Sales`



Rysunek 7: Wykres zależności zmiennej JP_Sales od Global_Sales

8. Wykonanie odpowiedniego testu statystycznego, który potwierdzi lub odrzuci hipotezę o zależności

W celu sprawdzenia korelacji między wybranymi zmiennymi a zmienną główną przeprowadzony został test Spearmana za pomocą poniższego kodu.

```
1 spearman_test = function(data) {
2   cor.test(Global_Sales, data, method = "spearman", exact = F)
3 }
4 spearman_test(NA_Sales)
5
6 spearman_test(EU_Sales)
7
8 spearman_test(JP_Sales)
```

Współczynniki dla poszczególnych przypadków przedstawiają się w następujący sposób:

$$\rho(\text{Global_Sales}, \text{NA_Sales}) = 0.7955717$$

$$\rho(\text{Global_Sales}, \text{EU_Sales}) = 0.6968457$$

$$\rho(\text{Global_Sales}, \text{JP_Sales}) = 0.1519311$$

9. Wnioski

Z wyników testu Spearmana wynika, że najbardziej skorelowaną zmienną ze zmienną główną jest NA_Sales. Wykres przedstawiony na rysunku 7 również zdaje się potwierdzać tę hipotezę (ze względu na widoczną niemalże liniową zależność). Ze statystyk opisanych w punkcie 6 można wywnioskować, że najmniej ze wszystkich gier w tym zestawieniu sprzedało się na rynku

japońskim. Wykres przedstawiony na rysunku 4 przedstawia najbardziej wahające się wartości ze wszystkich wybranych zmiennych. Może być to spowodowane niedostępnością niektórych z gier na rynku japońskim.

Wszystkie z przedstawionych zmiennych zdają się (w uśrednieniu) rosnać hiperbolicznie względem rangi.

Bibliografia

- [1] <https://www.kaggle.com/datasets/gregorut/videogamesales>
[dostęp: 13.06.2022 22:37 GMT+2]