

# Projekt WdAD

Filip Cebula 151410

19 grudnia 2024

# 1 Zadanie 1

## 1.1 Podpunkt A

Zapisujemy dane do naszego zadania.

```
1 earnings = c(45617,7166,18594,2236,1278,19828,4033,28151,2414,3800);
2 earnings.sd = 15000;
```

Estymujemy średnią z próbki, oraz błąd standardowy (korzystając z danego odchylenia standardowego).

```
1 earnings.mean = mean(earnings);
2 earnings.se = earnings.sd / sqrt(length(earnings));
3
4 # Output
5 # Earnings mean: 13311.7
6 # Earnings standard error: 4743.416
```

Obliczamy przedziały ufności 90%, przy założeniu rozkładu normalnego próbki.

```
1 earnings.interval1 = 0.1;
2 earnings.ci1.margin = qnorm(p = (1 - earnings.interval1 / 2)) * earnings.se;
3 earnings.ci1 = earnings.mean + c(-1,1) * earnings.ci1.margin;
4
5 # Output
6 # "Confidence interval for 0.9 assuming normal distribution:"
7 # 5509.474 21113.926
```

Obliczamy przedziały ufności 95%, przy założeniu rozkładu normalnego próbki.

```
1 earnings.interval2 = 0.05;
2 earnings.ci2.margin = qnorm(p = (1 - earnings.interval2 / 2)) * earnings.se;
3 earnings.ci2 = earnings.mean + c(-1,1) * earnings.ci2.margin;
4
5 # Output
6 # "Confidence interval for 0.95 assuming normal distribution:"
7 # 4014.775 22608.625
```

## 1.2 Podpunkt B

Odrzucamy dane nam wcześniej odchylenie standardowe i estymujemy nowe odchylenie standardowe, oraz błąd standardowy na podstawie danej próbki.

```
1 earnings.sd = sd(earnings);
2 earnings.se = earnings.sd / sqrt(length(earnings));
3
4 # Output
5 # Earnings standard derivation: 14662.04
6 # Earnings standard error: 4636.545
```

Używamy rozkładu t Studenta, żeby obliczyć przedziały ufności dla naszej próbki.

Przedział ufności 90%.

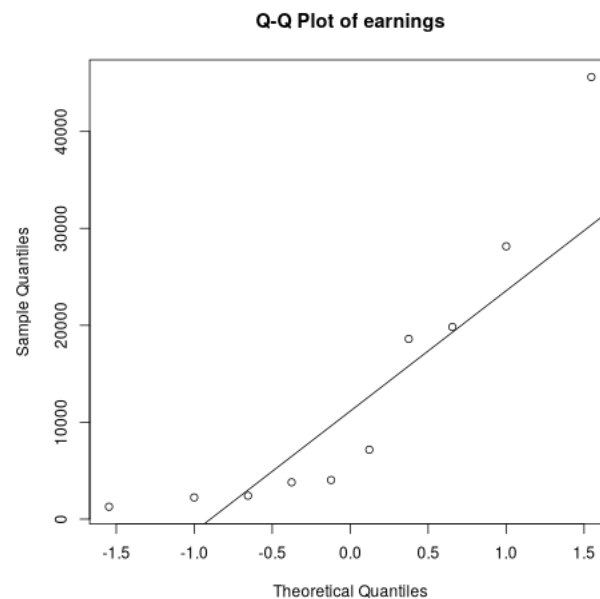
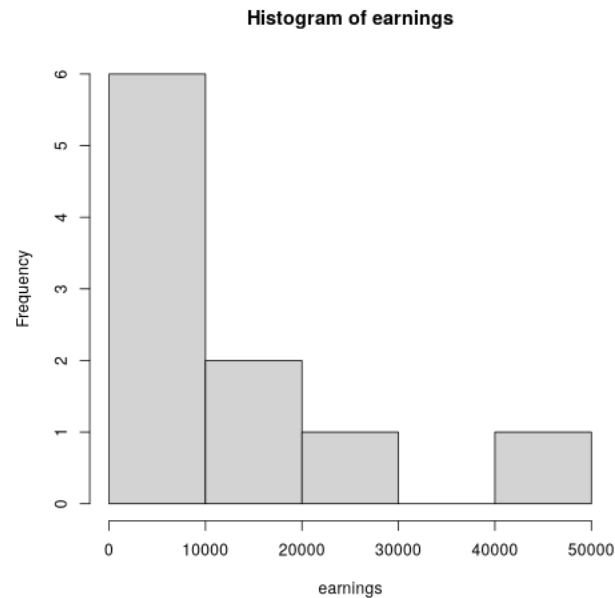
```
1 earnings.interval1 = 0.1;
2 earnings.ci1.margin = qt(p = (1 - earnings.interval1 / 2),
3                           df = length(earnings) - 1) * earnings.se;
4 earnings.ci1 = earnings.mean + c(-1,1) * earnings.ci1.margin;
5
6 # Output
7 # "Confidence interval for 0.9 assuming Student's t-distribution:"
8 # 4812.39 21811.01
```

Przedział ufności 95%.

```
1 earnings.interval2 = 0.05;
2 earnings.ci2.margin = qt(p = (1 - earnings.interval2 / 2),
3                           df = length(earnings) - 1) * earnings.se;
4 earnings.ci2 = earnings.mean + c(-1,1) * earnings.ci2.margin;
5
6 # Output
7 # "Confidence interval for 0.95 assuming Student's t-distribution:"
8 # 2823.107 23800.293
```

### 1.3 Podpunkt C

Odrzucamy założenie o normalności rozkładu naszej próbki. Możemy użyć histogramu, oraz wykresu kwantyl-kwantyl, aby sprawdzić jak blisko rozkładu normalnego jest rozkład naszej próbki.



Na podstawie powyższych wykresów, a szczególnie wykresu kwantyl-kwantyl, możemy zauważyć, że rozkład naszej próbki nie jest rozkładem normalnym. Wynika z tego, że obliczone przez nas powyżej przedziały ufności nie są dokładne. Aby otrzymać dokładniejsze wyniki, użyjemy metody bootstrap do obliczenia przedziałów ufności.

Tworzymy funkcję bootstrap i używamy jej na naszej próbce.

```
1 my_bootstrap = function(data) {  
2   n = length(data);  
3   means = c();  
4  
5   for (i in 1:10000) {  
6     rands = sample(1:n, n, replace = T);  
7     xs = data[rands];  
8  
9     means = append(means, mean(xs));  
10  }  
11  
12  return(means);  
13 }  
14  
15 earnings.bootstrap = my_bootstrap(earnings);
```

Wyliczamy przedział ufności 90%.

```
1 quantile(earnings.bootstrap, probs = c(0.05, 0.95));  
2  
3 # Output  
4 # "Confidence interval for 0.9:"  
5 #      5%      95%  
6 # 6641.84 20903.85
```

Wyliczamy przedział ufności 95%.

```
1 quantile(earnings.bootstrap, probs = c(0.025, 0.975));  
2  
3 # Output  
4 # "Confidence interval for 0.95:"  
5 #      2.5%      97.5%  
6 # 5638.485 22847.100
```

## 2 Zadanie 2

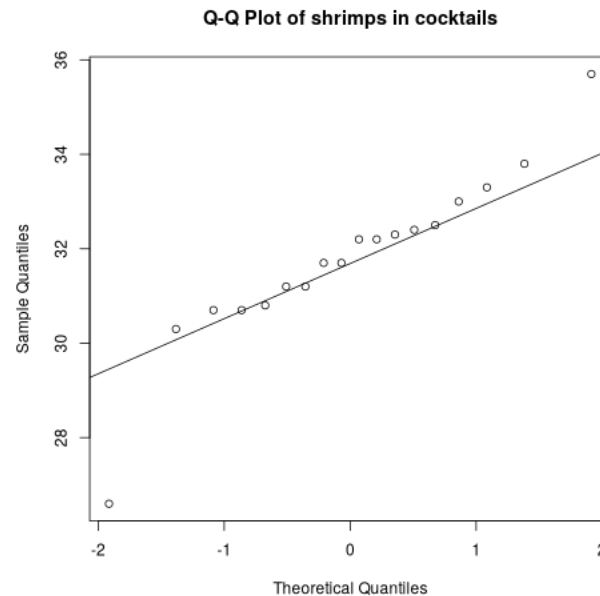
W tym zadaniu do zbadania wybrałem zestawy danych: shrimp, Sitka89 i quine.

Najpierw dodałem bibliotekę MASS, a następnie napisałem drugą metodę bootstrap, która ułatwi mi liczenie przedziałów ufności dla odchylenia standardowego i wariancji.

```
1 library(MASS);
2
3 adv_bootstrap = function(data, f) {
4   n = length(data);
5   result = c();
6
7   for (i in 1:10000) {
8     rand = sample(1:n, n, replace = T);
9     xs = data[rand];
10
11     result = append(result, f(xs));
12   }
13
14   return(result);
15 }
```

## 2.1 Shrimp

Shrimp to zestaw danych, który opisuje ilość krewetek (procent całkowitej masy) w koktajlu krewetkowym.



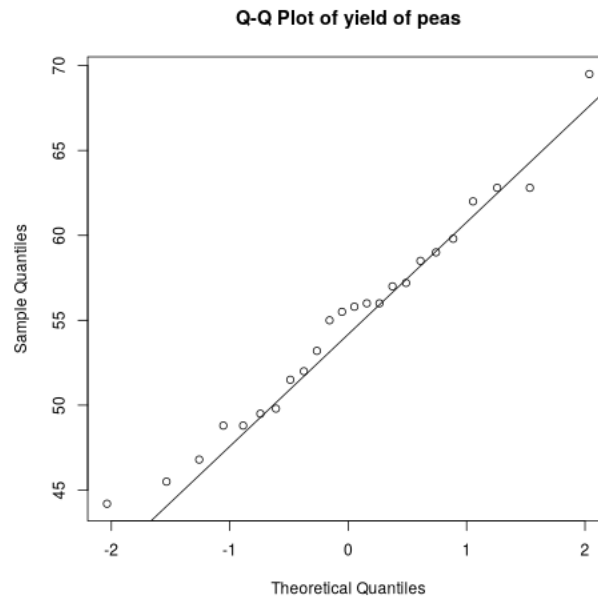
Na powyższym wykresie możemy zauważyć, że poza dwoma wartościami skrajnymi, rozkład naszych danych jest bardzo zbliżony do rozkładu normalnego, jednak o ciężkich ogonach (fat tails), co oznacza zbyt dużą kurtozę.

Obliczmy zatem przedziały ufności 95%, dla średniej, odchylenia standardowego i wariancji.

```
1 shrimps.bootstrap.mean = adv_bootstrap(shrimp, mean);
2 shrimps.bootstrap.sd = adv_bootstrap(shrimp, sd);
3 shrimps.bootstrap.var = adv_bootstrap(shrimp, var);
4
5 # Output
6 # "Confidence intervals for 0.95 for shrimps:"
7 # "Mean:"
8 #   2.5%   97.5%
9 # 30.92222 32.57222
10 # "Standard deviation:"
11 #   2.5%   97.5%
12 # 0.8732765 2.5980592
13 # "Standard error:"
14 #   2.5%   97.5%
15 # 0.7527917 6.7226168
```

## 2.2 Npk

Npk to zestaw danych z eksperymentu, który testował wpływ azotu, fosforu i potasu na wzrost groszku. Ja skupiłem się na statystyce yield, która zawiera informacje o tym, ile grochu zebrano na jednej działce.



Na powyższym wykresie widzimy, że rozkład naszych danych, jest bardzo podobny do rozkładu normalnego. Jednak nasz rozkład ma bardzo lekkie ogony, co oznacza ujemną kurtozę.

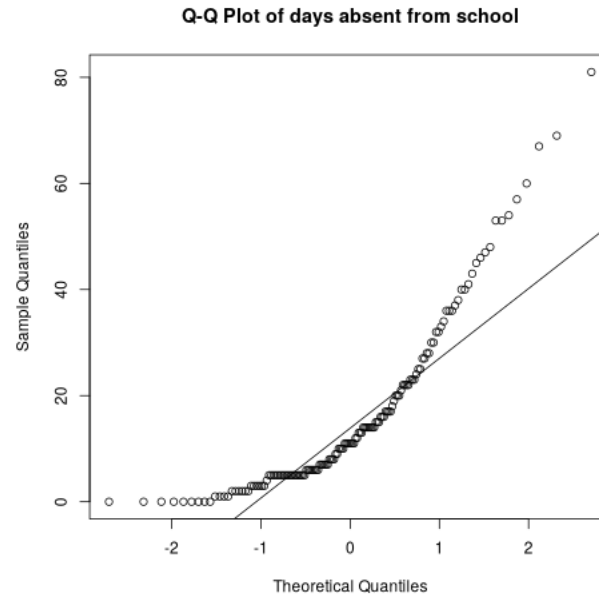
Obliczmy zatem przedziały ufności 95%, dla średniej, odchylenia standardowego i wariancji.

```
1 npkyield.bootstrap.mean = adv_bootstrap(npk$yield, mean);
2 npkyield.bootstrap.sd = adv_bootstrap(npk$yield, sd);
3 npkyield.bootstrap.var = adv_bootstrap(npk$yield, var);
4
5 # Output
6 # "Confidence intervals for 0.95 for npk yield:"
7 # "Mean:"
8 #   2.5%   97.5%
9 # 52.4625 57.3500
10 # "Standard deviation:"
11 #   2.5%   97.5%
12 # 4.422683 7.627811
13 # "Standard error:"
14 #   2.5%   97.5%
15 # 19.73285 57.91169
```



## 2.3 Quine

Quine to zestaw danych, który zawiera informacje o uczniach z Australii. Ja skupiłem się na statystyce, która przedstawi ile dni dane dziecko opuściło w trakcie roku szkolnego.



Na powyższym wykresie możemy zauważyć, że nasz rozkład znacząco odbiega od rozkładu normalnego, ponieważ nasz rozkład jest zbyt pozytywnie skośny (tj. pochyla się na lewą stronę).

Obliczmy zatem przedziały ufności 95%, dla średniej, odchylenia standardowego i wariancji.

```
1 quineabsent.bootstrap.mean = adv_bootstrap(quine$Days, mean);
2 quineabsent.bootstrap.sd = adv_bootstrap(quine$Days, sd);
3 quineabsent.bootstrap.var = adv_bootstrap(quine$Days, var);
4
5 # Output
6 # "Confidence intervals for 0.95 for absent days:"
7 # "Mean:"
8 #   2.5%   97.5%
9 # 13.96575 19.13031
10 # "Standard deviation:"
11 #   2.5%   97.5%
12 # 13.43493 18.74195
13 # "Standard error:"
14 #   2.5%   97.5%
15 # 181.1858 353.6281
```