

Analiza danych ankietowych - sprawozdanie 1

"Lista zadań nr 1"

Przedmiot i prowadzący:
Analiza danych ankietowych,
poniedziałki 9.15 - 11.00 (grupa nr 2),
Inż. Hubert Woszczek

Aleksandra Hodera (268733)

Aleksandra Polak (268786)

Spis treści

1	Wstęp	3
2	Użyte biblioteki	3
3	Część I	3
3.1	Zadanie 1	3
3.2	Zadanie 1.1	4
3.3	Zadanie 1.2	5
3.4	Zadanie 1.3	5
3.5	Zadanie 1.4	7
3.6	Zadanie 1.5	10
3.7	Zadanie 1.6	12
3.8	Zadanie 1.7	13
3.9	Zadanie 1.8	13
4	Część II	16
4.1	Zadanie 2	16
4.2	Zadanie 3	20
4.3	Zadanie 4	22
4.4	Zadanie 5	25
5	Część III i IV	27
5.1	Zadanie 6	27
5.2	Zadanie 7	28
5.3	Zadanie 8	28
5.4	Zadanie 9	29
6	Część V	37
6.1	Zadanie 10	37
6.2	Zadanie 11	38
6.3	Zadanie 12	41
7	Zadanie dodatkowe	45
7.1	Zadanie *1	45
8	Źródła	48

1 Wstęp

W poniższym sprawozdaniu zostały przedstawione wyniki listy zadań nr 1 przygotowanej w ramach laboratoriów z analizy danych ankietowych prowadzonych przez inż. Huberta Woszczek do wykładu dr inż. Aleksandry Grzesiek.

2 Użyte biblioteki

W tym punkcie zostały przedstawione wszystkie biblioteki, które użyliśmy podczas tworzenia raportu:

```
1 library(ggplot2)
2 library(tidyr)
3 library(dplyr)
4 library(kableExtra)
5 library(vcd)
6 library(knitr)
7 library(likert)
8 library(reshape2)
9 library(stats)
10 library(binom)
11 library(GenBinomApps)
```

3 Część I

3.1 Zadanie 1

W pewnej dużej agencji reklamowej przeprowadzono ankietę mającą na celu ocenę poziomu satysfakcji z pracy. Wzięło w niej udział dwieście losowo wybranych osób (losowanie proste ze zwracaniem). W pliku "ankieta.csv" umieszczono odpowiedzi na kilka z zadanych pytań:

- "W jakim dziale jesteś zatrudniony? zmienna **DZIAŁ** przyjmująca wartości: **HR** (Dział obsługi kadrowo-płacowej), **IT** (Dział utrzymania sieci i systemów informatycznych), **DK** (Dział Kreatywny) lub **DS** (Dział Strategii).
- "Jak długo pracujesz w firmie? zmienna **STAŻ** przyjmująca wartości: **1** (poniżej jednego roku), **2** (między jednym rokiem a trzema latami) lub **3** (powyżej trzech lat).
- "Czy pracujesz na stanowisku menedżerskim? zmienna **CZY_KIER** przyjmująca wartości: **Tak** (stanowisko menedżerskie) lub **Nie** (stanowisko inne niż menedżerskie).
- "Jak bardzo zgadzasz się ze stwierdzeniem, że firma pozwala na elastyczne godziny pracy tym samym umożliwiając zachowanie równowagi między pracą, a życiem prywatnym? zmienna **PYT_1** przyjmująca wartości: **-2** (zdecydowanie się nie zgadzam), **-1** (nie zgadzam się), **0** (nie mam zdania), **1** (zgadzam się), **2** (zdecydowanie się zgadzam).
- "Jak bardzo zgadzasz się ze stwierdzeniem, że twoje wynagrodzenie adekwatnie odzwierciedla zakres wykonywanych przez ciebie obowiązków? zmienna **PYT_2** przyjmująca wartości: **-2** (zdecydowanie się nie zgadzam), **-1** (nie zgadzam się), **1** (zgadzam się), **2** (zdecydowanie się zgadzam).

Dodatkowo w ramach metryczki ankietowani zostali poproszeni o wskazanie swojego wieku - zmienna **WIEK** przyjmująca wartości numeryczne, oraz wskazanie płci - zmienna **PŁEĆ** przyjmująca wartość **Kobieta** lub **Mężczyzna**.

Kilka tygodni później przeprowadzono rewizję wynagrodzeń, w wyniku której część pracowników otrzymała podwyżki. Ankietowanych biorących udział w badaniu poproszono wówczas o ponowną odpowiedź na pytanie dotyczące zadowolenia z wynagrodzenia - zmienna **PYT_3**.

3.2 Zadanie 1.1

Wczytaj dane i przygotuj je do analizy. Zadbaj o odpowiednie typy zmiennych, zweryfikuj czy przyjmują wartości zgodne z powyższym opisem, zbadaj czy nie występują braki w danych.

```
1 # wczytanie danych
2 data <- read.csv("ankieta.csv", fileEncoding = "Latin1", sep=";", na=c("
  "))
3
4 # zmiana nazw na zgodne z poleceniem
5 colnames(data) <- c('DZIAŁ', 'STAŻ', 'CZY_KIER', 'PYT_1', 'PYT_2', 'PYT_3',
  'PŁEĆ', 'WIEK')
6
7 # zbadanie czy dane mają odpowiednie typy
8 types <- sapply(data, class)
9 types_df <- data.frame(Zmienna = names(types), Typ = types, row.names =
  NULL)
10
11 # stworzenie tabeli z typami poszczególnych zmiennych
12 kable(types_df,
13       caption = "Typy zmiennych w kolumnach",
14       align = "c",
15       booktabs = TRUE) %>%
16 kable_styling(latex_options = c("striped", "hold_position"))
17
18 # pomocnicza zamiana na factor
19 data$PYT_1 <- factor(data$PYT_1, levels = c(-2, -1, 0, 1, 2),
20                     labels = c("zdecydowanie się nie zgadzam", "nie
  zgadzam się", "nie mam zdania", "zgadzam się", "zdecydowanie się
  zgadzam"))
21 data$PYT_2 <- factor(data$PYT_2, levels = c(-2, -1, 1, 2),
22                     labels = c("zdecydowanie się nie zgadzam", "nie
  zgadzam się", "zgadzam się", "zdecydowanie się zgadzam"))
23 data$PYT_3 <- factor(data$PYT_3, levels = c(-2, -1, 1, 2),
24                     labels = c("zdecydowanie się nie zgadzam", "nie
  zgadzam się", "zgadzam się", "zdecydowanie się zgadzam"))
```

Otrzymałyśmy tabelę, w której podane są typy danych.

ZMIENNA	TYP
DZIAŁ	character
STAŻ	integer
CZY_KIER	character
PYT_1	integer
PYT_2	integer
PYT_3	integer
PŁEĆ	character
WIEK	integer

Tabela 1: Typy danych

Jak możemy zauważyć w tabeli 1, zmienne **DZIAŁ**, **CZY_KIER** i **PŁEĆ** są zmiennymi kategorycznymi, natomiast **STAŻ**, **PYT_1**, **PYT_2**, **PYT_3** oraz **WIEK** to zmienne liczbowe, co jest zgodne z naszymi przypuszczeniami.

```
1 # zbadanie występowania braku danych
2 colSums(is.na(data))
```

Jako wynik działania powyższego kodu otrzymaliśmy tabelę prezentującą liczbę braków danych dla konkretnych zmiennych.

DZIAŁ	STAŻ	CZY_KIER	PYT_1	PYT_2	PYT_3	PŁEĆ	WIEK
0	0	0	0	0	0	0	0

Tabela 2: Występowanie braków danych

Jak możemy zauważyć w tabeli 2 w naszych danych nie obserwujemy braków.

3.3 Zadanie 1.2

Utwórz zmienną **WIEK_KAT** przeprowadzając kategoryzację zmiennej **WIEK** korzystając z następujących przedziałów: do 35 lat, między 36 a 45 lat, między 46 a 55 lat, powyżej 55 lat.

```
1 # Podział wieku na odpowiednie grupy
2 df <- mutate(data, WIEK_KAT = cut(WIEK, breaks = c(0, 35, 45, 55, max(
  WIEK)),
3                                     labels = c("0-35", "36-45", "46-55", "
  56+"))))
4 View(df)
```

3.4 Zadanie 1.3

Sporządź tablice liczości dla zmiennych: **DZIAŁ**, **STAŻ**, **CZY_KIER**, **PŁEĆ**, **WIEK_KAT**.

```
1 generate_frequency_table <- function(df, variable) {
2
3 # Tworzenie tablicy liczości
4   result <- table(df[[variable]])
5
6 # Wyświetlenie informacji o brakujących danych
```

```

7 kable(result,
8       caption = paste("Tablica liczności dla zmiennej", variable),
9       align = "c",
10      col.names = c(variable, "Liczba wystąpień"),
11      booktabs = TRUE) %>%
12 kable_styling(latex_options = c("striped", "hold_position"))
13 }
14
15 # Wyświetlanie tablic
16 generate_frequency_table(df, "DZIAŁ")
17 generate_frequency_table(df, "STAŻ")
18 generate_frequency_table(df, "CZY_KIER")
19 generate_frequency_table(df, "PŁEĆ")
20 generate_frequency_table(df, "WIEK_KAT")

```

DZIAŁ	Liczba wystąpień
DK	98
DS	45
HR	31
IT	26

Tabela 3: Tablica liczności dla zmiennej DZIAŁ

STAŻ	Liczba wystąpień
1	41
2	140
3	19

Tabela 4: Tablica liczności dla zmiennej STAŻ

CZY_KIER	Liczba wystąpień
Nie	173
Tak	27

Tabela 5: Tablica liczności dla zmiennej CZY_KIER

PŁEĆ	Liczba wystąpień
K	71
M	129

Tabela 6: Tablica liczności dla zmiennej PŁEĆ

WIEK_KAT	Liczba wystąpień
0-35	26
36-45	104
45-55	45
56+	25

Tabela 7: Tablica liczności dla zmiennej WIEK_KAT

W tabelach 3 - 7 widoczne są liczby wystąpień poszczególnych opcji, dla zmiennych **DZIAŁ** (tabela 3), **STAŻ** (tabela 4), **CZY_KIER** (tabela 5), **PŁEĆ** (tabela 6) oraz **WIEK_KAT** (tabela 7).

3.5 Zadanie 1.4

Sporządź wykresy kołowe oraz wykresy słupkowe dla zmiennych: **PYT_1** oraz **PYT_2**.

```

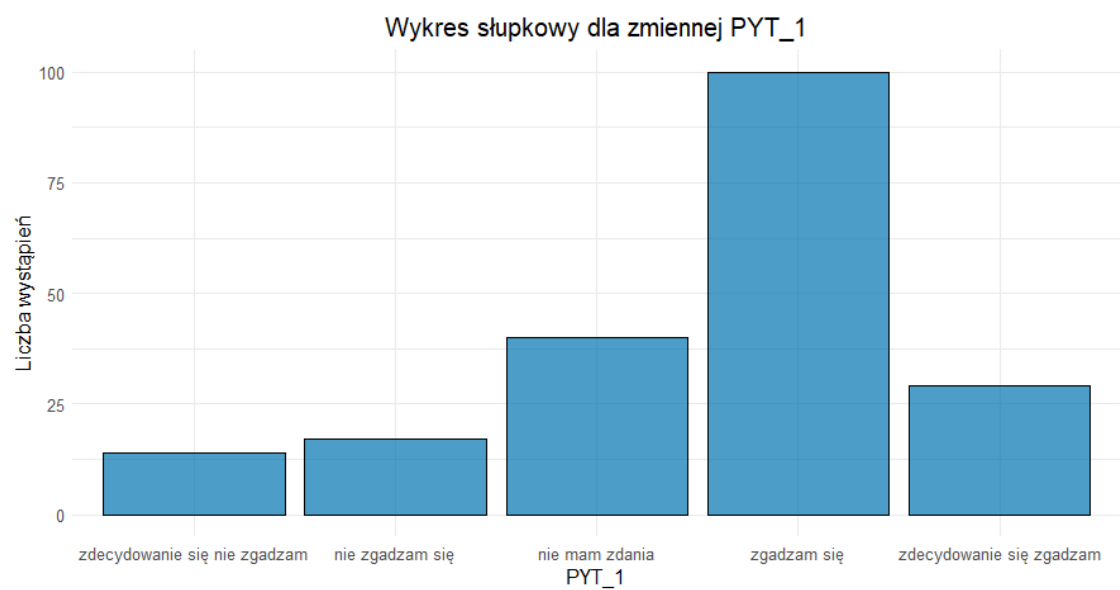
1 # ----- PYT_1 -----
2 # Wykres słupkowy dla pytania 1
3 ggplot(df, aes(x = PYT_1)) +
4   geom_bar(fill = "#0072B2", color = "black", alpha = 0.7) +
5   labs(title = "Wykres słupkowy dla zmiennej PYT_1", x = "PYT_1", y = "
6     Liczba wystąpień") +
7   theme_minimal() +
8   theme(
9     axis.title = element_text(size = 12),
10    axis.text = element_text(size = 10),
11    plot.title = element_text(size = 15, hjust = 0.5))
12
13 # Wykres kołowy dla pytania 1
14 ggplot(data = df, aes(x = "", fill = PYT_1)) +
15   geom_bar(color = "white", position="fill", alpha = 0.75) +
16   coord_polar("y", start=pi/2) +
17   scale_fill_viridis_d(option = "plasma", na.value = "gray50") +
18   labs(title = "Rozkład odpowiedzi dla zmiennej PYT_1", fill = "
19     Odpowiedzi") +
20   theme_void() +
21   theme(
22     plot.title = element_text(hjust = 0.5, size = 15),
23     axis.text = element_blank(),
24     axis.title = element_blank(),
25     legend.title = element_text(size = 12),
26     legend.text = element_text(size = 10)
27   )
28
29 # ----- PYT_2 -----
30 # Wykres słupkowy dla pytania 2
31 ggplot(df, aes(x = PYT_2)) +
32   geom_bar(fill = "#0072B2", color = "black", alpha = 0.7) +
33   labs(title = "Wykres słupkowy dla zmiennej PYT_2", x = "PYT_2", y = "
34     Liczba wystąpień") +
35   theme_minimal() +
36   theme(
37     axis.title = element_text(size = 12),
38     axis.text = element_text(size = 10),
39     plot.title = element_text(size = 15, hjust = 0.5))

```

```

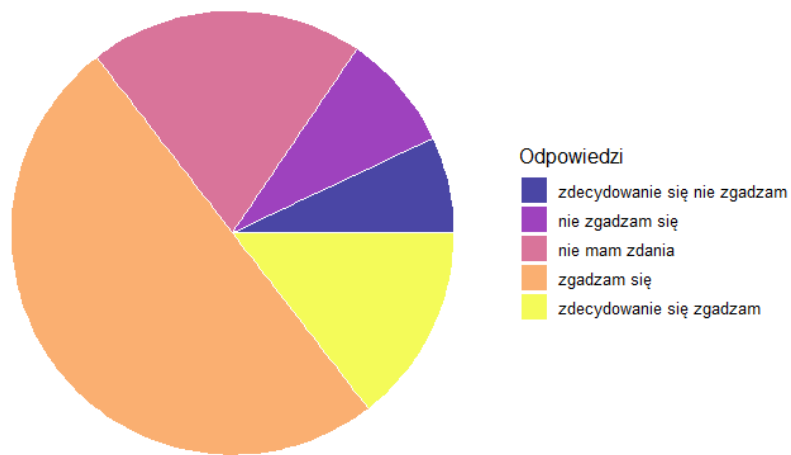
12 # Wykres kołowy dla pytania 2
13 ggplot(data = df, aes(x = "", fill = PYT_2)) +
14   geom_bar(color = "white", position="fill", alpha = 0.75) +
15   coord_polar("y", start=pi/2) +
16   scale_fill_viridis_d(option = "plasma", na.value = "gray50") +
17   labs(title = "Rozkład odpowiedzi dla zmiennej PYT_2", fill = "
18     Odpowiedzi") +
19   theme_void() +
20   theme(
21     plot.title = element_text(hjust = 0.5, size = 15),
22     axis.text = element_blank(),
23     axis.title = element_blank(),
24     legend.title = element_text(size = 12),
25     legend.text = element_text(size = 10)
26 )

```

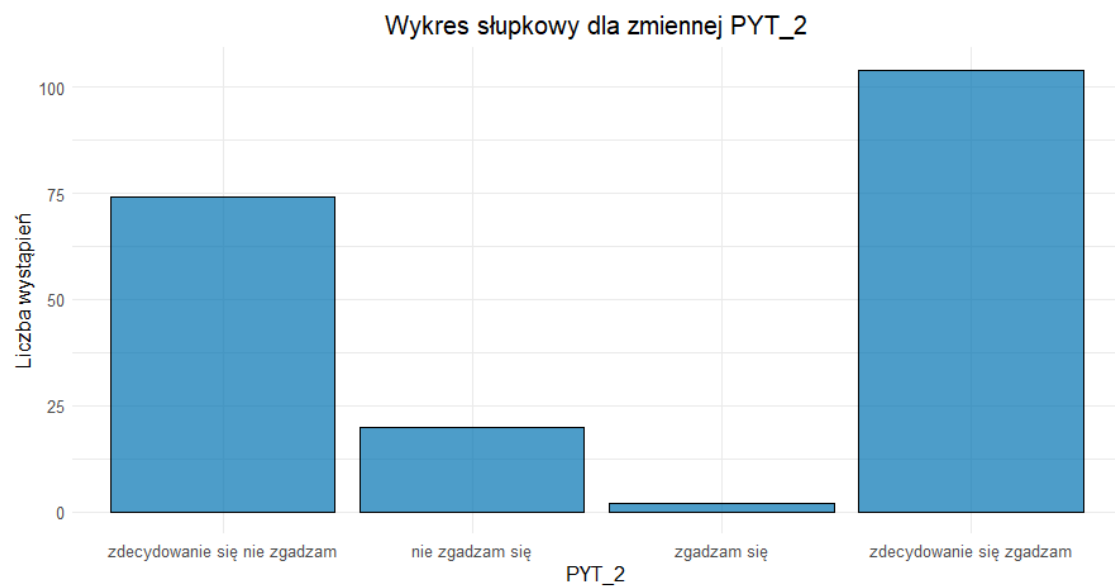


Rysunek 1: Wykres słupkowy dla zmiennej PYT_1

Rozkład odpowiedzi dla zmiennej PYT_1

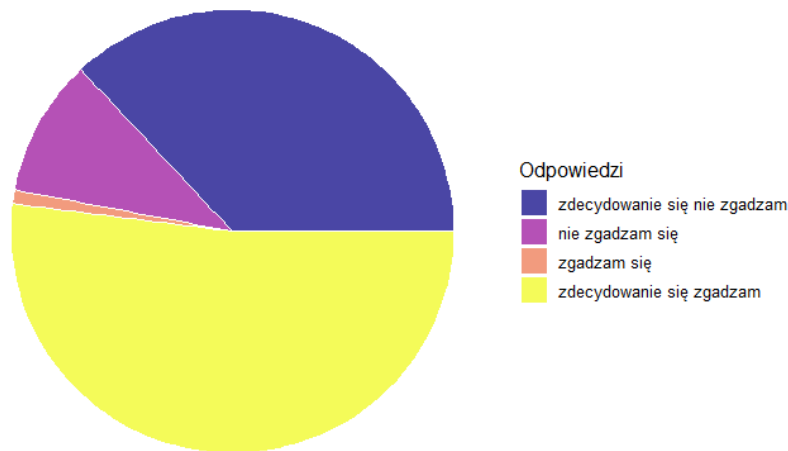


Rysunek 2: Wykres kołowy dla PYT_1



Rysunek 3: Wykres słupkowy dla PYT_2

Rozkład odpowiedzi dla zmiennej PYT_2



Rysunek 4: Wykres kołowy dla PYT_2

Na rysunkach 1 oraz 3 zostały przedstawione wykresy słupkowe odpowiednio dla zmiennych **PYT_1** i **PYT_2**. Natomiast na rysunkach 2 oraz 4 wykresy kołowe dla tych samych zmiennych. Jak możemy zauważyć na wykresie 1 i 2, większość ankietowanych w odpowiedzi na pytanie 1 zaznaczyła opcję "zgadzam się". Natomiast najmniej osób zaznaczyło negatywne odpowiedzi. W przypadku pytania 2 (wykresy 3, 4) przeważają skrajne opcje. Najwięcej osób zaznaczyło odpowiedź "zdecydowanie się zgadzam", jednak prawie równie popularna jest odpowiedź "zdecydowanie się nie zgadzam".

3.6 Zadanie 1.5

Sporządź tablice wielodzzielcze dla par zmiennych: **PYT_1** i **DZIAŁ**, **PYT_1** i **STAŻ**, **PYT_1** i **CZY_KIER**, **PYT_1** i **PŁEĆ** oraz **PYT_1** i **WIEK_KAT**.

```
1 generate_multiplication_table <- function(df, variable_1, variable_2) {
2   # Tworzenie tablicy wielodzzielczej
3   result <- table(df[[variable_1]], df[[variable_2]])
4
5   # Wyświetlenie tablic z odpowiednimi opcjami stylizacji
6   kable(result,
7     caption = paste("Tablica wielodzzielcza dla zmiennych", variable_
8       1, "i", variable_2),
9     align = "c",
10    booktabs = TRUE) %>%
11    kable_styling(latex_options = c("striped", "hold_position"))
12  }
13 # Wyświetlanie tablic
14 generate_multiplication_table(df, "PYT_1", "DZIAŁ")
15 generate_multiplication_table(df, "PYT_1", "STAŻ")
16 generate_multiplication_table(df, "PYT_1", "CZY_KIER")
17 generate_multiplication_table(df, "PYT_1", "PŁEĆ")
18 generate_multiplication_table(df, "PYT_1", "WIEK_KAT")
```

W wyniku działania powyższego kodu otrzymaliśmy 5 tablic wielodzielczych, które są widoczne poniżej.

	DK	DS	HR	IT
zdecydowanie się nie zgadzam	9	3	2	0
nie zgadzam się	10	3	2	2
nie mam zdania	17	14	5	4
zgadzam się	51	15	19	15
zdecydowanie się zgadzam	11	10	3	5

Tabela 8: Tablica wielodzielcza dla zmiennych PYT_1 i DZIAŁ

	1	2	3
zdecydowanie się nie zgadzam	5	5	4
nie zgadzam się	6	10	1
nie mam zdania	8	26	6
zgadzam się	19	75	6
zdecydowanie się zgadzam	3	24	2

Tabela 9: Tablica wielodzielcza dla zmiennych PYT_1 i STAŻ

	Nie	Tak
zdecydowanie się nie zgadzam	10	4
nie zgadzam się	14	3
nie mam zdania	34	6
zgadzam się	88	12
zdecydowanie się zgadzam	27	2

Tabela 10: Tablica wielodzielcza dla zmiennych PYT_1 i CZY_KIER

	K	M
zdecydowanie się nie zgadzam	3	11
nie zgadzam się	7	10
nie mam zdania	14	26
zgadzam się	36	64
zdecydowanie się zgadzam	11	18

Tabela 11: Tablica wielodzielcza dla zmiennych PYT_1 i PŁEĆ

	0-35	36-45	46-55	56+
zdecydowanie się nie zgadzam	1	11	2	0
nie zgadzam się	6	7	1	3
nie mam zdania	3	24	5	8
zgadzam się	13	50	25	12
zdecydowanie się zgadzam	3	12	12	2

Tabela 12: Tablica wielodzielcza dla zmiennych PYT_1 i WIEK_KAT

W tabelach 8 - 12 widoczne są tablice wielodzielcze odpowiednio dla par zmiennych:

- PYT_1 i DZIAŁ - tabela 8,
- PYT_1 i STAŻ - tabela 9,
- PYT_1 i CZY_KIER - tabela 10,
- PYT_1 i PŁEĆ - tabela 11,
- PYT_1 i WIEK_KAT - tabela 12.

Możemy z nich wyczytać liczbę wystąpień dla poszczególnych par z danej kategorii.

3.7 Zadanie 1.6

Sporządź tablicę wielodzielczą dla pary zmiennych: PYT_2 i PYT_3.

```

1 # zwykła tablica wielodzielcza
2 generate_multiplication_table(df, "PYT_2", "PYT_3")
3
4 # tablica wielodzielcza z sumami
5 generate_multiplication_table_2 <- function(df, variable_1, variable_2) {
6
7   # Tworzenie tablicy licznosci
8   result <- table(df[[variable_1]], df[[variable_2]])
9
10  # Policzenie sumy dla każdego wiersza
11  row_sum <- rowSums(result)
12
13  # Policzenie sumy dla każdej kolumny
14  col_sum <- colSums(result)
15
16  # Dodanie sum do tabeli jako dodatkowych wierszy i kolumn
17  sum_table <- cbind(result, "Suma" = row_sum)
18  sum_table <- rbind(sum_table, "Suma" = c(col_sum, ""))
19
20  # Zwrócenie tabeli
21  kable(sum_table,
22        caption = paste("Tabela wielodzielcza z sumami dla zmiennych",
23                          variable_1, "i", variable_2),
24        align = "c",
25        booktabs = TRUE) %>%
26    kable_styling(latex_options = c("striped", "hold_position"))
27
28 generate_multiplication_table_2(df, "PYT_2", "PYT_3")

```

W wyniku działania powyższego kodu otrzymaliśmy tablicę wielodzielczą z sumami dla zmiennych **PYT_2** i **PYT_3**, która jest widoczna poniżej.

	zdecydowanie się nie zgadzam	nie zgadzam się	zgadzam się	zdecydowanie się zgadzam	Suma
zdecydowanie się nie zgadzam	49	16	5	4	74
nie zgadzam się	3	6	10	1	20
zgadzam się	0	0	2	0	2
zdecydowanie się zgadzam	0	8	15	81	104
Suma	52	30	32	86	

Tabela 13: Tablica wielodzielcza z sumami dla zmiennych **PYT_2** i **PYT_3**

W tabeli 13 widzimy jak wygląda liczba wystąpień dla poszczególnych par odpowiedzi na **PYT_2** i **PYT_3**.

3.8 Zadanie 1.7

Utwórz zmienną **CZY_ZADOW** na podstawie zmiennej **PYT_2** łącząc kategorie "nie zgadzam się" i "zdecydowanie się nie zgadzam" oraz "zgadzam się" i "zdecydowanie się zgadzam".

Wybraliśmy nowszą wersję zadania, czyli tą, w której tworzymy na podstawie zmiennej **PYT_2**. Poniżej widoczny jest kod, który wykorzystaliśmy do tego celu.

```
1 # Tworzenie zmiennej CZY_ZADOW na podstawie zmiennej PYT_2
2 df$CZY_ZADOW <- ifelse(df$PYT_2 %in% c("zdecydowanie się nie zgadzam", "nie zgadzam się"), "NIE", "TAK")
```

Powyższy kod tworzy nową zmienną **CZY_ZADOW**, którą wykorzystamy w późniejszych zadaniach.

3.9 Zadanie 1.8

Korzystając z funkcji *mosaic* z biblioteki *vcd*, sporządź wykresy mozaikowe odpowiadające parom zmiennych: **CZY_ZADOW** i **DZIAŁ**, **CZY_ZADOW** i **STAŻ**, **CZY_ZADOW** i **CZY_KIER**, **CZY_ZADOW** i **PŁEĆ** oraz **CZY_ZADOW** i **WIEK_KAT**. Czy na podstawie uzyskanych wykresów można postawić pewne hipotezy dotyczące realizacji między powyższymi zmiennymi? Spróbuj sformułować kilka takich hipotez.

W tym zadaniu wykorzystaliśmy funkcję *mosaic* z biblioteki *vcd*. Funkcja ta służy do tworzenia wykresów mozaikowych, które mają zastosowanie w wizualizacji zależności pomiędzy dwoma zmiennymi kategorycznymi.

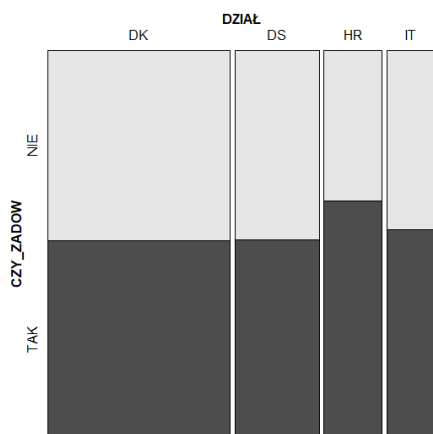
Na podstawie danych wejściowych funkcja *mosaic* tworzy prostokątne obszary (mozaikę), które są proporcjonalne do liczby obserwacji w każdej kombinacji kategorii.

Tego typu wykresy pozwalają nam łatwo zauważyć, które kombinacje kategorii są bardziej lub mniej powszechne, co może pomóc nam dostarczyć informacji na temat zależności między poszczególnymi zmiennymi.

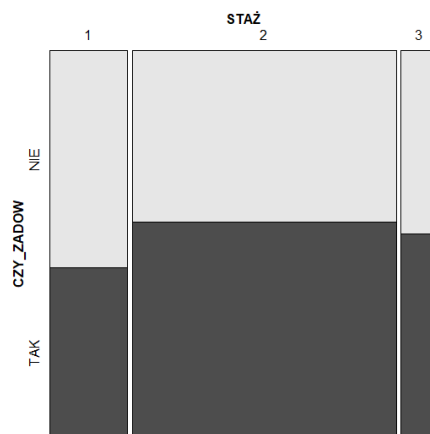
```
1 mosaic( ~ CZY_ZADOW + DZIAŁ, data = df, highlight = TRUE)
2 mosaic( ~ CZY_ZADOW + STAŻ, data = df, highlight = TRUE)
3 mosaic( ~ CZY_ZADOW + CZY_KIER, data = df, highlight = TRUE)
4 mosaic( ~ CZY_ZADOW + PŁEĆ, data = df, highlight = TRUE)
5 mosaic( ~ CZY_ZADOW + WIEK_KAT, data = df, highlight = TRUE)
```

Korzystając z powyższego kodu sporządziliśmy 5 wykresów mozaikowych dla odpowiednich par zmiennych:

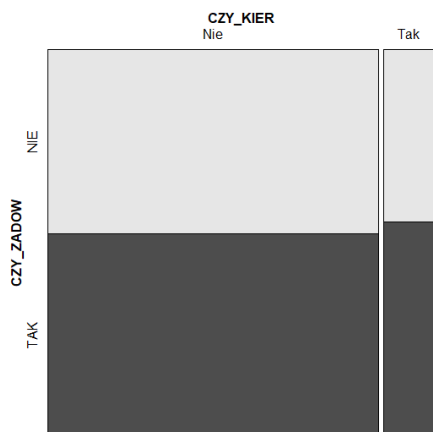
- **CZY_ZADOW** i **DZIAŁ** - rysunek 2,
- **CZY_ZADOW** i **STAŻ** - rysunek 3,
- **CZY_ZADOW** i **CZY_KIER** - rysunek 4,
- **CZY_ZADOW** i **PŁEĆ** - rysunek 5,
- **CZY_ZADOW** i **WIEK_KAT** - rysunek 6.



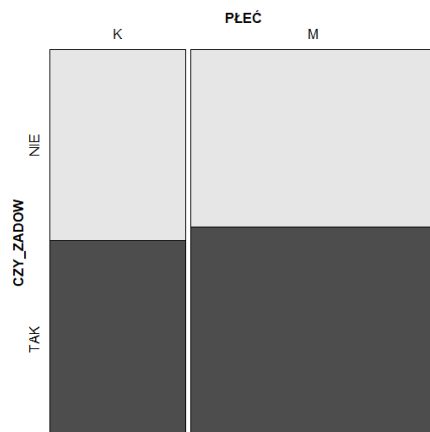
Rysunek 5: CZY_ZADOW i DZIAŁ



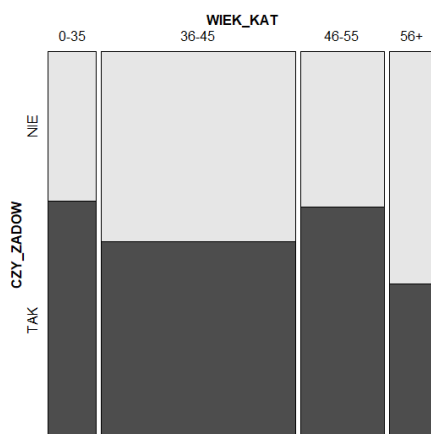
Rysunek 6: CZY_ZADOW i STAŻ



Rysunek 7: CZY_ZADOW i CZY_KIER



Rysunek 8: CZY_ZADOW i PŁEĆ



Rysunek 9: CZY_ZADOW i WIEK_KAT

Analizując powyższe wykresy możemy wyciągnąć kilka hipotez, między innymi:

1. Istnieje zależność pomiędzy poziomem zadowolenia pracowników, a działem w jakim oni pracują.
2. Zazwyczaj pracownicy z dłuższym stażem mają tendencję do bycia bardziej zadowolonymi niż ci ze stażem krótszym.
3. Pracownicy zajmujący kierownicze stanowiska są bardziej zadowoleni niż ci na niższych stanowiskach.
4. Płeć pracowników może mieć wpływ na ich zadowolenie z pracy.
5. Istnieje dość znaczący związek pomiędzy wiekiem pracowników, a ich poziomem zadowolenia z pracy.

Możemy także zauważyć, że najwięcej osób pracuje w dziale kreatywnym (DK), prze-
ważają osoby w wieku 36-45, większość osób nie pracuje na stanowisku kierowniczym.
Dodatkowo widzimy, że w firmie pracuje więcej mężczyzn.

4 Część II

4.1 Zadanie 2

Zapoznaj się z biblioteką *likert* i dostępnymi tam funkcjami *summary* oraz *plot* (wykresy
typu "bar", "heat" oraz "density"), a następnie zilustruj odpowiedzi na pytanie "Jak bar-
dzo zgadzasz się ze stwierdzeniem, że firma pozwala na (...)?" (zmienna **PYT_1**) w całej
badanej grupie oraz w podgrupach ze względu na zmienną **CZY_KIER**.

Do tego zadania wykorzystaliśmy funkcję *summary* oraz *plot* z biblioteki *likert*.

Funkcja *summary* generuje podsumowanie danych związane ze skalami Likerta, w formie
tabeli zawierającej następujące kolumny:

- Item - pytanie, dla którego są prezentowane wyniki,
- low - ankietowani, którzy odpowiedzieli na zadane pytanie na niższym końcu skali:
"zdecydowanie się nie zgadzam" lub "nie zgadzam się",
- neutral - ankietowani, którzy odpowiedzieli neutralnie na zadane pytanie "nie mam
zdania",
- high - ankietowani, którzy odpowiedzieli na zadane pytanie na wyższym końcu skali:
"zgadzam się" lub "zdecydowanie się zgadzam",
- mean - wartość średnia,
- sd - odchylenie standardowe z odpowiedzi.

Natomiast funkcja *plot* służy do wizualizacji tych danych. Za jej pomocą możemy stwo-
rzyć między innymi:

- wykres słupkowy (bar plot),
- wykres rozkładu gęstości (density plot),
- wykres mapy ciepła/macierzy kolorów (heatmap).

W pierwszej części tego zadania zilustrowaliśmy odpowiedź na **PYT_1** w całej badanej
grupie za pomocą wykresów typu "bar", "heat" oraz "density". W tym celu wykorzystaliśmy
kod widoczny poniżej.

```
1 df$CZY_KIER <- factor(df$CZY_KIER, levels = c("Nie", "Tak"))
2 head(df)
3
4 # ----- bez grupowania -----
5 # Przekształcenie podsumowania na ramkę danych
6 likert_pyt_1 <- likert(df[, "PYT_1", drop=FALSE])
7
8 summ_pyt_1 <- summary(likert_pyt_1)
9
```



```

10 # Wyświetlenie podsumowania w formie tabeli
11 kable(summ_pyt_1,
12       caption = "Podsumowanie odpowiedzi na pytanie PYT_1",
13       align = "c",
14       booktabs = TRUE) %>%
15 kable_styling(latex_options = c("striped", "hold_position"))

```

W wyniku działania funkcji *summary* otrzymaliśmy poniższą tabelę.

Item	low	neutral	high	mean	sd
PYT_1	15.5	20	64.5	3.565	1.063688

Tabela 14: Podsumowanie odpowiedzi na pytanie PYT_1 bez grupowania

Jak możemy zauważyć w tabeli 14, większość osób (64.5%) uważa, że firma pozwala na elastyczne godziny pracy, 20% osób nie ma na ten temat zdania, natomiast pozostałe 15.5% uważa, że firma nie umożliwia zachowania równowagi pomiędzy pracą, a życiem prywatnym. Warto również zaznaczyć, że wyniki zostały opracowane na podstawie 5-stopniowej skali, gdzie:

- 1 - zdecydowanie się nie zgadzam,
- 2 - nie zgadzam się,
- 3 - nie mam zdania,
- 4 - zgadzam się,
- 5 - zdecydowanie się zgadzam.

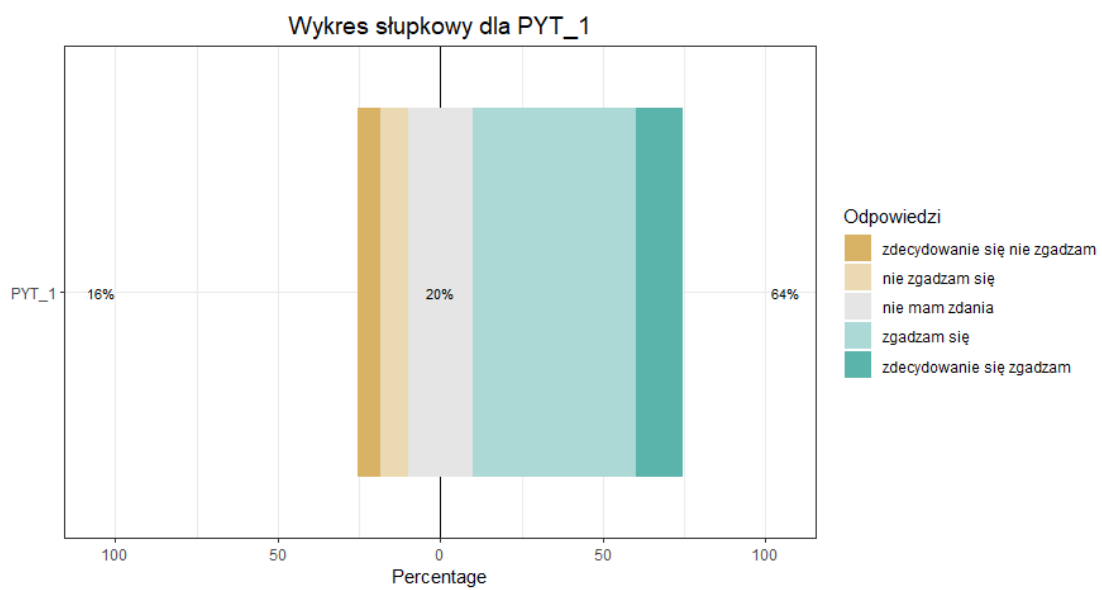
Widzimy, że średnia wynosi 3.565 co pozwala nam stwierdzić, że większość osób uważa, że firma pozwala na elastyczne godziny pracy. Odchylenie standardowe jest równe 1.063688 co pozwala stwierdzić, że niewiele osób zaznaczyło skrajne odpowiedzi.

Następnie, z pomocą funkcji *plot* sporządziliśmy wykres słupkowy, mapy ciepła oraz gęstości.

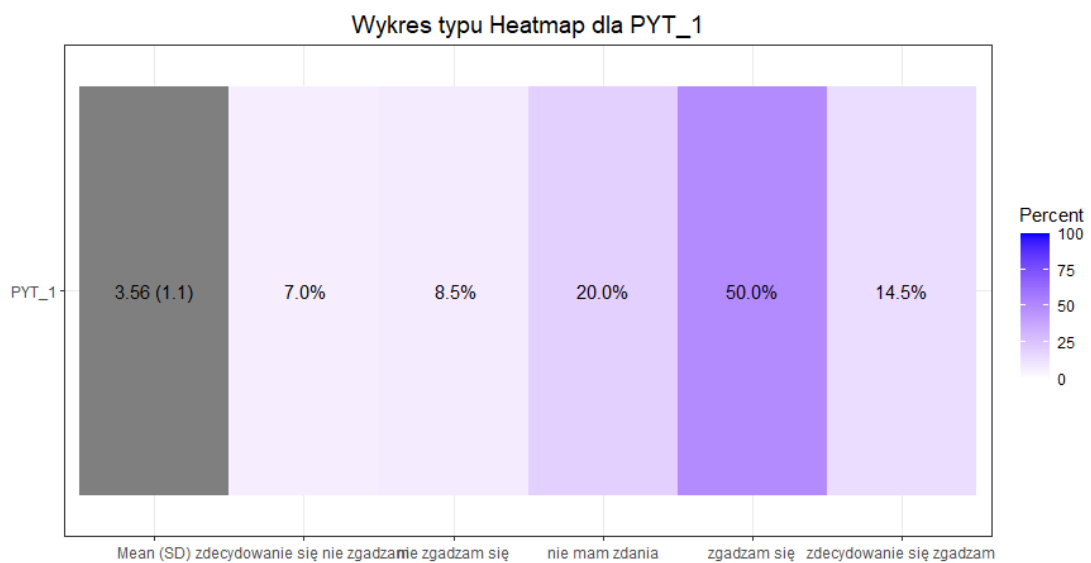
```

1 plot(likert_pyt_1, type = "bar", auto.key = list(columns = 2)) +
2   labs(title = "Wykres słupkowy dla PYT_1") +
3   guides(fill = guide_legend(title = "Odpowiedzi")) +
4   theme_bw() +
5   theme(plot.title = element_text(size = 14,, hjust = 0.5))
6
7 plot(likert_pyt_1, type = "heat") +
8   labs(title = "Wykres typu Heatmap dla PYT_1") +
9   theme_bw() +
10  theme(plot.title = element_text(size = 14, hjust = 0.5))
11
12 plot(likert_pyt_1, type = "density") +
13   theme_minimal() +
14   labs(title = "Wykres typu Density") +
15   theme(plot.title = element_text(size = 14, hjust = 0.5))

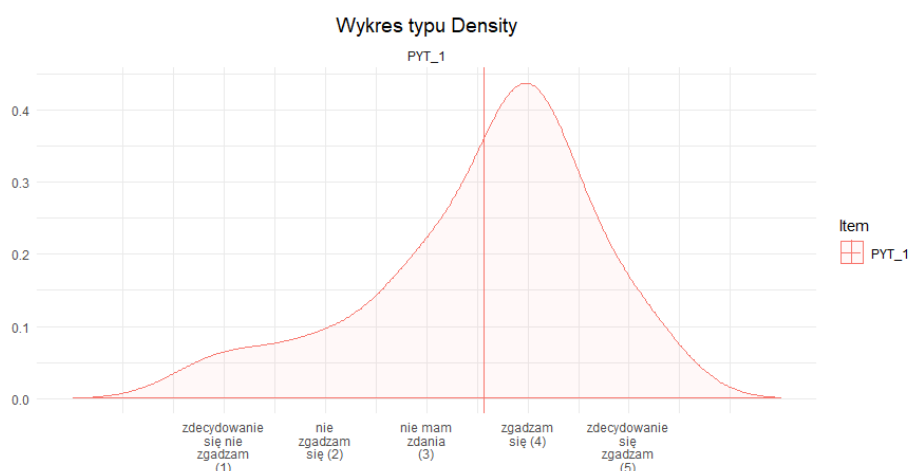
```



Rysunek 10: Wykres słupkowy dla zmiennej PYT_1



Rysunek 11: Wykres typu Heatmap dla zmiennej PYT_1



Rysunek 12: Wykres typu Density dla zmiennej PYT_1

Na rysunkach 10 - 12 możemy zauważyć, że nasza teza mówiąca o tym, że większość osób pracujących w firmie uważa, że oferuje ona elastyczny grafik pracy, jest prawdziwa. Widzimy także, że największy odsetek osób zaznaczyło opcję "zgadzam się" na wykresie 10 pole to jest najszersze, na wykresie 11 kolor jest najciemniejszy, natomiast na wykresie 12 w tym miejscu występuje pik funkcji gęstości. Na wszystkich powyższych wykresach widzimy także, że najmniej osób zaznaczyło odpowiedź "zdecydowanie się nie zgadzam".

Następnie pogrupowałyśmy ankietowanych ze względu na fakt zajmowania przez nich kierowniczych stanowisk (zmienna **CZY_KIER**) i zwizualizowałyśmy wyniki na wykresie typu "bar".

```
1 # ----- z grupowaniem -----
2 likert_py1_grouped <- likert(df[, "PYT_1", drop = FALSE], grouping = df
  $CZY_KIER)
3
4 summ_py1_grouped <- summary(likert_py1_grouped)
5 print(summ_py1_grouped)
6
7 # Wyświetlenie podsumowania w formie tabeli
8 kable(summ_py1_grouped,
9       caption = "Podsumowanie odpowiedzi na pytanie PYT_1",
10      align = "c",
11      booktabs = TRUE) %>%
12 kable_styling(latex_options = c("striped", "hold_position"))
```

Group	Item	low	neutral	high	mean	sd
Nie	PYT_1	13.87283	19.65318	66.47399	3.624277	1.030291
Tak	PYT_1	25.92593	22.22222	51.85185	3.185185	1.210119

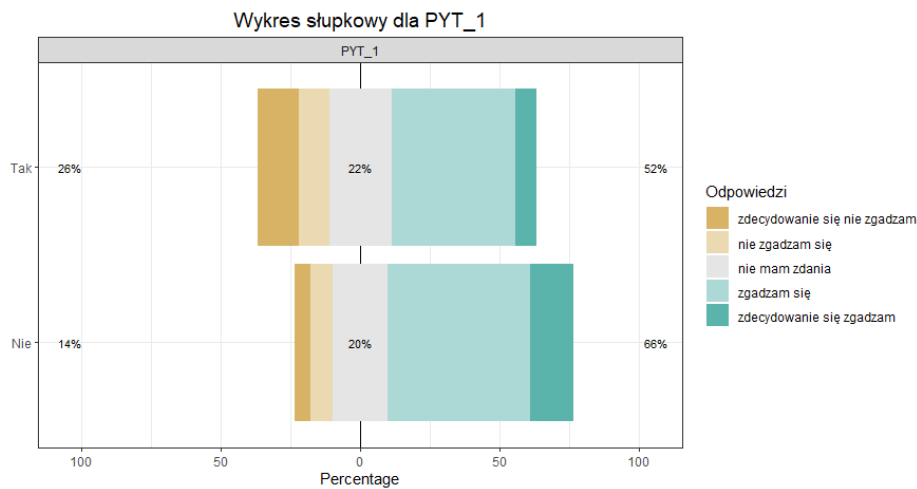
Tabela 15: Podsumowanie odpowiedzi na pytanie PYT_1 z grupowaniem

Jak możemy zauważyć w tabeli 15, także sporządzonej z pomocą funkcji *summary*, osoby pracujące na stanowisku kierowniczym w większym stopniu nie zgadzają się ze stwierdzeniem, że firma umożliwia zachowanie równowagi pomiędzy pracą i życiem prywatnym

(jednak ciągle są zadowoleni z elastyczności grafiku). W tym przypadku, podobnie jak w zadaniu 2 została użyta 5-stopniowa skala (4.1). Analizując zachowanie średniej potwierdza się nasza teza mówiąca o tym, że osoby nie pracujące na stanowisku kierowniczym w większym stopniu zgadzają się ze stwierdzeniem, że firma oferuje elastyczny grafik pracy. Patrząc na odchylenie standardowe widzimy, że wśród osób zajmujących kierownicze stanowiska odpowiedzi są bardziej zróżnicowane (jednak z przewagą odpowiedzi pozytywnych).

W tym przypadku także podjęliśmy próbę zwizualizowania odpowiedzi na wykresach.

```
1 plot(likert_pyt_1_grouped, type = "bar", auto.key = list(columns = 2)) +
2   labs(title = "Wykres słupkowy dla PYT_1") +
3   guides(fill = guide_legend(title = "Odpowiedzi")) +
4   theme_bw() +
5   theme(plot.title = element_text(size = 14, hjust = 0.5))
```



Rysunek 13: Wykres słupkowy dla zmiennej PYT_2

W tym przypadku udało nam się sporządzić wykres słupkowy (bar plot), który widoczny jest na rysunku 13. Analizując jego wygląd możemy zauważyć, że nasza teza potwierdziła się. W obu podgrupach zdecydowanie przeważa odpowiedź "zgadzam się" (najszerzy przedział). W zależności od zajmowanego stanowiska widoczna jest dość znacząca różnica w skrajnych odpowiedziach: "zdecydowanie się zgadzam" i "zdecydowanie się nie zgadzam". Widzimy także, że osoby zaznaczające odpowiedź "nie mam zdania" stanowią w obu podgrupach podobny odsetek. Podsumowując, możemy stwierdzić, że osoby nie zajmujące kierowniczych stanowisk w większej części zgadzają się ze stwierdzeniem, że firma zapewnia elastyczny grafik pracy. Jednak wszyscy, bez względu na zajmowane stanowisko, są raczej zadowoleni z elastyczności grafiku.

4.2 Zadanie 3

Zapoznaj się z funkcją *sample* z biblioteki *stats*, a następnie wylosuj próbkę o liczności 10% wszystkich rekordów z pliku "ankieta.csv" w dwóch wersjach: ze zwracaniem oraz bez zwracania.

W tym zadaniu wykorzystaliśmy funkcję *sample* z biblioteki *stats*. Służy ona do losowego próbkowania elementów z zestawu danych. Pozwala na wybieranie losowych próbek danych bez zwracania i ze zwracaniem.

```

1 # Liczność próbki (10% wszystkich rekordów)
2 sample_size <- round(0.1 * nrow(df))
3
4 # losowanie ze zwracaniem
5 with_returning <- sample(1:nrow(df), size = sample_size, replace = TRUE)
6
7 # losowanie bez zwracania
8 without_returning <- sample(1:nrow(df), size = sample_size, replace =
  FALSE)
9
10 # Wyświetlenie wylosowanych próbek
11 kable(with_returning,
12       caption = "Wylosowana próbka ze zwracaniem",
13       align = "c",
14       col.names = c("Próbka"),
15       booktabs = TRUE) %>%
16   kable_styling(latex_options = c("striped", "hold_position"))
17
18 kable(without_returning,
19       caption = "Wylosowana próbka bez zwracania",
20       align = "c",
21       col.names = c("Próbka"),
22       booktabs = TRUE) %>%
23   kable_styling(latex_options = c("striped", "hold_position"))

```

W poniższej tabeli przedstawiliśmy wyniki działania powyższego kodu. W obu przypadkach wylosowałyśmy 20 rekordów, co stanowi 10% wszystkich obserwacji. Warto zaznaczyć, że są to jedynie przykładowe wyniki. Ze względu na fakt, że losowanie odbywa się w sposób losowy wyniki w przypadku każdorazowego wywołania kodu będą inne.

Losowanie ze zwracaniem	Losowanie bez zwracania
80	65
43	200
157	53
128	167
29	1
154	36
88	58
82	125
159	130
155	96
24	180
66	162
77	144
108	50
45	145
18	83
155	6
73	179
52	177
97	40

Tabela 16: Przykładowe wylosowane próbki ze zwracaniem i bez zwracania

W tabeli 16 widzimy przykładowe wyniki losowania ze zwracaniem i bez.

4.3 Zadanie 4

Zaproponuj metodę symulowania zmiennych losowych z rozkładu dwumianowego. Napisz funkcję do generowania realizacji, a następnie zaprezentuj jej działanie porównując wybrane teoretyczne i empiryczne charakterystyki dla przykładowych wartości parametrów rozkładu: n i p .

W celu wygenerowania zmiennych losowych z rozkładu dwumianowego wykorzystaliśmy poniższy algorytm:

1. Inicjalizacja wyników - tworzymy pusty wektor do przechowywania późniejszych wyników.
2. Dla każdej próbki:
 - (a) Inicjalizujemy zmienną służącą do przechowywania liczby sukcesów w danej próbce.
 - (b) Dla każdej z prób:
 - i. Generujemy losową liczbę z rozkładu jednostajnego $U(0, 1)$.
 - ii. Sprawdzamy czy liczba ta jest mniejsza niż zadane prawdopodobieństwo sukcesu p . Jeśli tak, to zwiększamy licznik sukcesów.
 - (c) Zapisujemy liczbę sukcesów do wektora wynikowego.
3. Zwracamy wektor wynikowy.

Poniżej widoczny jest kod obrazujący działanie powyższego algorytmu.

```
1 generate_binomial <- function(n, p, size) {  
2   result <- numeric(size)  
3   for (i in 1:size) {  
4     successes <- 0  
5     for (j in 1:n) {  
6       if (runif(1) < p) {  
7         successes <- successes + 1  
8       }  
9     }  
10    result[i] <- successes  
11  }  
12  return(result)  
13 }
```

W celu sprawdzenia poprawności działania zaproponowanej metody porównaliśmy ze sobą teoretyczne i empiryczne wartości średniej i wariancji. Narysowaliśmy także teoretyczne i empiryczne rozkłady prawdopodobieństw oraz dystrybuant. Zarówno liczba próbek ("size"), jak i liczba prób ("n") wynosiła 1000. Wartości prawdopodobieństw sukcesu p ("p-values") należą do przedziału 0.1, 0.5, 0.9.

```
1 n <- 1000  
2 p_values <- c(0.1, 0.5, 0.9)  
3 size <- 1000  
4  
5 empirical_p <- numeric(length(p_values))
```

```

6
7 theoretical_mean <- numeric(length(p_values))
8 theoretical_variance <- numeric(length(p_values))
9 empirical_mean <- numeric(length(p_values))
10 empirical_variance <- numeric(length(p_values))
11
12 for (i in 1:length(p_values)) {
13   p <- p_values[i]
14
15   empirical_binomial <- generate_binomial(n, p, size)
16   empirical_p[i] <- mean(empirical_binomial) / n
17
18   theoretical_mean[i] <- n * p
19   theoretical_variance[i] <- n * p * (1 - p)
20
21   empirical_mean[i] <- mean(empirical_binomial)
22   empirical_variance[i] <- var(empirical_binomial)
23
24 }
25
26 results_binomial <- data.frame(
27   "Charakterystyka" = c("Teoretyczna średnia", "Empiryczna średnia", "
   Teoretyczna wariancja", "Empiryczna wariancja"),
28   "p = 0.1" = c(theoretical_mean[1], empirical_mean[1], theoretical_
   variance[1], empirical_variance[1]),
29   "p = 0.5" = c(theoretical_mean[2], empirical_mean[2], theoretical_
   variance[2], empirical_variance[2]),
30   "p = 0.9" = c(theoretical_mean[3], empirical_mean[3], theoretical_
   variance[3], empirical_variance[3])
31 )
32
33
34 kable(results_binomial,
35       caption = "Porównanie empirycznych statystyk z teoretycznymi",
36       align = "c",
37       col.names = c("Charakterystyka", "p = 0.1", "p = 0.5", "p = 0.9"),
38       booktabs = TRUE) %>%
39   kable_styling(latex_options = c("striped", "hold_position"))

```

Charakterystyka	$p = 0.1$	$p = 0.5$	$p = 0.9$
Teoretyczna średnia	100.0000	500.0000	900.0000
Empiryczna średnia	99.5460	499.3590	899.5260
Teoretyczna wariancja	90.0000	250.0000	90.0000
Empiryczna wariancja	88.4143	259.1272	92.85418

Tabela 17: Porównanie teoretycznych i empirycznych charakterystyk

W tabeli numer 17 zaprezentowano przykładowe wyniki wysymulowania zmiennych z rozkładu dwumianowego $B(1000, 0.1)$, $B(1000, 0.5)$, $B(1000, 0.9)$. Dla tych rozkładów porównaliśmy teoretyczne i empiryczne średnie oraz wariancje. Jak możemy zauważyć, są one do siebie zbliżone co pozwala nam wnioskować, że zaproponowana przez nas metoda jest najprawdopodobniej poprawna. W celu zweryfikowania tej tezy sporządziliśmy wykresy teoretycznych i empirycznych rozkładów prawdopodobieństwa oraz dystrybuant dla każdego p . Poniżej widoczny jest kod oraz wyniki.

```

1 results_comparison <- data.frame(
2   "Teoretyczne" = p_values,
3   "Empiryczne" = empirical_p
4 )
5 ggplot(results_comparison, aes(x = factor(0:(length(Teoretyczne) - 1))))
6   +
7   geom_bar(aes(y = Empiryczne, fill = "Empiryczna"), stat = "identity",
8     width = 0.05) +
9   geom_bar(aes(y = Teoretyczne, fill = "Teoretyczna"), stat = "identity",
10    width = 0.02, alpha = 0.5) +
11   geom_point(aes(y = Empiryczne, color = "Empiryczna"), size = 6) +
12   geom_point(aes(y = Teoretyczne, color = "Teoretyczna"), size = 4, alpha
13     = 0.5) +
14   scale_fill_manual(values = c("Teoretyczna" = "red", "Empiryczna" = "
15     skyblue"), name = "") +
16   scale_color_manual(values = c("Teoretyczna" = "red", "Empiryczna" = "
17     skyblue"), name = "") +
18   labs(x = "X", y = "Prawdopodobieństwo") +
19   ggtitle("Porównanie teoretycznych i empirycznych prawdopodobieństw\
20     rozkładu dwumianowego") +
21   theme_minimal() +
22   theme(
23     axis.title = element_text(size = 12),
24     axis.text = element_text(size = 10),
25     plot.title = element_text(size = 15, hjust = 0.5))

```



Rysunek 14: Porównanie teoretycznych i empirycznych prawdopodobieństw

```

1 sorted_empirical <- sort(empirical_p)
2
3 ggplot() +
4   geom_step(data = data.frame(x = c(0, 0, 1, 2, 3), y = c(0, sorted_
5     empirical, max(sorted_empirical))), aes(x, y),
6     color = "skyblue", linetype = "solid", size = 1.5) +
7   geom_step(data = data.frame(x = c(0, 0, 1, 2, 3), y = c(0, p_values,
8     max(sorted_empirical))), aes(x, y),
9     color = "red", linetype = "dashed", size = 1) +
10   labs(x = "X", y = "Dystrybuanta",
11     title = "Porównanie teoretycznej i empirycznej dystrybuanty dla p"
12   ) +

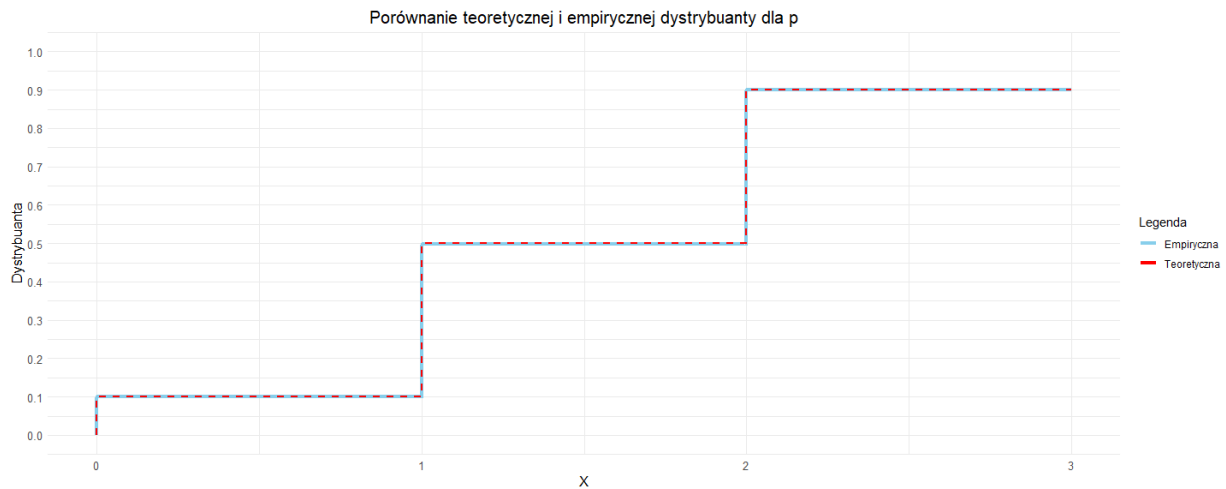
```



```

10 scale_x_continuous(limits = c(0, 3), breaks = c(0, 1, 2, 3)) +
11 scale_y_continuous(limits = c(0, 1), breaks = seq(0, 1, by = 0.1)) +
12 theme_minimal() +
13 theme(
14   axis.title = element_text(size = 12),
15   axis.text = element_text(size = 10),
16   plot.title = element_text(size = 15, hjust = 0.5))

```



Rysunek 15: Porównanie dystrybuanty teoretycznej i empirycznej

Jak możemy wnioskować z wykresu 14 (porównanie teoretycznych i empirycznych prawdopodobieństw) oraz z wykresu 15 (porównanie teoretycznej i empirycznej dystrybuanty) zaproponowana przez nas metoda symulacji najprawdopodobniej jest poprawna. Widzimy, że teoretyczne i empiryczne wartości się dość dobrze pokrywają dla każdej z badanych wartości prawdopodobieństwa sukcesu p .

4.4 Zadanie 5

Zaproponuj metodę symulowania wektorów losowych z rozkładu wielomianowego. Napisz funkcję do generowania realizacji, a następnie zaprezentuj jej działanie porównując wybrane teoretyczne i empiryczne charakterystyki dla przykładowych wartości parametrów rozkładu: n i p .

W celu wygenerowania wektorów losowych z rozkładu wielowymiarowego wykorzystaliśmy poniższy algorytm:

1. Inicjalizacja wyników - tworzymy pusty wektor do przechowywania późniejszych wyników.
2. Inicjalizacja pomocniczych zmiennych określających pozostałą liczbę prób oraz sumującą prawdopodobieństwa p .
3. Dla każdego p (z wyjątkiem ostatniego):
 - (a) Jeśli suma prawdopodobieństw jest różna od 0 to generujemy losową zmienną z rozkładu dwumianowego.

- (b) Aktualizujemy liczbę prób pozostałych do wykonania oraz sumę prawdopodobieństw.
4. Do ostatniej próbki przypisujemy pozostałe próby.
 5. Zwracamy wektor wynikowy zawierający liczbę sukcesów dla każdej próbki.

Poniżej widoczny jest kod obrazujący działanie powyższego algorytmu.

```
1 generate_multinomial <- function(n, p) {
2   result <- numeric(length(p))
3   remaining_trials <- n
4   prob_sum <- 1
5
6   for (i in 1:(length(p) - 1)) {
7     if (prob_sum != 0) {
8       result[i] <- rbinom(1, remaining_trials, p[i]/prob_sum)
9     } else {
10      result[i] <- 0
11    }
12    remaining_trials <- remaining_trials - result[i]
13    prob_sum <- prob_sum - p[i]
14  }
15
16  result[length(p)] <- remaining_trials
17  return(result)
18 }
```

W celu sprawdzenia poprawności działania zaproponowanej metody porównaliśmy ze sobą teoretyczne i empiryczne wartości prawdopodobieństwa sukcesów p . Liczba prób wynosiła 10000. Wartości prawdopodobieństw sukcesu p wynoszą odpowiednio: (0.1, 0.3, 0.4, 0.2), (0.5, 0.3, 0.2), (0.1, 0.1, 0.3, 0.2, 0.3).

```
1 p_values <- list(
2   c(0.1, 0.3, 0.4, 0.2),
3   c(0.5, 0.3, 0.2),
4   c(0.1, 0.1, 0.3, 0.2, 0.3)
5 )
6
7 n <- 10000
8
9 results_df <- data.frame()
10
11
12 for (p in p_values) {
13   sample <- generate_multinomial(n, p)
14
15   empirical_p <- sample / n
16
17   result <- data.frame(
18     Theoretical_p = as.character(p),
19     Empirical_p = as.character(empirical_p)
20   )
21
22   results_df <- rbind(results_df, result)
23 }
24
25 kable(results_df,
26       caption = "Wylosowana próbka bez zwracania",
```

```

27 align = "c",
28 col.names = c("Teoretyczne p", "Empiryczne p"),
29 booktabs = TRUE) %>%
30 kable_styling(latex_options = c("striped", "hold_position"))

```

Teoretyczne p	Empiryczne p
0.1	0.1056
0.3	0.3037
0.4	0.3914
0.2	0.1993
0.5	0.499
0.3	0.2993
0.2	0.2017
0.1	0.1002
0.1	0.1015
0.3	0.2994
0.2	0.1979
0.3	0.3010

Tabela 18: Porównanie teoretycznych i empirycznych wartości p

Jak możemy zauważyć w tabeli 18, wszystkie teoretyczne i empiryczne wartości prawdopodobieństw sukcesu p są do siebie zbliżone. Pozwala nam to wnioskować, że zaproponowana przez nas metoda symulacji najprawdopodobniej jest poprawna.

5 Część III i IV

5.1 Zadanie 6

Napisz funkcję do wyznaczania realizacji przedziału ufności Cloppera-Pearsona. Niech argumentem wejściowym będzie poziom ufności, liczba sukcesów i liczba prób lub poziom ufności i wektor danych (funkcja powinna obsługiwać oba przypadki).

W zadaniu tym zajęliśmy się wyznaczeniem realizacji przedziału ufności Cloppera-Pearsona. W tym celu wykorzystaliśmy poniższy algorytm, który jako argumenty wejściowe przyjmuje poziom ufności, liczbę sukcesów i liczbę prób:

1. Sprawdzamy czy została podana liczba prób. Jeśli nie to za dane wejściowe przyjmujemy wektor, w którym sukcesy są oznaczone jako "TAK".
 - (a) Liczymy liczbę sukcesów, a następnie liczbę prób.
2. Obliczamy przedziały ufności Cloppera-Pearsona korzystając z rozkładu beta z parametrami zależnymi od liczby sukcesów i prób.
3. Zwracamy otrzymane przedziały.

```

1 intervals_clopper_pearson <- function(alpha, successes, trials = NULL) {
2
3   if (is.null(trials)) {
4     data <- successes

```

```

5     successes <- sum(data == "TAK")
6     trials <- length(data)
7 }
8
9 result_lower <- qbeta(alpha/2, successes, trials - successes + 1)
10 result_upper <- qbeta(1 - (alpha/2), successes + 1, trials - successes)
11
12 result <- data.frame(
13     lower = result_lower,
14     upper = result_upper
15 )
16
17 return(result)
18 }

```

Zaproponowaną przez nas funkcję wykorzystamy w dalszej części raportu.

5.2 Zadanie 7

Korzystając z funkcji napisanej w zadaniu 6 (5.1). wyznacz realizacje przedziałów ufności dla prawdopodobieństwa, że pracownik jest zadowolony z wynagrodzenia w pierwszym badanym okresie oraz w drugim badanym okresie. Skorzystaj ze zmiennych **CZY_ZADW** oraz **CZY_ZADW_2** (utwórz zmienną analogicznie jak w zadaniu 3.8). Przyjmij $1 - \alpha = 0.95$.

```

1 df$CZY_ZADW_2 <- ifelse(df$PYT_3 %in% c("zdecydowanie się nie zgadzam",
2     "nie zgadzam się"), "NIE",
3     "TAK")
4
5 alph <- 0.05
6
7 clopper_pearson_PYT_3 <- intervals_clopper_pearson(alph, df$CZY_ZADW_2)
8 clopper_pearson_PYT_2 <- intervals_clopper_pearson(alph, df$CZY_ZADW)
9
10 print(clopper_pearson_PYT_3)
11 print(clopper_pearson_PYT_2)

```

Wyniki działania funkcji z zadania 5.1 dla zmiennych **PYT_2** i **PYT_3** zostały przedstawione w poniższej tabeli.

Zmienna	Dolna granica przedziału ufności	Górna granica przedziału ufności
PYT_2	0.4583305	0.6007671
PYT_3	0.5184216	0.6588694

Tabela 19: Przedziały ufności dla zmiennej PYT_2 i PYT_3

5.3 Zadanie 8

Zapoznaj się z funkcjami *rbinom* z biblioteki *stats* oraz *binom.confint* z biblioteki *binom*.

```

1 help(rbinom)
2 help(binom.confint)

```

W zadaniu 9 (5.4) będziemy używać funkcji *rbinom* z biblioteki *stats* oraz *binom.confint* z biblioteki *binom*.

Pierwsza z nich służy do generowania losowych próbek z rozkładu dwumianowego opisuującego liczbę sukcesów w sekwencji niezależnych prób Bernoulliego. Jako wynik zwraca wektor losowych liczb całkowitych, które reprezentują liczbę sukcesów w każdej próbie. Składnia funkcji *rbinom* ma następującą postać: *rbinom(n, size, prob)*, gdzie:

- *n* - rozmiar próby (np. ilość pytań w ankiecie),
- *size* - ilość powtórzeń (np. ilość ankietowanych)
- *prob* - prawdopodobieństwo sukcesu w każdej z prób.

Natomiast funkcja *binom.confint* z biblioteki *binom* służy do obliczania przedziałów ufności dla parametrów rozkładu dwumianowego. Jako wynik zwraca przedział lub przedziały ufności dla prawdopodobieństwa sukcesu w rozkładzie dwumianowym. Składnia funkcji *rbinom* ma następującą postać: *binom.confint(x, n, conf.level, methods)*, gdzie:

- *x* - liczba sukcesów w próbie,
- *n* - rozmiar próby,
- *conf.level* - poziom ufności dla przedziałów ufności,
- *methods* - metoda lub metody obliczania przedziałów ufności.

Metody obliczania przedziałów ufności dostępne w funkcji *binom.confint* to między innymi: metoda Clopper-Pearsona, metoda Wilsona czy metoda Wilsona-Score'a.

5.4 Zadanie 9

Przeprowadź symulacje, których celem jest porównanie prawdopodobieństwa pokrycia i długości przedziałów ufności Cloppera-Pearsona, Walda i trzeciego dowolnego typu zaimplementowanego w funkcji *binom.confint*. Rozważ $1 - \alpha = 0.95$, rozmiar próby $n \in \{30, 100, 1000\}$ i różne wartości prawdopodobieństwa p . Wyniki umieść na wykresach i sformułuj wnioski, które dla konkretnych danych ułatwią wybór konkretnego typu przedziału ufności.

W celu przeprowadzenia symulacji, których celem jest porównanie prawdopodobieństwa pokrycia i długości przedziałów ufności Cloppera-Pearsona, Walda i Wilsona wykorzystaliśmy poniższy algorytm:

1. Inicjalizacja zmiennych:
 - (a) Ustalamy poziom istotności α .
 - (b) Określamy liczbę prób n .
 - (c) Definiujemy wektor p zawierający prawdopodobieństwa sukcesu, dla których będziemy liczyć przedziały ufności.
 - (d) Określamy liczbę symulacji M .
2. Właściwa symulacja:

(a) Dla każdego $n \in 30, 100, 1000$:

i. Dla każdej wartości prawdopodobieństwa sukcesu p :

A. Dla każdej z M symulacji:

- Generujemy próbę z rozkładu dwumianowego.
- Obliczamy przedziały ufności Cloppera-Pearsona, Walda i Wilsona dla liczby sukcesów w próbce.
- Obliczamy długość każdego przedziału ufności, jako różnicę między górną i dolną granicą przedziału.
- Sprawdzamy, czy prawdopodobieństwo sukcesu znajduje się w wyznaczonym przedziale ufności.

B. Obliczamy średnią długość przedziałów ufności dla każdej metody.

C. Obliczamy średnie prawdopodobieństwo pokrycia dla każdej metody.

3. Zwracamy średnie długości przedziałów ufności i średnie prawdopodobieństwo pokrycia dla każdej metody.

Poniżej widoczny jest kod obrazujący działanie powyższego algorytmu.

```
1 calculate_cf <- function(alpha, n, p, M) {
2   size <- length(p)
3
4   n_cp <- numeric(size)
5   n_wald <- numeric(size)
6   n_wilson <- numeric(size)
7
8   p_cp <- numeric(size)
9   p_wald <- numeric(size)
10  p_wilson <- numeric(size)
11
12  for (k in 1:size) {
13    length_cp <- numeric(M)
14    length_wald <- numeric(M)
15    length_wilson <- numeric(M)
16
17    coverage_cp <- numeric(M)
18    coverage_wald <- numeric(M)
19    coverage_wilson <- numeric(M)
20
21    for (i in 1:M) {
22      sample_binom <- rbinom(n, 1, p[k])
23      successes <- sum(sample_binom)
24
25      clopper_pearson <- binom.confint(successes, n, conf.level = 1 -
26      alpha, methods = "exact")
27      wald <- binom.confint(successes, n, conf.level = 1 - alpha,
28      methods = "asymptotic")
29      wilson <- binom.confint(successes, n, conf.level = 1 - alpha,
30      methods = "wilson")
31      length_cp[i] <- clopper_pearson[, "upper"] - clopper_pearson[, "
32      lower"]
33      length_wald[i] <- wald[, "upper"] - wald[, "lower"]
34      length_wilson[i] <- wilson[, "upper"] - wilson[, "lower"]
35
36      coverage_cp[i] <- (clopper_pearson[, "lower"] <= p[k]) & (p[k] <=
37      clopper_pearson[, "upper"])
```

```

33     coverage_wald[i] <- (wald[, "lower"] <= p[k]) & (p[k] <= wald[, "
upper"])
34     coverage_wilson[i] <- (wilson[, "lower"] <= p[k]) & (p[k] <=
wilson[, "upper"])
35   }
36
37   n_cp[k] <- mean(length_cp)
38   n_wald[k] <- mean(length_wald)
39   n_wilson[k] <- mean(length_wilson)
40
41   p_cp[k] <- mean(coverage_cp)
42   p_wald[k] <- mean(coverage_wald)
43   p_wilson[k] <- mean(coverage_wilson)
44 }
45
46 result_n <- data.frame(
47   "p" = p,
48   "test Cloppera-Pearsona" = n_cp,
49   "test Walda" = n_wald,
50   "test Wilsona" = n_wilson
51 )
52
53 result_p <- data.frame(
54   "p" = p,
55   "test Cloppera-Pearsona" = p_cp,
56   "test Walda" = p_wald,
57   "test Wilsona" = p_wilson
58 )
59
60 return(list(result_n = result_n, result_p = result_p))
61 }
62
63 alpha <- 0.05
64 n <- c(30, 100, 1000)
65 p <- seq(0.01, 0.99, by = 0.01)
66 M <- 100
67
68 results_30 <- calculate_cf(alpha, n[1], p, M)
69 results_100 <- calculate_cf(alpha, n[2], p, M)
70 results_1000 <- calculate_cf(alpha, n[3], p, M)
71
72 df_length_30 <- melt(results_30$result_n, id.vars = "p", variable.name =
"Method", value.name = "Length")
73 df_coverage_30 <- melt(results_30$result_p, id.vars = "p", variable.name
= "Method", value.name = "Coverage")
74 df_length_100 <- melt(results_100$result_n, id.vars = "p", variable.name
= "Method", value.name = "Length")
75 df_coverage_100 <- melt(results_100$result_p, id.vars = "p", variable.
name = "Method", value.name = "Coverage")
76 df_length_1000 <- melt(results_1000$result_n, id.vars = "p", variable.
name = "Method", value.name = "Length")
77 df_coverage_1000 <- melt(results_1000$result_p, id.vars = "p", variable.
name = "Method", value.name = "Coverage")

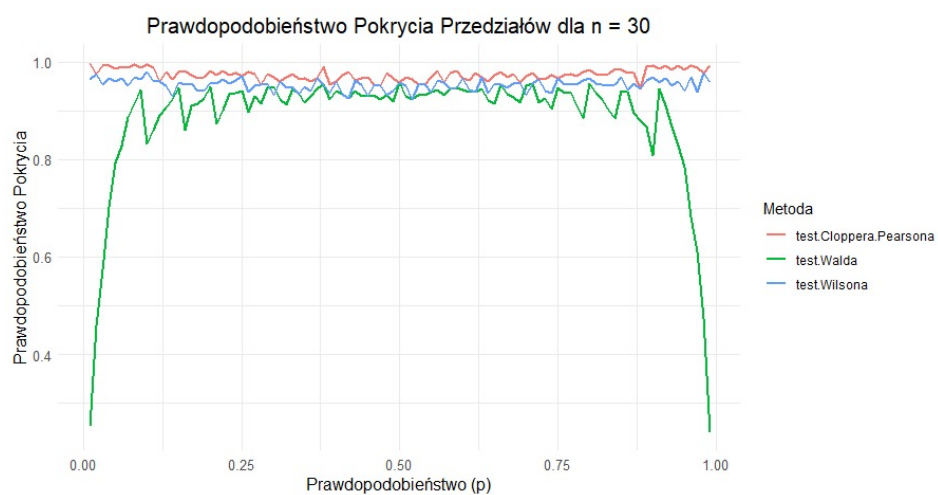
```

Poniżej widoczny jest kod służący do wizualizacji prawdopodobieństw pokrycia przedziałów ufności Cloppera-Pearsona, Walda i Wilsona dla liczby prób n równej odpowiednio 30, 100 i 1000 i parametru $p = 0.01, 0.02, \dots, 0.99$.

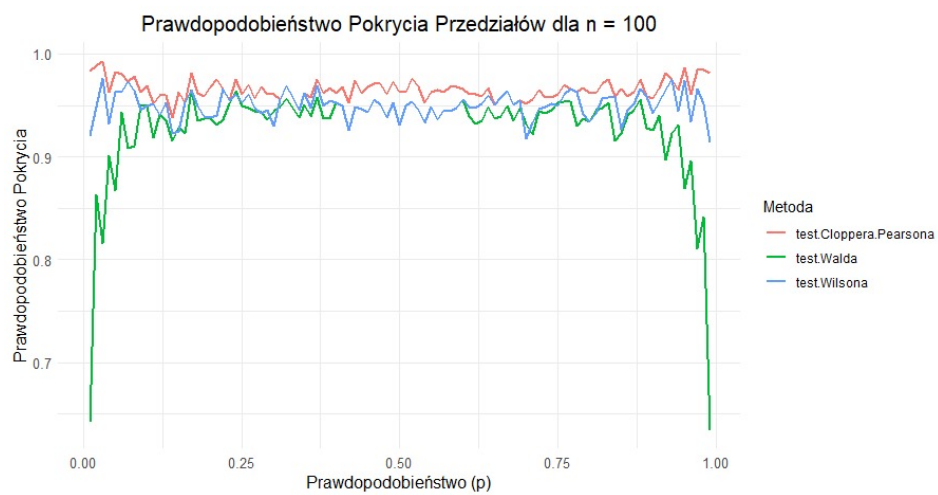
```

1 # wykresy prawdopodobieństwa pokrycia (czy)
2 ggplot(df_coverage_30, aes(x = p, y = Coverage, color = Method)) +
3   geom_line(linewidth = 0.8) +
4   labs(title = "Prawdopodobieństwo Pokrycia Przedziałów dla n = 30",
5         y = "Prawdopodobieństwo Pokrycia",
6         x = "Prawdopodobieństwo (p)",
7         color = "Metoda") +
8   theme_minimal() +
9   theme(
10     axis.title = element_text(size = 12),
11     axis.text = element_text(size = 10),
12     plot.title = element_text(size = 15, hjust = 0.5))
13
14
15 ggplot(df_coverage_100, aes(x = p, y = Coverage, color = Method)) +
16   geom_line(linewidth = 0.8) +
17   labs(title = "Prawdopodobieństwo Pokrycia Przedziałów dla n = 100",
18         y = "Prawdopodobieństwo Pokrycia",
19         x = "Prawdopodobieństwo (p)",
20         color = "Metoda") +
21   theme_minimal() +
22   theme(
23     axis.title = element_text(size = 12),
24     axis.text = element_text(size = 10),
25     plot.title = element_text(size = 15, hjust = 0.5))
26
27
28 ggplot(df_coverage_1000, aes(x = p, y = Coverage, color = Method)) +
29   geom_line(linewidth = 0.8) +
30   labs(title = "Prawdopodobieństwo Pokrycia Przedziałów dla n = 1000",
31         y = "Prawdopodobieństwo Pokrycia",
32         x = "Prawdopodobieństwo (p)",
33         color = "Metoda") +
34   theme_minimal() +
35   theme(
36     axis.title = element_text(size = 12),
37     axis.text = element_text(size = 10),
38     plot.title = element_text(size = 15, hjust = 0.5))

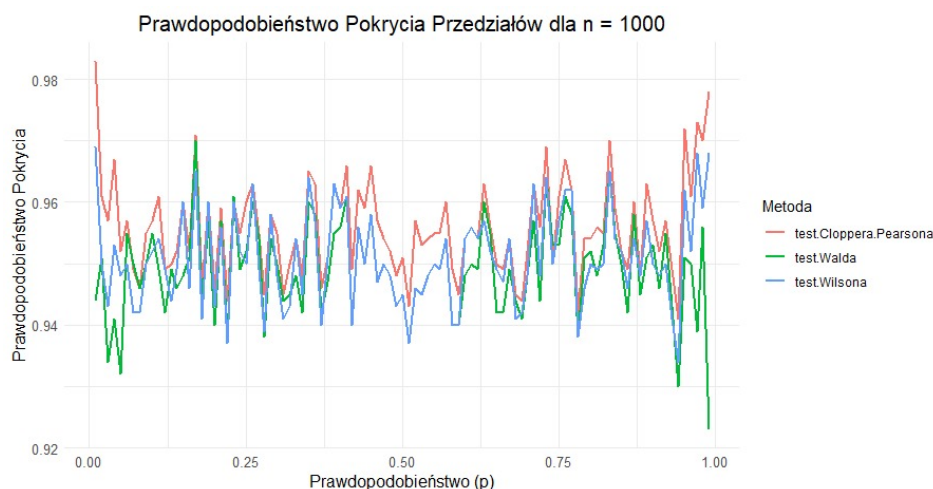
```

Rysunek 16: Prawdopodobieństwo pokrycia przedziałów dla $n = 30$



Rysunek 17: Prawdopodobieństwo pokrycia przedziałów dla $n = 100$



Rysunek 18: Prawdopodobieństwo pokrycia przedziałów dla $n = 1000$

Na wykresach 16, 17, 18 zwizualizowane zostały prawdopodobieństwa pokrycia przedziałów ufności Cloppera-Pearsona, Walda i Wilsona dla liczby prób n równej odpowiednio 30, 100 i 1000 i parametru $p = 0.01, 0.02, \dots, 0.99$.

Natomiast poniżej widoczny jest kod służący do wizualizacji średniej długości przedziałów ufności Cloppera-Pearsona, Walda i Wilsona dla liczby prób n równej odpowiednio 30, 100 i 1000 i parametru $p = 0.01, 0.02, \dots, 0.99$.

```

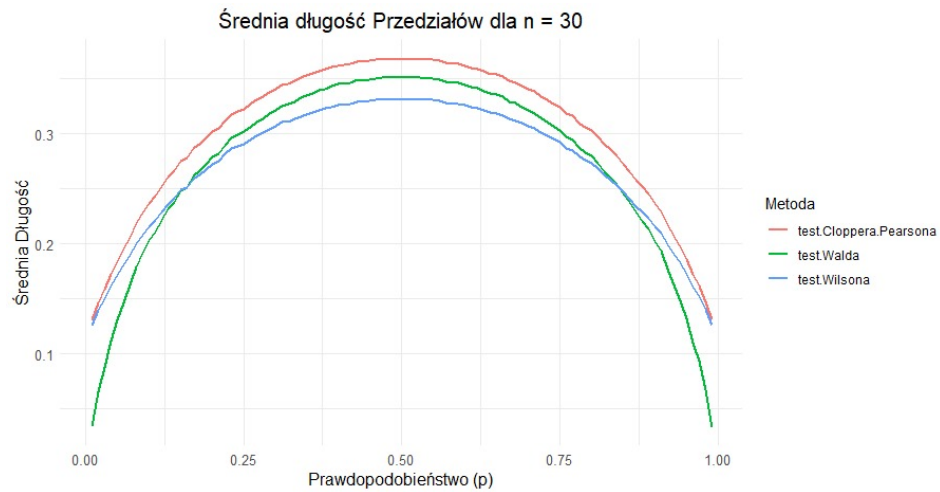
1 # wykresy średniej długości
2 ggplot(df_length_30, aes(x = p, y = Length, color = Method)) +
3   geom_line(linewidth = 0.8) +
4   labs(title = "Średnia długość Przedziałów dla n = 30",
5         y = "Średnia Długość",
6         x = "Prawdopodobieństwo (p)",
7         color = "Metoda") +
8   theme_minimal()+
9   theme(
10     axis.title = element_text(size = 12),
11     axis.text = element_text(size = 10),
12     plot.title = element_text(size = 15, hjust = 0.5))
13
14 ggplot(df_length_100, aes(x = p, y = Length, color = Method)) +
15   geom_line(linewidth = 0.8) +
16   labs(title = "Średnia długość Przedziałów dla n = 100",
17         y = "Średnia Długość",
18         x = "Prawdopodobieństwo (p)",
19         color = "Metoda") +
20   theme_minimal()+
21   theme(
22     axis.title = element_text(size = 12),
23     axis.text = element_text(size = 10),
24     plot.title = element_text(size = 15, hjust = 0.5))
25
26 ggplot(df_length_1000, aes(x = p, y = Length, color = Method)) +
27   geom_line(linewidth = 0.8) +
28   labs(title = "Średnia długość Przedziałów dla n = 1000",
29         y = "Średnia Długość",
30         x = "Prawdopodobieństwo (p)",

```

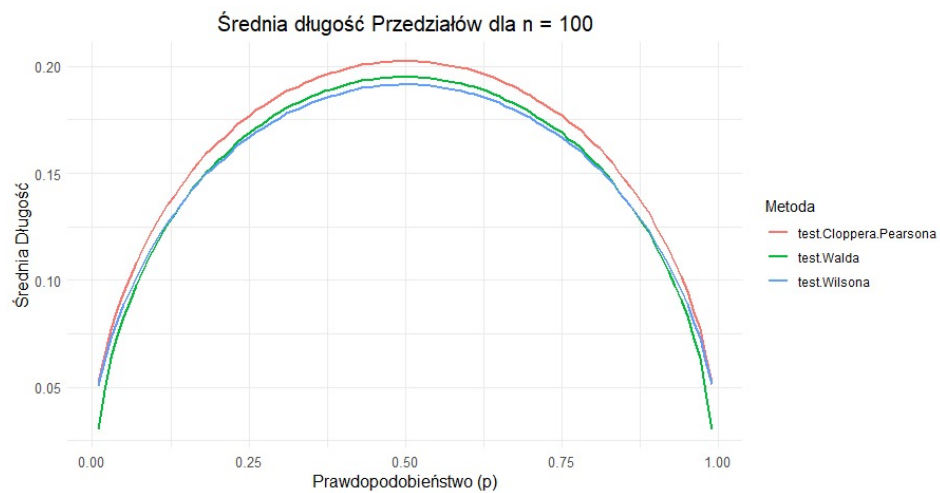
```

31     color = "Metoda") +
32     theme_minimal() +
33     theme(
34       axis.title = element_text(size = 12),
35       axis.text = element_text(size = 10),
36       plot.title = element_text(size = 15, hjust = 0.5))

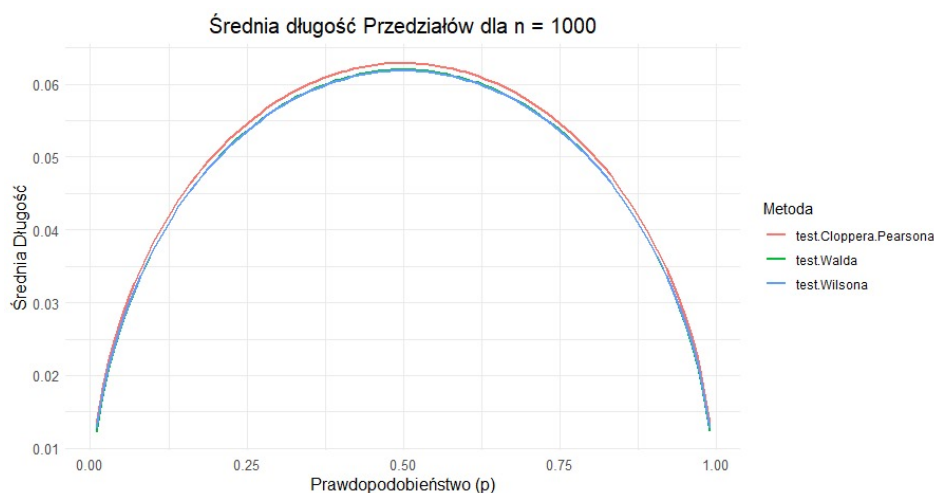
```



Rysunek 19: Średnia długość przedziałów przedziałów dla $n = 30$



Rysunek 20: Średnia długość przedziałów przedziałów dla $n = 100$



Rysunek 21: Średnia długość przedziałów przedziałów dla $n = 1000$

Na wykresach 19, 20, 21 zwizualizowane zostały średnie długości przedziałów ufności Cloppera-Pearsona, Walda i Wilsona dla liczby prób n równej odpowiednio 30, 100 i 1000 i parametru $p = 0.01, 0.02, \dots, 0.99$.

Na podstawie analizy wykresów 16 - 21, można stwierdzić, że metoda konstruowania przedziałów ufności Cloppera-Pearsona wykazuje największą dokładność dla środkowych wartości prawdopodobieństwa p oraz dla małych prób. W porównaniu do innych metod, ta metoda często osiąga wyższe prawdopodobieństwo pokrycia, co świadczy o uzyskaniu wysokiej dokładności w konstruowaniu przedziałów ufności, szczególnie dla małych prób.

Metoda Wilsona również okazała się skuteczna, prawdopodobieństwo pokrycia przedziału ufności jest dla niej również dość wysokie. Metoda ta nie jest tak konserwatywna jak metoda Cloppera-Pearsona, częściej niedoszacowuje wartości.

Z kolei metoda Walda okazała się najmniej skuteczna dla małych prób i skrajnych wartości prawdopodobieństwa, ale dla większych prób osiąga podobną skuteczność jak pozostałe metody.

Analizując średnie długości przedziałów ufności wyznaczone za pomocą metod Cloppera-Pearsona, Walda i Wilsona, dla różnych wielkości prób n i wartości prawdopodobieństw p , można zauważyć, że średnie długości przedziałów są najmniejsze dla skrajnych wartości prawdopodobieństwa oraz maleją wraz ze wzrostem wielkości próby (n). Największa różnica pomiędzy średnimi długościami przedziałów występuje dla małych prób ($n = 30$), natomiast dla dużych prób ($n = 1000$) różnice te są mniejsze.

Długości przedziałów różnią się między metodami, ale są do siebie zazwyczaj dość zbliżone. Wartości prawdopodobieństw również są do siebie zbliżone. Oznacza to, że różne metody symulacji przedziałów ufności prowadzą do podobnych wyników średnich (z przewagą w stronę metody Cloppera-Pearsona).

6 Część V

6.1 Zadanie 10

Zapoznaj się z funkcjami *binom.test* oraz *prop.test* z biblioteki *stats*.

```
1 help(binom.test) # test dokładny
2 help(prop.test) # test asymptotyczny
3 # correcter = TRUE - poprawione
```

W zadaniu 11 (6.2) oraz 12 (6.3) będziemy używać funkcji *binom.test* oraz *prop.test* z biblioteki *stats*.

Pierwsza z nich służy do przeprowadzania testu hipotez dotyczących rozkładu dwumianowego. Jej składnia ma następującą postać: *binom.test(x, n, p, alternative, conf.level)*, gdzie:

- *x* - liczba sukcesów w próbie,
- *n* - rozmiar próby,
- *p* - spodziewane prawdopodobieństwo sukcesu (wartość opcjonalna, jeśli nie jest podana to zostaje użyta wartość próbkowa (x/n),
- *alternative* - hipoteza alternatywna przyjmująca możliwe wartości: "two.sided" dwustronny, "less" jedenstronny mniejszy i "greater" jedenstronny większy,
- *conf.level* - poziom ufności dla przedziału ufności (domyślnie 0.95).

W wyniku działania tej funkcji otrzymujemy wyniki testu wraz z wartościami *p* i przedziałem ufności.

Natomiast funkcja *prop.test* służy do przeprowadzania testu hipotez dotyczących różnicy proporcji w dwóch grupach. Jej składnia ma następującą postać:

$$\text{prop.test}(x, n, \text{correct}, \text{alternative}, \text{conf.level}),$$

gdzie:

- *x* - wektor zawierający liczby sukcesów w próbie,
- *n* - wektor zawierający rozmiary prób,
- *correct* - określa czy ma być stosowana poprawka na ciągłość zalecana dla małych prób,
- *alternative* - hipoteza alternatywna przyjmująca możliwe wartości: "two.sided" dwustronny, "less" jedenstronny mniejszy i "greater" jedenstronny większy,
- *conf.level* - poziom ufności dla przedziału ufności (domyślnie 0.95).

6.2 Zadanie 11

Dla danych z pliku "ankieta.csv" korzystając z funkcji z zadania 10 (6.1), przyjmując $1 - \alpha = 0.95$, zweryfikuj następujące hipotezy i sformułuj wnioski:

1. Prawdopodobieństwo, że w firmie pracuje kobieta wynosi 0.5.
2. Prawdopodobieństwo, że pracownik jest zadowolony ze swojego wynagrodzenia w pierwszym badanym okresie jest większe bądź równe 0.7.
3. Prawdopodobieństwo, że kobieta pracuje na stanowisku menedżerskim jest równe prawdopodobieństwu, że mężczyzna pracuje na stanowisku menedżerskim.
4. Prawdopodobieństwo, że kobieta jest zadowolona ze swojego wynagrodzenia w pierwszym badanym okresie jest równe prawdopodobieństwu, że mężczyzna jest zadowolony ze swojego wynagrodzenia w pierwszym badanym okresie.
5. Prawdopodobieństwo, że kobieta pracuje w dziale obsługi kadrowo-płacowej jest większe lub równe prawdopodobieństwu, że mężczyzna pracuje w dziale obsługi kadrowo-płacowej.

Dla każdej z hipotez przeprowadziliśmy *binom.test*, *prop.test* bez poprawki oraz *prop.test* z poprawką na ciągłość. Następnie sprawdziliśmy czy otrzymane p-wartości są większe od poziomu istotności α , który w naszym przypadku wynosi 0.05.

Przypominając: jeżeli poziom krytyczny (p-wartość) jest mniejszy bądź równy od poziomu istotności α to hipotezę odrzucamy, w przeciwnym przypadku nie mamy podstaw do jej odrzucenia.

```
1 # dla części z podpunktów
2 females <- sum(df$PŁEĆ == "K")
3 males <- sum(df$PŁEĆ == "M")
4 length_sex <- c(females, males)
5 all_rows <- nrow(df)
```

- HIPOTEZA 1: Prawdopodobieństwo, że w firmie pracuje kobieta wynosi 0.5.

```
1 # 1.Prawdopodobieństwo, że w firmie pracuje kobieta wynosi 0.5.
2 binom_test <- binom.test(females, all_rows, p=0.5, alternative = "two.
  sided", conf.level = 0.95)
3 prop_test_1 <- prop.test(females, all_rows, p=0.5, alternative = "two.
  sided", conf.level = 0.95, correct = TRUE)
4 prop_test_2 <- prop.test(females, all_rows, p=0.5, alternative = "two.
  sided", conf.level = 0.95, correct = FALSE)
5
6 print(binom_test$p.value)
7 print(prop_test_1$p.value)
8 print(prop_test_2$p.value)
```

W przypadku tej hipotezy otrzymaliśmy następujące p-wartości:

- *binom.test*: 4.972973×10^{-5}
- *prop.test* bez poprawki: 5.565628×10^{-5}
- *prop.test* z poprawką: 4.109788×10^{-5}

Dla każdego z testów p-wartość jest mniejsza niż poziom istotności $\alpha = 0.05$ co sugeruje, że należy odrzucić hipotezę zerową. Sugeruje to, że prawdopodobieństwo, że w firmie pracuje kobieta jest istotnie statystycznie różne od 0.5.

- HIPOTEZA 2: Prawdopodobieństwo, że pracownik jest zadowolony ze swojego wynagrodzenia w pierwszym badanym okresie jest większe bądź równe 0.7.

```

1 # 2. Prawdopodobieństwo, że pracownik jest zadowolony ze swojego
  wynagrodzenia w pierwszym badanym okresie jest większe bądź równe 0.7.
2 content_1 <- sum(df$CZY_ZADOW == "YES")
3
4 binom_test <- binom.test(content_1, all_rows, p=0.7, alternative = "
  greater", conf.level = 0.95)
5 prop_test_1 <- prop.test(content_1, all_rows, p=0.7, alternative = "
  greater", conf.level = 0.95, correct = TRUE)
6 prop_test_2 <- prop.test(content_1, all_rows, p=0.7, alternative = "
  greater", conf.level = 0.95, correct = FALSE)
7
8 print(binom_test$p.value)
9 print(prop_test_1$p.value)
10 print(prop_test_2$p.value)
11
12 content_2 <- sum(df$CZY_ZADOW_2 == "YES")
13
14 binom_test <- binom.test(content_2, all_rows, p=0.7, alternative = "
  greater", conf.level = 0.95)
15 prop_test_1 <- prop.test(content_2, all_rows, p=0.7, alternative = "
  greater", conf.level = 0.95, correct = TRUE)
16 prop_test_2 <- prop.test(content_2, all_rows, p=0.7, alternative = "
  greater", conf.level = 0.95, correct = FALSE)
17
18 print(binom_test$p.value)
19 print(prop_test_1$p.value)
20 print(prop_test_2$p.value)

```

W przypadku tej hipotezy, zarówno dla zmiennej **CZY_ZADOW**, jak i dla zmiennej **CZY_ZADOW_2**, dla wszystkich testów otrzymaliśmy p-wartość równą 1. Pozwala nam to na przyjęcie hipotezy zerowej.

- HIPOTEZA 3: Prawdopodobieństwo, że kobieta pracuje na stanowisku menedżerskim jest równe prawdopodobieństwu, że mężczyzna pracuje na stanowisku menedżerskim.

```

1 # 3. Prawdopodobieństwo, że kobieta pracuje na stanowisku menedżerskim
  jest równe prawdopodobieństwu, że mężczyzna pracuje na stanowisku
  menedżerskim.
2 manager <- filter(df, CZY_KIER == "Tak")
3 female_manager <- sum(manager$PŁEĆ == "K")
4 male_manager <- sum(manager$PŁEĆ == "M")
5
6 fem_male_manager <- c(female_manager, male_manager)
7
8 prop_test_1 <- prop.test(fem_male_manager, length_sex, alternative = "two
  .sided", conf.level = 0.95, correct = TRUE)
9 prop_test_2 <- prop.test(fem_male_manager, length_sex, alternative = "two
  .sided", conf.level = 0.95, correct = FALSE)
10
11 print(prop_test_1$p.value)
12 print(prop_test_2$p.value)

```

W przypadku tej hipotezy otrzymaliśmy następujące p-wartości:

- *prop.test* bez poprawki: 0.6389361

- *prop.test* z poprawką: 0.4930904

Możemy zauważyć, że p-wartości dla obu testów są większe od poziomu istotności $\alpha = 0.05$. Pozwala nam to przyjąć hipotezę zerową i stwierdzić, że najprawdopodobniej nie istnieje istotna statystycznie różnica pomiędzy zajmowaniem stanowisk menadżerskich przez kobiety i mężczyzn.

- HIPOTEZA 4: Prawdopodobieństwo, że kobieta jest zadowolona ze swojego wynagrodzenia w pierwszym badanym okresie jest równe prawdopodobieństwu, że mężczyzna jest zadowolony ze swojego wynagrodzenia w pierwszym badanym okresie.

```

1 # 4. Prawdopodobieństwo, że kobieta jest zadowolona ze swojego
  wynagrodzenia w pierwszym badanym okresie jest równe prawdopodobień
  stwu, że mężczyzna jest zadowolony ze swojego wynagrodzenia w
  pierwszym badanym okresie.
2 content <- filter(df, CZY_ZADOW == "TAK")
3 female_content <- sum(content$PŁEĆ == "K")
4 male_content <- sum(content$PŁEĆ == "M")
5
6 fem_male_content <- c(female_content, male_content)
7
8 prop_test_1 <- prop.test(fem_male_content, length_sex, alternative = "two
  .sided", conf.level = 0.95, correct = TRUE)
9 prop_test_2 <- prop.test(fem_male_content, length_sex, alternative = "two
  .sided", conf.level = 0.95, correct = FALSE)
10
11 print(prop_test_1$p.value)
12 print(prop_test_2$p.value)

```

W przypadku tej hipotezy otrzymaliśmy następujące p-wartości:

- *prop.test* bez poprawki: 0.7379521
- *prop.test* z poprawką: 0.6293763

Możemy zauważyć, że p-wartości dla obu testów są większe od poziomu istotności $\alpha = 0.05$. Pozwala nam to przyjąć hipotezę zerową i stwierdzić, że najprawdopodobniej nie istnieje istotna statystycznie różnica pomiędzy zadowoleniem ze swojego wynagrodzenia kobiet oraz mężczyzn w pierwszym badanym okresie.

- HIPOTEZA 5: Prawdopodobieństwo, że kobieta pracuje w dziale obsługi kadrowo-płacowej jest większe lub równe prawdopodobieństwu, że mężczyzna pracuje w dziale obsługi kadrowo-płacowej.

```

1 # 5. Prawdopodobieństwo, że kobieta pracuje w dziale obsługi kadrowo-pł
  acowej jest większe lub równe prawdopodobieństwu, że mężczyzna pracuje
  w dziale obsługi kadrowo-płacowej.
2 HR <- filter(df, DZIAŁ == "HR")
3 female_HR <- sum(HR$PŁEĆ == "K")
4 male_HR <- sum(HR$PŁEĆ == "M")
5
6 fem_male_HR <- c(female_HR, male_HR)
7
8 prop_test_1 <- prop.test(fem_male_HR, length_sex, alternative = "greater"
  , conf.level = 0.95, correct = TRUE)
9 prop_test_2 <- prop.test(fem_male_HR, length_sex, alternative = "greater"
  , conf.level = 0.95, correct = FALSE)
10 print(prop_test_1$p.value)
11 print(prop_test_2$p.value)

```


W przypadku tej hipotezy otrzymaliśmy następujące p-wartości:

- *prop.test* bez poprawki: 0.9960475
- *prop.test* z poprawką: 0.9978835

Możemy zauważyć, że p-wartości dla obu testów są większe od poziomu istotności $\alpha = 0.05$ i bliskie 1. Pozwala nam to przyjąć hipotezę zerową i stwierdzić, że prawdopodobieństwo, że kobieta pracuje w dziale obsługi kadrowo-płacowej jest większe lub równe prawdopodobieństwu, że mężczyzna pracuje w tym dziale.

6.3 Zadanie 12

Wyznacz symulacyjnie moc testu dokładnego oraz moc testu asymptotycznego w przypadku weryfikacji hipotezy zerowej $H_0 : p = 0.9$ przeciwko $H_1 : p \neq 0.9$ przyjmując wartość $1 - \alpha = 0.95$. Uwzględnij różne wartości alternatyw i różne rozmiary próby. Sformułuj wnioski.

W tym zadaniu zajęliśmy się wyznaczeniem mocy testu dokładnego oraz asymptotycznego z poprawką i bez dla $n \in \{30, 100, 1000\}$ i $p = 0.01, 0.02, \dots, 0.89, 0.91, \dots, 0.99$. Dla przypomnienia: Moc testu to prawdopodobieństwo odrzucenia fałszywej hipotezy zerowej i przyjęcia prawdziwej hipotezy alternatywnej.

W celu obliczenia mocy testów wykorzystaliśmy poniższy algorytm:

1. Inicjalizacja zmiennych:

- (a) Określenie wartości hipotezy zerowej H_0 .
- (b) Ustalenie poziomu istotności α .
- (c) Określenie wielkości próby n .
- (d) Zdefiniowanie wektora prawdopodobieństw p .
- (e) Określenie liczby symulacji M .

2. Obliczenie mocy testu:

- (a) Dla każdego n :
 - i. Dla każdej wartości prawdopodobieństwa sukcesu p :
 - A. Dla każdej z M symulacji:
 - Generujemy próbkę o rozmiarze n z rozkładu dwumianowego, w tym celu używamy funkcji *rbinom*.
 - Przeprowadzamy testy z użyciem funkcji *binom.test* i *prop.test*.
 - Zapisujemy czy odrzuciliśmy hipotezę zerową czy nie.
 - B. Obliczamy średnią z liczby odrzuconych hipotez zerowych dla każdego z testów.
 - ii. Zapisujemy wyniki
- (b) Przedstawiamy wyniki za pomocą wykresów.

```

1 calculate_power <- function(H_0, alpha, n, p, M) {
2   size <- length(p)
3
4   result_binom <- numeric(size)
5   result_prop_1 <- numeric(size)
6   result_prop_2 <- numeric(size)
7
8   for (k in 1:size) {
9     rejection_binom <- 0
10    rejection_prop_1 <- 0
11    rejection_prop_2 <- 0
12
13    for (i in 1:M) {
14      sample_binom <- rbinom(n, 1, p[k])
15      summ <- sum(sample_binom)
16
17      binom <- binom.test(summ, n, p=H_0, alternative = "two.sided", conf
.level = 1 - alpha)
18      prop_1 <- prop.test(summ, n, p=H_0, alternative = "two.sided", conf
.level = 1 - alpha, correct = TRUE)
19      prop_2 <- prop.test(summ, n, p=H_0, alternative = "two.sided", conf
.level = 1 - alpha, correct = FALSE)
20
21      binom_pv <- binom$p.value
22      prop_pv_1 <- prop_1$p.value
23      prop_pv_2 <- prop_2$p.value
24
25      if (binom_pv < alpha) {
26        rejection_binom <- rejection_binom + 1
27      }
28      if (prop_pv_1 < alpha) {
29        rejection_prop_1 <- rejection_prop_1 + 1
30      }
31      if (prop_pv_2 < alpha) {
32        rejection_prop_2 <- rejection_prop_2 + 1
33      }
34
35    }
36    result_binom[k] <- rejection_binom/M
37    result_prop_1[k] <- rejection_prop_1/M
38    result_prop_2[k] <- rejection_prop_2/M
39  }
40  return(list(result_binom = result_binom, result_prop_1 = result_prop_1,
    result_prop_2 = result_prop_2))
41 }
42
43 alpha <- 0.05
44 H_0 <- 0.9
45 n <- c(30, 100, 1000)
46 p <- seq(0.01, 0.99, by = 0.01)
47 M <- 1000
48
49 result_30 <- calculate_power(H_0, alpha, n[1], p, M)
50 result_100 <- calculate_power(H_0, alpha, n[2], p, M)
51 result_1000 <- calculate_power(H_0, alpha, n[3], p, M)
52
53 df_30 <- data.frame(p = p, moc_binom = result_30$result_binom, moc_prop_1
  = result_30$result_prop_1, moc_prop_2 = result_30$result_prop_2)
54 df_100 <- data.frame(p = p, moc_binom = result_100$result_binom, moc_prop

```

```

    _1 = result_100$result_prop_1, moc_prop_2 = result_100$result_prop_2)
55 df_1000 <- data.frame(p = p, moc_binom = result_1000$result_binom, moc_
    prop_1 = result_1000$result_prop_1, moc_prop_2 = result_1000$result_
    prop_2)

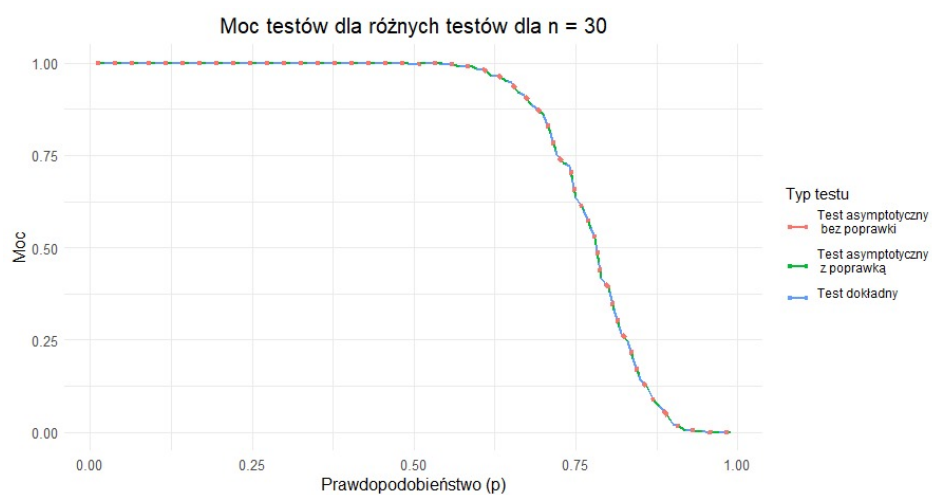
```

Poniżej widoczne są kody służące do wizualizacji wyników odpowiednio dla $n = 30$, $n = 100$ i $n = 1000$.

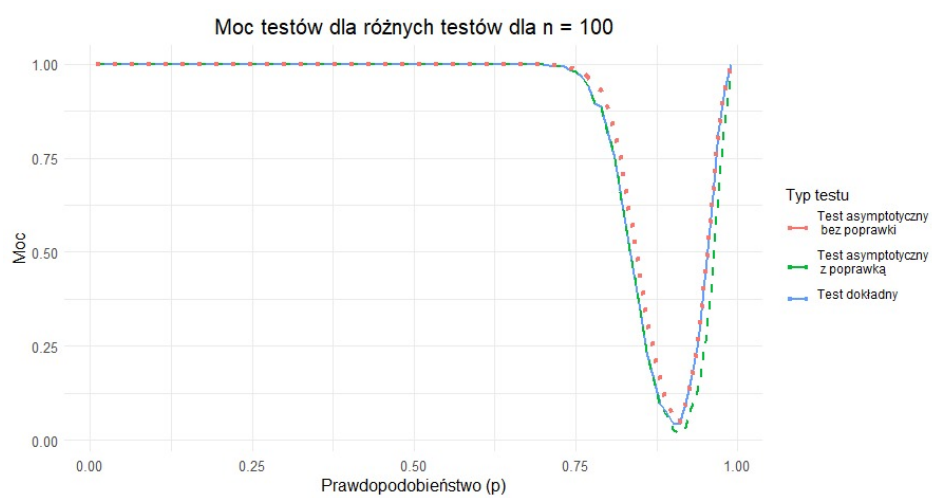
```

1 ggplot(df_30, aes(x = p)) +
2   geom_line(aes(y = moc_binom, color = "Test dokładny\n"), linewidth =
   0.8) +
3   geom_line(aes(y = moc_prop_1, color = "Test asymptotyczny\n z poprawką\
   n"), linetype = "dashed", linewidth = 1) +
4   geom_line(aes(y = moc_prop_2, color = "Test asymptotyczny\n z bez
   poprawki\n"), linetype = "dotted", linewidth = 1.5) +
5   labs(title = "Moc testów dla różnych testów dla n = 30",
6         x = "Prawdopodobieństwo (p)",
7         y = "Moc",
8         color = "Typ testu") +
9   theme_minimal() +
10  theme(
11    axis.title = element_text(size = 12),
12    axis.text = element_text(size = 10),
13    plot.title = element_text(size = 15, hjust = 0.5))
14
15 ggplot(df_100, aes(x = p)) +
16   geom_line(aes(y = moc_binom, color = "Test dokładny\n"), linewidth =
   0.8) +
17   geom_line(aes(y = moc_prop_1, color = "Test asymptotyczny\n z poprawką\
   n"), linetype = "dashed", linewidth = 1) +
18   geom_line(aes(y = moc_prop_2, color = "Test asymptotyczny\n z bez
   poprawki\n"), linetype = "dotted", linewidth = 1.5) +
19   labs(title = "Moc testów dla różnych testów dla n = 100",
20         x = "Prawdopodobieństwo (p)",
21         y = "Moc",
22         color = "Typ testu") +
23   theme_minimal() +
24   theme(
25     axis.title = element_text(size = 12),
26     axis.text = element_text(size = 10),
27     plot.title = element_text(size = 15, hjust = 0.5))
28
29 ggplot(df_1000, aes(x = p)) +
30   geom_line(aes(y = moc_binom, color = "Test dokładny\n"), linewidth =
   0.8) +
31   geom_line(aes(y = moc_prop_1, color = "Test asymptotyczny\n z poprawką\
   n"), linetype = "dashed", linewidth = 1) +
32   geom_line(aes(y = moc_prop_2, color = "Test asymptotyczny\n z bez
   poprawki\n"), linetype = "dotted", linewidth = 1.5) +
33   labs(title = "Moc testów dla różnych testów dla n = 1000",
34         x = "Prawdopodobieństwo (p)",
35         y = "Moc",
36         color = "Typ testu") +
37   theme_minimal() +
38   theme(
39     axis.title = element_text(size = 12),
40     axis.text = element_text(size = 10),
41     plot.title = element_text(size = 15, hjust = 0.5))

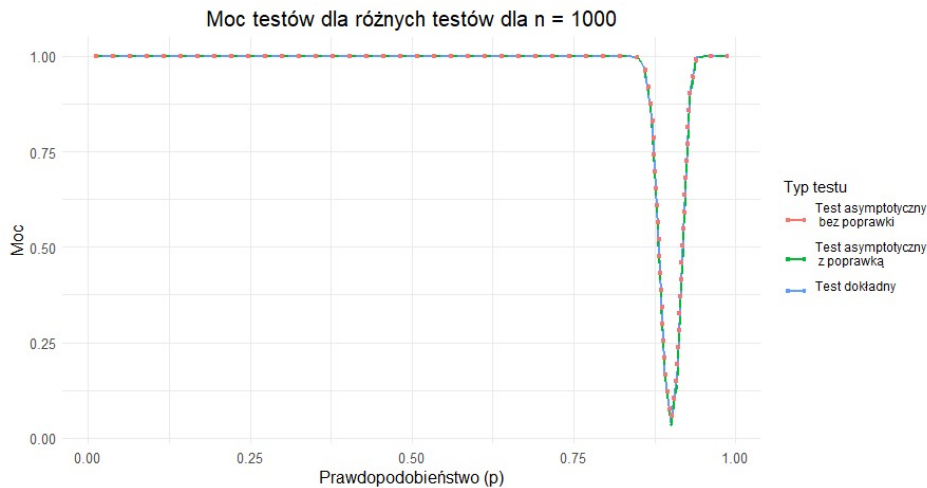
```



Rysunek 22: Moc testu dla różnych testów dla $n = 30$



Rysunek 23: Moc testu dla różnych testów dla $n = 100$



Rysunek 24: Moc testu dla różnych testów dla $n = 1000$

Analizując wygląd wykresów 22, 23, 24 możemy zauważyć, że wraz ze wzrostem n dołek dokładniej pokrywa się z wartością prawdopodobieństwa z hipotezy zerowej, równą 0.9.

Podsumowując, na podstawie analizy wyników możemy wnioskować, że nie istnieje test najlepszy w każdej sytuacji. Wybór pomiędzy testem dokładnym i testem asymptotycznym zależy od konkretnych warunków badawczych. Test dokładny będzie lepszym wyborem, jeśli zależy nam na kontrolowaniu ryzyka popełnienia błędu I-go rodzaju. Natomiast w przypadku większych prób oraz gdy zależy nam na większej mocy powinniśmy wybrać testy asymptotyczne.

7 Zadanie dodatkowe

7.1 Zadanie *1

Wyznacz granice asymptotycznego przedziału ufności dla prawdopodobieństwa sukcesu bazując na przekształceniu logit korzystając z metody delta. Zaimplementuj metodę oraz porównaj wyniki z funkcją *binom.confint*.

W tym zadaniu użyliśmy przekształcenia logitowego, metody delta oraz sposobu, w jaki są one wykorzystywane do wyznaczania asymptotycznych przedziałów ufności dla prawdopodobieństwa sukcesu.

Poniżej znajduje się wyjaśnienie kluczowych pojęć:

- **Przekształcenie logitowe** - przekształcenie logitowe prawdopodobieństwa sukcesu p jest zdefiniowane jako:

$$\text{logit}(p) = \log \left(\frac{p}{1-p} \right),$$

gdzie p - prawdopodobieństwem sukcesu.

- **Metoda delta** - jest to technika wykorzystywana do aproksymacji rozkładu funkcji estymatora. Jeśli mamy estymator $\hat{\theta}$ z wariancją $\text{Var}(\hat{\theta})$ i funkcję g , której chcemy

użyć na $\hat{\theta}$, to wariancja $g(\hat{\theta})$ może być aproksymowana przez:

$$\text{Var}(g(\hat{\theta})) \approx [g'(\hat{\theta})]^2 \text{Var}(\hat{\theta}),$$

gdzie $g'(\hat{\theta})$ - pierwsza pochodna g względem $\hat{\theta}$.

- **Asymptotyczny przedział ufności dla p** - założmy, że mamy n prób binarnych (np. sukces/porażka), z których X to liczba sukcesów. Estymatorem p jest $\hat{p} = \frac{X}{n} = \frac{1}{n} \sum_{i=1}^n X_i$, $X_i \sim B(n, p)$. Przekształcenie logitowe tego estymatora to $\text{logit}(\hat{p})$. Aby zastosować metodę delta do wyznaczenia przedziału ufności dla p , potrzebujemy najpierw pochodnej przekształcenia logitowego względem p , co daje:

$$\frac{d}{dp} \text{logit}(p) = \frac{1}{p(1-p)}$$

Zatem wariancja $\text{logit}(\hat{p})$ może być aproksymowana przez:

$$\text{Var}(\text{logit}(\hat{p})) \approx \left[\frac{1}{\hat{p}(1-\hat{p})} \right]^2 \text{Var}(\hat{p}) = \left[\frac{1}{\hat{p}(1-\hat{p})} \right]^2 \frac{\hat{p}(1-\hat{p})}{n} = \frac{1}{n\hat{p}(1-\hat{p})}$$

$$(*) \text{Var}(\hat{p}) = \text{Var}\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n^2} \sum_{i=1}^n \text{Var}(X_i) \stackrel{iid}{=} \frac{1}{n^2} n\hat{p}(1-\hat{p}) = \frac{\hat{p}(1-\hat{p})}{n}$$

Wtedy:

$$Q = \frac{g(\hat{\theta}) - g(\theta)}{g'(\theta) \sqrt{\frac{\text{Var}(\hat{\theta})}{n}}},$$

gdzie: $Q \sim N(0, 1)$

Mając to, możemy wyznaczyć asymptotyczny przedział ufności dla $\text{logit}(p)$ używając kwantyli rozkładu normalnego $q_{1-\frac{\alpha}{2}}$, gdzie α to poziom istotności.

$$P(\theta \in [-q_{1-\frac{\alpha}{2}}; q_{1-\frac{\alpha}{2}}]) = 1 - \alpha$$

Dokonując odpowiednich przekształceń otrzymujemy:

$$\text{logit}(\hat{p}) \pm q_{1-\alpha/2} \sqrt{\frac{1}{n\hat{p}(1-\hat{p})}}$$

Po rozwinięciu $\text{logit}(\hat{p})$ do powyższego wzoru otrzymujemy:

$$\text{logit}(p) \in \left[\ln \left(\frac{\hat{p}}{1-\hat{p}} \right) \pm q_{1-\alpha/2} \frac{1}{\sqrt{n\hat{p}(1-\hat{p})}} \right]$$

Aby otrzymać przedział ufności dla p , odwracamy przekształcenie logitowe i podstawiamy do wzoru powyżej:

$$p = \frac{e^{\text{logit}(p)}}{1 + e^{\text{logit}(p)}}$$

Następnie, w celu sprawdzenia poprawności zaimplementowanej przez nas metody dokonaliśmy porównania z funkcją *binom.confint*.

Funkcja ta oferuje kilka metod wyznaczania przedziałów ufności dla p , w tym metodę asymptotyczną, która jest zbliżona do opisanej powyżej, ale może korzystać z różnych aproksymacji. Zaimplementowana przez nas metoda oparta na przekształceniu logitowym i metodzie delta jest jedną z metod asymptotycznych w związku z czym przypuszczamy, że powinna dawać podobne wyniki do tych uzyskanych za pomocą *binom.confint* z odpowiednią opcją metody.

Poniżej widoczny jest kod prezentujący metody opisane powyżej.

```

1 logit_confint_delta <- function(successes, trials, alpha) {
2   p_hat <- successes / trials
3   logit_p_hat <- log(p_hat / (1 - p_hat))
4
5   # Granice dla przekształcenia logitowego
6   q <- qnorm(1 - (alpha/2))
7
8   part_logit = sqrt(1/(trials*p_hat*(1 - p_hat)))
9
10  logit_lcl <- logit_p_hat - q * part_logit
11  logit_ucl <- logit_p_hat + q * part_logit
12
13
14  # Odwrócenie przekształcenia logitowego do p
15  lcl <- exp(logit_lcl) / (1 + exp(logit_lcl))
16  ucl <- exp(logit_ucl) / (1 + exp(logit_ucl))
17
18  return(c(lower = lcl, upper = ucl))
19 }
20
21
22 alpha <- 0.05
23 conf_level <- 1 - alpha
24
25 trials <- 1000
26 sample_trial <- sample(c(1, 0), trials, replace = TRUE)
27 successes <- sum(sample_trial)
28
29 delta_confint <- logit_confint_delta(successes, trials, alpha)
30 binom_confint <- binom.confint(successes, trials, conf_level, methods = c
31   ("exact", "asymptotic"))
32 cat("Metoda delta (logit):\n")
33 print(delta_confint)
34 cat("\nFunkcja binom.confint:\n")
35 print(binom_confint)

```

W wyniku działania powyższego kodu otrzymaliśmy przedziały ufności dla zaimplementowanej przez nas metody oraz dla funkcji wbudowanej.

Metoda	Dolna granica	Górna granica
Zaimplementowana metoda	0.4800152	0.5419005
Wbudowana - asymptotic	0.4800177	0.5419823
Wbudowana - exact	0.4795242	0.5424109

Tabela 20: Granice asymptotyczne przedziałów ufności

Jak możemy zauważyć w tabeli 20 granice asymptotyczne przedziałów ufności wyznaczone metodą zaimplementowaną przez nas oraz funkcją wbudowaną *binom.confint* są do

siebie bardzo zbliżone. Pozwala nam to wnioskować, że metoda zaimplementowana przez nas najprawdopodobniej jest poprawna.

8 Źródła

- Wykłady dr hab. inż. Krzysztofa Burneckiego oraz laboratoria dr inż. Aleksandry Grzesiek z przedmiotu „Statystyka stosowana”.
- Wykłady dr inż. Aleksandry Grzesiek oraz laboratoria inż. Huberta Woszcza z przedmiotu „Analiza danych ankietowych”.
- Wykłady dr inż. Andrzeja Giniewicza oraz laboratoria dr inż. Agnieszki Kamińskiej z przedmiotu "Pakiety statystyczne".
- <https://www.rdocumentation.org/>
- <https://www.r-project.org/other-docs.html>