

Lista 3

Celem niniejszego raportu jest pogłębiona analiza danych ankietowych z wykorzystaniem wybranych metod statystycznych. Zakres analizy obejmuje zarówno testy symetrii dla danych sparowanych, jak i symulacyjne porównanie mocy testów statystycznych, a także modelowanie zależności między zmiennymi za pomocą modeli log-liniowych.

W pierwszej części skupiliśmy się na testach symetrii, w tym na implementacji warunkowego testu symetrii oraz jego zastosowaniu do rzeczywistych danych dotyczących skuteczności leków i oceny szkoleń. Druga część obejmuje porównanie mocy dwóch testów przy różnych rozmiarach prób, co pozwala ocenić ich efektywność w różnych warunkach badawczych.

W dalszej części przeanalizowaliśmy wybrane zbiory danych z wcześniejszych list zadań, weryfikując hipotezy o symetrii rozkładu odpowiedzi oraz przeprowadzając testy związane z możliwymi zmianami opinii respondentów w czasie. W kolejnych sekcjach zastosowaliśmy modele log-liniowe do opisu zależności między zmiennymi takimi jak stanowisko kierownicze, opinia o szkoleniach i staż pracy, uwzględniając również alternatywne modele oraz ich porównania za pomocą kryteriów AIC i BIC.

Raport kończy analiza zjawiska paradoksu Simpsona na podstawie danych dotyczących skuteczności dwóch metod leczenia oraz zadania dodatkowe związane z dokładnymi testami symetrii i wyborem najlepszego modelu dla wskazanych zmiennych.

Część I i II

Zadanie 1

Napisz funkcję, która zwraca p-wartość w omówionym na wykładzie warunkowym teście symetrii w przypadku tabeli 2×2 .

```
p <- function(n12, n21){
  part <- 0
  if(n12 < (n12+n21)/2){
    for(i in 0:n12){
      part <- part + choose(n12+n21,i)*(1/2)^i*(1/2)^(n12+n21-i)
    }
    part <- 2*part
  }
  if(n12 > (n12+n21)/2){
    for(i in 0:n12+n21){
      part <- part + choose(n12+n21,i)*(1/2)^i*(1/2)^(n12+n21-i)
    }
  }
}
```

```

    part <- 2*part
  }
  if(n12==(n12+n21)/2){
    part <- 1
  }
  return(part)
}

```

Zadanie 2

W tabeli 1 umieszczono dane dotyczące reakcji na lek po godzinie od jego przyjęcia dla dwóch różnych leków przeciwbólowych stosowanych w migrenie. Leki zostały zaaplikowane grupie pacjentów w dwóch różnych atakach bólowych. Na podstawie danych zweryfikuj hipotezę, że leki te są jednakowo skuteczne korzystając z testu.

- McNemara z poprawką na ciągłość,

- warunkowego (korzystając z funkcji zadeklarowanej w zadaniu 1.).

Zadanie 2.1

Lek A	Lek B	
	Negatywna	Pozytywna
Negatywna	1	5
Pozytywna	2	4

McNemar's Chi-squared test with continuity correction

data: tabela

McNemar's chi-squared = 0.57143, df = 1, p-value = 0.4497

W celu porównania skuteczności dwóch leków przeciwbólowych stosowanych podczas dwóch różnych ataków migreny u tych samych pacjentów, zastosowano test McNemara z poprawką na ciągłość. Uzyskałyśmy w ten sposób p-value = 0.4447, co jest większe niż $\alpha = 0.05$, więc brakuje podstaw do odrzucenia hipotezy zerowej mówiącej, że leki są jednakowo skuteczne.

Zadanie 2.2

[1] 0.453125

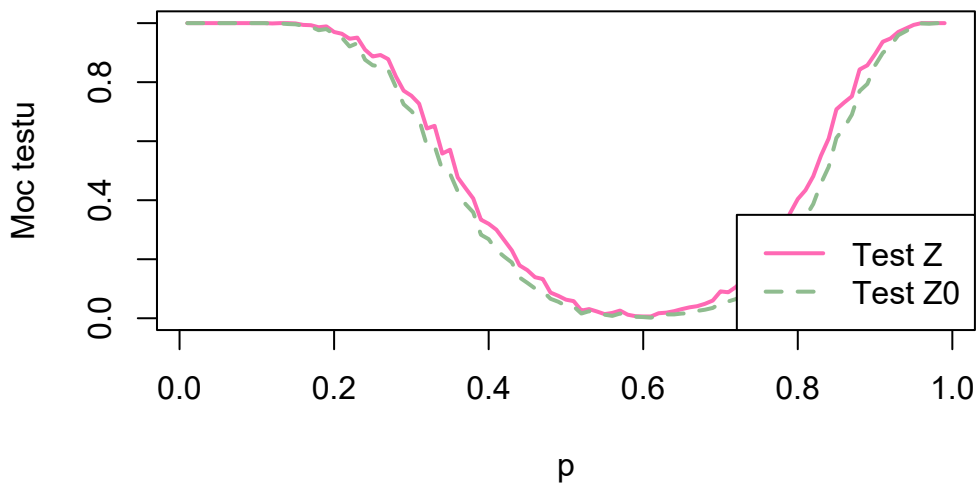
[1] 0.453125

W celu porównania skuteczności dwóch leków przeciwbólowych zastosowano test warunkowy w postaci testu dwumianowego (na podstawie 7 przypadków rozbieżnych odpowiedzi). Otrzymano wartość $p\text{-value} = 0.4531$, co jest znacznie większe od przyjętego poziomu istotności $= 0.05$. Nie ma podstaw do odrzucenia hipotezy zerowej, zakładającej jednakową skuteczność leków. Na podstawie testu warunkowego nie wykazano istotnej różnicy w skuteczności obu leków przeciwbólowych stosowanych w migrenie.

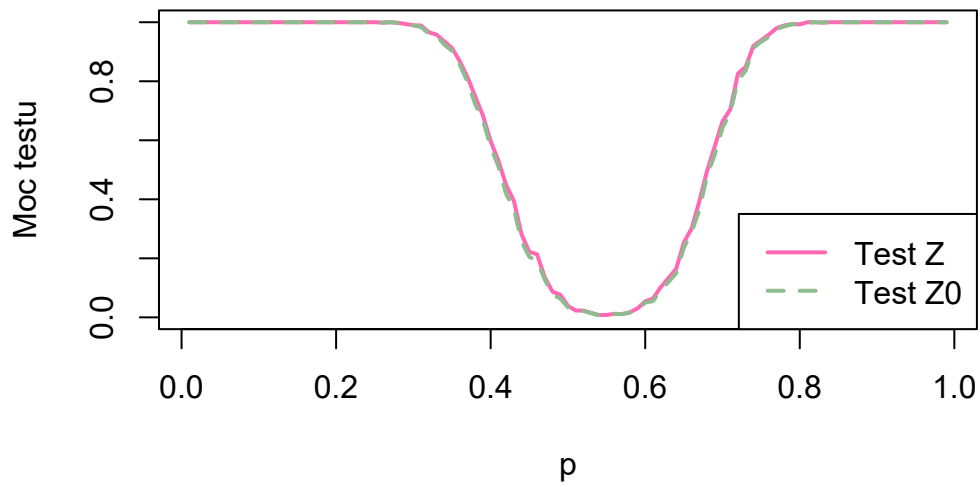
Zadanie 3

Przeprowadź symulacje w celu porównania mocy testu Z i testu Z_0 przedstawionych na wykładzie. Rozważ różne długości prób.

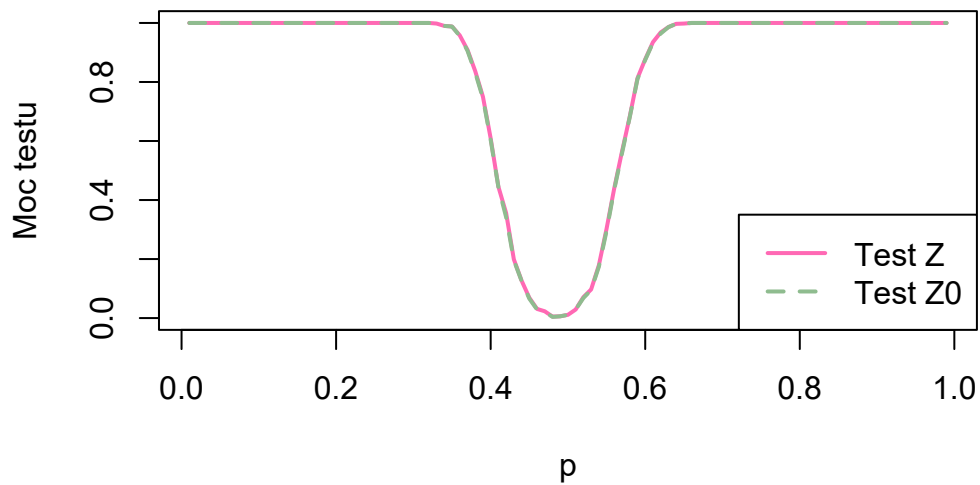
Wykres mocy testów Z i Z_0 dla $n=30$



Wykres mocy testów Z i Z_0 dla $n=100$



Wykres mocy testów Z i Z_0 dla $n=300$



Widzimy, że dla mniejszych n test Z ma większą moc od testu Z_0 . Dla większych n moce testów zbliżają się do siebie oraz rosną, szczególnie wokół $p = 0.5$.

Zadanie 4

Dla danych dołączonych do pierwszej listy zadań, na podstawie zmiennych CZY_ZADW oraz CZY_ZADW_2, zweryfikuj hipotezę, że zadowolenie ze szkoleń w pierwszym badanym okresie i w drugim badanym okresie pierwszego badania odpowiada modelowi symetrii. Czy na podstawie uzyskanych wyników możemy wnioskować, że poziom zadowolenia ze szkoleń nie uległ zmianie? Przyjmij poziom istotności 0.05.

Cel zadania Sprawdzić, czy dwa różne leki przeciwbólowe podawane tym samym pacjentom w dwóch atakach migreny są jednakowo skuteczne. Testowana hipoteza zerowa: oba leki są jednakowo skuteczne.

McNemar's Chi-squared test with continuity correction

data: tabela

McNemar's chi-squared = 4.3214, df = 1, p-value = 0.03764

Ponieważ $p < 0.05$, istnieją statystyczne podstawy do odrzucenia hipotezy zerowej, zakładającej brak zmiany w poziomie zadowolenia. Można przyjąć, że zadowolenie respondentów ze szkoleń uległo istotnej zmianie między pierwszym a drugim okresem badania.

Zadanie 5

Zadanie 5 dotyczy oceny skuteczności działań wdrożonych w firmie mających na celu poprawę komfortu pracy oraz efektywności wykorzystania wiedzy zdobytej na szkoleniach. W tym celu przeprowadzono badanie ankietowe w dwóch okresach: przed wdrożeniem zmian oraz po ich wprowadzeniu. Respondenci zostali poproszeni o ocenę podejścia firmy do umożliwiania praktycznego wdrażania zdobytej wiedzy.

Zebrane dane zostały przedstawione w postaci tablicy dwudzielczej, w której odpowiedzi z pierwszego okresu zostały zestawione z odpowiedziami z drugiego okresu. Oceny przyjmują wartości od -2 do 2, gdzie wyższe wartości oznaczają bardziej pozytywną ocenę. Celem analizy jest weryfikacja hipotezy, że rozkład odpowiedzi w obu okresach jest symetryczny, co odpowiadałoby brakowi zmiany w ocenach.

	-2	-1	0	1	2
-2	10	2	1	1	0
-1	0	15	1	1	0
0	1	1	32	6	0
1	0	0	1	96	3
2	1	1	0	1	26

McNemar's Chi-squared test

```
data: dane  
McNemar's chi-squared = NaN, df = 10, p-value = NA
```

W celu weryfikacji hipotezy o symetrii odpowiedzi w dwóch okresach badania zastosowaliśmy rozszerzoną wersję testu McNemara, czyli test Bowkera, odpowiednią dla tablic większych niż 2×2 . W analizie wykorzystano funkcję `mcnemar.test()`, jednak test zwrócił wartość statystyki NaN oraz `p-value = NA`. Taki wynik jest zgodny z teoretycznymi założeniami testu Bowkera, który opiera się na porównaniu częstości odpowiedzi w pozycjach symetrycznych względem głównej przekątnej (n_{ij} vs n_{ji}). W przypadku, gdy suma tych par (tj. $n_{ij} + n_{ji}$) wynosi zero, powstaje niedozwolona operacja dzielenia przez zero w obliczeniach statystyki testowej.

Z tego względu klasyczny test Bowkera nie może być zastosowany w tej sytuacji i posłużymy się alternatywnym podejściem, testem LW, który lepiej radzi sobie w obecności zerowych komórek poza przekątną.

```
$statistic  
[1] 13.32669
```

```
$df  
[1] 10
```

```
$p.value  
[1] 0.2059752
```

```
$method  
[1] "Test LW"
```

```
$data.name  
[1] "dane"
```

Wnioski P-value wynosi 0.206, co jest większe od przyjętego poziomu istotności $\alpha = 0.05$. Oznacza to, że brak podstaw do odrzucenia hipotezy o symetrii rozkładu odpowiedzi w dwóch badanych okresach. Nie stwierdzono istotnej statystycznie zmiany oceny podejścia firmy do umożliwiania wdrażania wiedzy zdobytej na szkoleniach. Odpowiedzi respondentów przed i po wdrożeniu działań poprawiających komfort pracy można uznać za zgodne z modelem symetrii.

Część III

Zadanie 6

W pewnym badaniu porównywano skuteczność dwóch metod leczenia: Leczenie A to nowa procedura, a Leczenie B to stara procedura. Przeanalizuj dane przedstawione w Tabeli 3 (wyniki dla całej grupy pacjentów) oraz w Tabelach 4 i 5 (wyniki w podgrupach ze względu na dodatkową zmienną) i odpowiedz na pytanie, czy dla danych występuje paradoks Simpsona.

Table 1: Tabela 3: Dane dla całej grupy

Metoda	Poprawa	Brak poprawy
Leczenie A	117	104
Leczenie B	177	44

Table 2: Tabela 4: Dane dla pacjentów z chorobami współistniejącymi.

Metoda	Poprawa	Brak poprawy
Leczenie A	17	101
Leczenie B	2	36

Table 3: Tabela 5: Dane dla pacjentów bez chorób współistniejących.

Metoda	Poprawa	Brak poprawy
Leczenie A	100	3
Leczenie B	175	8

```
wszyscy z chorobami bez chorób
A 0.5294118 0.14406780 0.9708738
B 0.8009050 0.05263158 0.9562842
```

Chociaż leczenie B “wygrywa” patrząc na całą grupę badanych, po podziale na grupy ze względu na obecność chorób współistniejących możemy zauważyć, że to leczenie A ma większy odsetek wyzdrowień.

```
wszyscy z chorobami bez chorób
2.740007e-09 2.248419e-01 7.675118e-01
```

W przeprowadzonym teście niezależności χ^2 dla całej grupy p-value jest bardzo małe, więc odrzucamy hipotezę H_0 o niezależności. Jednak ten sam test wykonany osobno dla badanych grup - z chorobami współistniejącymi oraz bez chorób - w obu przypadkach daje p-value większą od poziomu istotności, a więc nie mamy podstaw do odrzucania hipotezy zerowej o niezależności zmiennych, to znaczy wyniku leczenia (poprawy) od przyjętego leczenia. To znaczy, że pozorny związek dla całej badanej grupy nie przekłada się na zależność w podgrupach - a więc jest to klasyczny przypadek paradoksu Simpsona.

Zadanie 7

Dla danych z listy 1, przyjmując za zmienną 1 zmienną CZY_KIER, za zmienną 2 – zmienną PYT_2 i za zmienną 3 – zmienną STAŻ, podaj interpretacje następujących modeli log-liniowych: [1 3], [13], [1 2 3], [12 3], [12 13] oraz [1 23].

[1 3] - zmienne CZY_KIER oraz STAŻ są niezależne,

[13] - zmienne CZY_KIER oraz STAŻ nie są niezależne,

[1 2 3] - zmienne CZY_KIER, PYT_2 oraz STAŻ są niezależne,

[12 3] - zmienne CZY_KIER i PYT_2 nie są niezależne, a zmienna STAŻ jest niezależna od nich obu,

[12 13] - zmienne CZY_KIER i PYT_2 nie są niezależne, CZY_KIER i STAŻ nie są niezależne, a PYT_2 i STAŻ są warunkowo niezależne,

[1 23] - zmienna CZY_KIER jest niezależna od pozostałych dwóch, PYT_2 i STAŻ, które nie są od siebie niezależne.

Część IV i V

Zadanie 8

Przyjmując model log-liniowy [123] dla zmiennych opisanych w zadaniu 7 oszacuj prawdopodobieństwa:

- że osoba pracująca na stanowisku kierowniczym jest zdecydowanie zadowolona ze szkoleń;
- że osoba o stażu pracy krótszym niż rok pracuje na stanowisku kierowniczym;
- że osoba o stażu pracy powyżej trzech lat nie pracuje na stanowisku kierowniczym.

Jakie byłyby oszacowania powyższych prawdopodobieństw przy założeniu modelu [12 23]?

Zaczynamy od modelu [123]:


```
# A tibble: 4 x 5
  PYT_2 freq_sum fitted_sum p_dane p_model
  <fct>   <int>      <dbl> <dbl>   <dbl>
1 -2         10      10.0  0.370   0.370
2 -1          2       2.00  0.0741  0.0741
3 1           2       2.00  0.0741  0.0741
4 2          13      13.0  0.481   0.481
```

```
# A tibble: 2 x 5
  CZY_KIER freq_sum fitted_sum p_dane p_model
  <fct>      <int>      <dbl> <dbl>   <dbl>
1 Nie        40      40.0  0.976   0.976
2 Tak         1       1.00  0.0244  0.0244
```

```
# A tibble: 2 x 5
  CZY_KIER freq_sum fitted_sum p_dane p_model
  <fct>      <int>      <dbl> <dbl>   <dbl>
1 Nie        10      10.0  0.526   0.526
2 Tak         9       9.00  0.474   0.474
```

Model [123] dobrze oszacował potrzebne prawdopodobieństwa (w 1. tabeli interesuje nas wiersz z odpowiedzią “2” na PYT_2, w 2. i 3. odpowiedź “Tak” w kolumnie CZY_KIER). Zarówno szacowane licznosci jak i prawdopodobieństwa są równe dla modelu i danych.

```
# A tibble: 4 x 5
  PYT_2 freq_sum fitted_sum p_dane p_model
  <fct>   <int>      <dbl> <dbl>   <dbl>
1 -2         10      10.0  0.370   0.370
2 -1          2       2.00  0.0741  0.0741
3 1           2       2.00  0.0741  0.0741
4 2          13      13.0  0.481   0.481
```

```
# A tibble: 2 x 5
  CZY_KIER freq_sum fitted_sum p_dane p_model
  <fct>      <int>      <dbl> <dbl>   <dbl>
1 Nie        40      35.7  0.976   0.872
2 Tak         1       5.25  0.0244  0.128
```

```
# A tibble: 2 x 5
  CZY_KIER freq_sum fitted_sum p_dane p_model
  <fct>      <int>      <dbl> <dbl>   <dbl>
```

1 Nie	10	14.8	0.526	0.778
2 Tak	9	4.22	0.474	0.222

Dla modelu [12 23] odpowiedź na pierwszy podpunkt się zgadza - wartości w danych są równe przewidywanym przez model. Jednak przy pytaniach, które łączą zmienne `CZY_KIER` oraz `STAZ` (podpunkt 2. i 3.) model przeszacował wyniki z dla osób o krótkim stażu oraz niedoszacował odpowiedzi w dla osób o długim stażu - wynika to z braku powiązanie między tymi zmiennymi. Jak widzimy, złe dobranie modelu skutkuje złym oszacowaniem badanych prawdopodobieństw.

Zadanie 9

Dla danych wskazanych w zadaniu 7 zweryfikuj następujące hipotezy:

- Zmienne losowe **CZY_KIER**, **PYT_2** i **STAZ** są wzajemnie niezależne,
- Zmienna losowa **PYT_2** jest niezależna od pary zmiennych **CZY_KIER** i **STAZ**,
- Zmienna losowa **PYT_2** jest niezależna od zmiennej **CZY_KIER**, przy ustalonej wartości zmiennej **STAZ**.

Celem zadania jest weryfikacja zależności pomiędzy trzema zmiennymi kategorycznymi: `CZY_KIER` (czy osoba pracuje na stanowisku kierowniczym), `PYT_2` (ocena szkolenia) oraz `STAZ` (długość stażu pracy). Analiza oparta jest na modelach log-liniowych, które umożliwiają badanie zarówno pełnej niezależności pomiędzy wszystkimi zmiennymi, jak i zależności marginalnych oraz warunkowych.

a) Wzajemna niezależność [1][2][3]

Rozważmy 3 nadmodele, np. [12 23], [13 23], [123]

Warning: glm.fit: dopasowane stosunki numerycznie okazały się być 0

Analysis of Deviance Table

Model 1: LICZBA ~ CZY_KIER + PYT_2 + STAZ

Model 2: LICZBA ~ CZY_KIER * PYT_2 + PYT_2 * STAZ

	Resid. Df	Resid. Dev	Df	Deviance	Pr(>Chi)
1	17	42.242			
2	8	15.642	9	26.601	0.001628 **

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Analysis of Deviance Table

Model 1: LICZBA ~ CZY_KIER + PYT_2 + STAZ

Model 2: LICZBA ~ CZY_KIER * STAZ + PYT_2 * STAZ

	Resid. Df	Resid. Dev	Df	Deviance	Pr(>Chi)
1	17	42.242			
2	9	4.880	8	37.362	9.871e-06 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Analysis of Deviance Table

Model 1: LICZBA ~ CZY_KIER + PYT_2 + STAZ

Model 2: LICZBA ~ CZY_KIER * PYT_2 * STAZ

	Resid. Df	Resid. Dev	Df	Deviance	Pr(>Chi)
1	17	42.242			
2	0	0.000	17	42.242	0.0006187 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

1. [12 23] vs [1][2][3]

- Deviance: 26.601
- df: 9
- p-value: 0.001628

2. [13 23] vs [1][2][3]

- Deviance: 37.362
- df: 8
- p-value: 9.871e-06

3. [123] vs [1][2][3]

- Deviance: 42.242
- df: 17
- p-value: 0.0006187

Ponieważ we wszystkich porównaniach p-value są mniejsze niż 0.05, odrzucamy hipotezę zerową o wzajemnej niezależności zmiennych. Wnioskujemy, że zmienne te nie są wzajemnie niezależne – występują między nimi istotne statystycznie zależności.

- b) Zmienna PYT_2 jest niezależna od pary zmiennych CZY_KIER i STAŻ, czyli model ma postać log-liniową: [2 13]

Warning: glm.fit: dopasowane stosunki numerycznie okazały się być 0

Analysis of Deviance Table

```
Model 1: LICZBA ~ PYT_2 + CZY_KIER * STAZ
Model 2: LICZBA ~ CZY_KIER * PYT_2 + PYT_2 * STAZ
  Resid. Df Resid. Dev Df Deviance Pr(>Chi)
1         15      23.152
2          8      15.642  7    7.5105  0.3777
```

Analysis of Deviance Table

```
Model 1: LICZBA ~ PYT_2 + CZY_KIER * STAZ
Model 2: LICZBA ~ CZY_KIER * STAZ + PYT_2 * STAZ
  Resid. Df Resid. Dev Df Deviance Pr(>Chi)
1         15      23.152
2          9       4.880  6    18.272 0.005587 **
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Analysis of Deviance Table

```
Model 1: LICZBA ~ PYT_2 + CZY_KIER * STAZ
Model 2: LICZBA ~ CZY_KIER * PYT_2 * STAZ
  Resid. Df Resid. Dev Df Deviance Pr(>Chi)
1         15      23.152
2          0       0.000 15    23.152 0.08096 .
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

1. [12 23] vs [2 13]

- Deviance: 7.5105
- df: 7
- p-value: 0.3777 → brak podstaw do odrzucenia hipotezy zerowej

2. [13 23] vs [2 13]

- Deviance: 18.272
- df: 6
- p-value: 0.0056 → odrzucamy hipoteze zerową

3. [123] vs [2 13]

- Deviance: 23.152
- df: 10
- p-value: 0.08096 → brak podstaw do odrzucenia hipotezy zerowej

Na podstawie porównań modelu zakładającego niezależność zmiennej PYT_2 od pary zmiennych CZY_KIER i STAZ z trzema bardziej złożonymi nadmodelami, można stwierdzić, że: - tylko w jednym przypadku (model [13 23]) uzyskano istotną statystycznie poprawę dopasowania ($p = 0.0056 < 0.05$), - natomiast dla modeli [12 23] oraz [123] wartości p były większe niż 0.05, co oznacza brak istotnej poprawy dopasowania.

Nie ma wystarczających statystycznych podstaw do odrzucenia hipotezy, że zmienna PYT_2 jest niezależna od pary zmiennych CZY_KIER i STAZ na poziomie istotności $\alpha = 0.05$.

c) [13 23]

Warning: glm.fit: dopasowane stosunki numerycznie okazały się być 0

Analysis of Deviance Table

Model 1: LICZBA ~ CZY_KIER * STAZ + PYT_2 * STAZ

Model 2: LICZBA ~ CZY_KIER * PYT_2 + PYT_2 * STAZ + CZY_KIER * STAZ

	Resid. Df	Resid. Dev	Df	Deviance	Pr(>Chi)
1	9	4.8800			
2	6	1.5967	3	3.2832	0.35

Analysis of Deviance Table

Model 1: LICZBA ~ CZY_KIER * STAZ + PYT_2 * STAZ

Model 2: LICZBA ~ CZY_KIER * PYT_2 * STAZ

	Resid. Df	Resid. Dev	Df	Deviance	Pr(>Chi)
1	9	4.88			
2	0	0.00	9	4.88	0.8446

1. [12 13 23] vs [12 23]

- Deviance: 3.2832
- df: 3
- p-value: 0.35 → brak podstaw do odrzucenia hipotezy zerowej

2. [123] vs [1 23]

- Deviance: 4.88
- df: 9
- p-value: 0.8446 → brak podstaw do odrzucenia hipotezy zerowej

W obu przypadkach wartości p są znacznie większe od przyjętego poziomu istotności $\alpha = 0.05$, co oznacza, że nie ma statystycznych podstaw do odrzucenia hipotezy zerowej. Innymi słowy, nie stwierdzono istotnych zależności pomiędzy odpowiedziami na pytanie PYT_2 a statusem kierowniczym, jeżeli uwzględnimy staż pracy.

Zadania dodatkowe

Zadanie 1*

W przypadku zadania 5 występuje problem z zastosowaniem testu Bowkera ze względu na występowanie zer na określonych miejscach w tabeli z danymi. Zastosuj w tym przypadku dokładny test symetrii i opisz, w jaki sposób wyznaczana jest wartość poziomu krytycznego w tym teście.

Dlatego, że w macierzy mamy obecność zer, to test Bowkera jest niewłaściwy. Stosujemy test dokładny, który sprawdza, czy odpowiedzi są symetryczne względem przekątnej.

	-2	-1	0	1	2
-2	10	2	1	1	0
-1	0	15	1	1	0
0	1	1	32	6	0
1	0	0	1	96	3
2	1	1	0	1	26

[1] 0.2948

W przypadku, gdy klasyczny test Bowkera nie może zostać zastosowany ze względu na obecność zer w tabeli kontyngencji, stosuje się tzw. dokładny test symetrii. Test ten opiera się na permutacyjnym podejściu do weryfikacji hipotezy zerowej mówiącej, że badana tablica jest symetryczna względem głównej przekątnej.

W tym podejściu najpierw obliczana jest statystyka testowa W , która mierzy całkowitą niesymetryczność tabeli — w tym przypadku jest to suma bezwzględnych różnic pomiędzy odpowiadającymi sobie elementami symetrycznymi względem przekątnej. Następnie, zakładając prawdziwość hipotezy zerowej (czyli że rozkład odpowiedzi jest symetryczny), przeprowadza się symulację wielu możliwych losowych, symetrycznych tablic. W każdej iteracji dla każdej pary komórek symetrycznych (i,j) i (j,i) rozdziela się ich łączną sumę na dwa składniki zgodnie z rozkładem dwumianowym o parametrze $p = 0.5$ (czyli zakładając równą szansę odpowiedzi po każdej stronie symetrii). Dla każdej takiej zasymulowanej tablicy oblicza się wartość statystyki testowej.

P-value wynosi 0.2948 co jest większe niż 0.05, a więc nie mamy podstaw do odrzucenia hipotezy symetrii. Na podstawie tego testu nie ma podstaw, żeby twierdzić, że opinie pracowników zmieniły się po działaniach firmy.

Zadanie 2*

Na podstawie danych z listy 1 dokonaj wyboru modelu rozważając uwzględnienie zmiennych PYT_1, PYT_2 i PŁEĆ w oparciu o:

- testy,
- kryterium AIC,
- kryterium BIC.

Będzimy roważać modele [1 2 3], [12 13 23] oraz [123].

Analysis of Deviance Table

```
Model 1: Freq ~ PYT_1 + PYT_2 + PŁEĆ
Model 2: Freq ~ PYT_1 * PYT_2 + PYT_1 * PŁEĆ + PYT_2 * PŁEĆ
Model 3: Freq ~ PYT_1 * PYT_2 * PŁEĆ
  Resid. Df Resid. Dev Df Deviance Pr(>Chi)
1         31    226.057
2         12      7.365 19   218.692 <2e-16 ***
3          0      0.000 12     7.365  0.8326
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

W tej analizie testujemy, czy model prostszy wystarcza (H_0), czy potrzebny jest model bardziej złożony (H_1). Testujemy, czy sensownie jest zmieniać model z [1 2 3] na [12 13 23] oraz [12 13 23] na [123]. Wykonujemy test istotności χ^2 .

Interpretacja:

Model 1 (niezależność) jest zbyt prosty, ponieważ po dodaniu interakcji dwójkowych (Model 2) dopasowanie znacznie się poprawia ($p < 2e-16$ więc odrzucamy H_0).

Model 3 (pełny) nie poprawia istotnie dopasowania względem Modelu 2 ($p = 0.8326$, nie ma podstaw do odrzucenia H_0), więc interakcja trójkowa nie jest potrzebna.

Ostateczny wybór: Model 2 – zawiera wszystkie istotne interakcje (dwójkowe), a jest prostszym modelem niż pełny.

	Model	AIC	BIC
1	model_full	150.1856	217.7408
2	model_12_13_23	133.5509	180.8396
3	model_indep	314.2426	329.4425

Porównując **AIC** oraz **BIC** widzimy, że dla obu kryteriów model [12 13 23] przyjmuje najmniejsze wartości, więc dla tego porównania jest najlepszy.

Zarówno testy chi-kwadrat, jak i kryteria AIC/BIC wskazują, że najlepszym modelem jest model z interakcjami dwójkowymi: $\text{Freq} \sim \text{PYT}_1 * \text{PYT}_2 + \text{PYT}_1 * \text{PŁEĆ} + \text{PYT}_2 * \text{PŁEĆ}$, oznaczany jako **[12 13 23]**.