

TABLE NO. IV.

Cause of Death.	Admissions.	Deaths.	1839—1853.	Officers		Non-Com
				Strength	Deaths	
(continued).						
...	264	23				
...	590	161				
Aspnoea	55	1				
Diseases	—	—				
Order 4.						
...	29	8				
...	36	11				
...	16	9				
...	—*	—*				
Intestines	—*	—*				
...	101	2				
...	1,862	5				
...	1	1				
Intestines	—*	—*				
...	129	3				
...	906	3				
...	358	0				
...	15	2				
...	1	0				
...	—*	—*				
...	251	17				
...	878	22				
Diseases	—*	—*				
...	9	1				
Order 5.						
Nephria)	26	2				
...	39	0				
(above)	—	—				
...	8	1				
...	—*	—*				
...	1	0				
...	9	1				
Urethra	139	2				
...	15	0				
...	2	0				
Order 6.						
(able to the Army.)						
Order 7.						
...	87	0				
Prostitis	7	0				
...	25	1				
...	2	0				

Tomasz Przechlewski

# Podstawy statystyki

Podręcznik dla studentów wydziałów nauk o zdrowiu

Powiślańska Szkoła Wyższa • Kwidzyn-Gdańsk

Na okładce: Statystyki zdrowia armii brytyjskiej z okresu wojny krymskiej  
Notes on matters affecting the health, efficiency,  
and hospital administration of the British Army  
by Nightingale, Florence, 1820-1910  
<https://archive.org/details/b20387118/>

Książka została zredagowana w formacie Rmarkdown/bookdown .

Kolorowa wersja podręcznika znajduje się pod adresem:  
[https://hrpunio.github.io/SMI\\_Bookdown/wstęp.html](https://hrpunio.github.io/SMI_Bookdown/wstęp.html)

Recenzent: Michał Pietrzak

# Spis treści

<b>Wstęp</b>	<b>5</b>
<b>1 Przedmiot i metody badań statystycznych</b>	<b>6</b>
1.1 Przedmiot statystyki . . . . .	6
1.2 Podstawowe pojęcia . . . . .	6
1.3 Pomiar . . . . .	7
1.4 Rodzaje i sposoby analizy danych . . . . .	7
1.5 Sposoby pomiaru danych i organizacja badania . . . . .	8
1.6 Miary częstości chorób . . . . .	12
1.7 Oprogramowanie . . . . .	13
<b>2 Analiza jednej zmiennej</b>	<b>14</b>
2.1 Tablice statystyczne . . . . .	14
2.2 Wykresy . . . . .	17
2.3 Statystyczka Florence Nightingale . . . . .	20
2.4 Analiza parametryczna . . . . .	20
2.5 Porównanie wielu rozkładów . . . . .	25
2.6 Zestawienie metod opisu statystycznego . . . . .	29
<b>3 Wprowadzenie do wnioskowania statystycznego</b>	<b>30</b>
3.1 Masa ciała uczestników PŚ w rugby . . . . .	30
3.2 Wiek kandydatów na radnych . . . . .	32
3.3 Rozkład normalny . . . . .	35
3.4 Odsetek kobiet wśród kandydatów na radnych . . . . .	36
3.5 Wnioskowanie statystyczne . . . . .	38
3.6 Statystyk Carl Pearson . . . . .	40
3.7 Słownik terminów, które warto znać . . . . .	40
<b>4 Analiza współzależności pomiędzy zmiennymi</b>	<b>42</b>
4.1 Dwie zmienne nominalne . . . . .	43
4.2 Zmienna liczbowa i zmienna nominalna . . . . .	48
4.3 Dwie zmienne liczbowe . . . . .	53
4.4 Zmienna liczbowa i zmienne liczbowe lub nominalne . . . . .	61
4.5 Przypadek specjalny: regresja logistyczna . . . . .	67
4.6 Przypadek specjalny: co najmniej dwie zmienne porządkowe . . . . .	72
4.7 Podsumowanie . . . . .	72
<b>5 Przykłady badań ankietowych</b>	<b>73</b>
5.1 Jak zacząć badanie? . . . . .	73
5.2 Wiedza na temat szkodliwości palenia i jej uwarunkowania wśród studentów PSW . . . . .	76
5.3 Depresja i jej uwarunkowania wśród studentów PSW . . . . .	82
5.4 Satysfakcja, przywiązanie i zamiar odejścia . . . . .	87
5.5 Formularze ankiet . . . . .	90

<b>6</b>	<b>Praca z programem Jamovi</b>	<b>97</b>
6.1	Podstawy pracy z Jamovi . . . . .	97
6.2	Analiza ankiety: satysfakcja – wiedza o paleniu – zamiar odejścia . . .	98
	<b>Literatura</b>	<b>111</b>

# Wstęp

Przygotowując się kilka lat temu do zajęć z przedmiotu Statystyka Medyczna ze zdumieniem odkryłem, że nie istnieje coś takiego jak polski podręcznik do statystyki medycznej powszechnie dostępny w księgarni, który mogę polecić swoim studentom. Byłem tym tym bardziej zdziwiony, że wiedziałem już wtedy co oznacza akronim EBM. W tamtym czasie prawdę mówiąc była tylko jedna książka ze statystyką medyczną w tytule – absolutnie moim zdaniem nie nadająca się do polecenia jej komukolwiek kto chciałby się nauczyć statystyki. Jest to książka tak **beznadziejna**, że nie wymienię nawet jej tytułu – tłumaczenie z języka angielskiego zresztą.

W trakcie zajęć odkryłem także, że studenci w ogromnej większości są przekonani, że przedmiot ten jest trudny i mają w związku z tym ogromną – graniczącą z pewnością – obawę, że nie są w stanie go opanować.

W moim przekonaniu żeby ten stan rzeczy zmienić należy zademonstrować – teoretycznie oraz praktycznie – że statystyka nie jest wcale aż taka trudna. Teoretycznie oznaczałoby, że część matematyczna powinna być zredukowana do absolutnego minimum, a praktycznie, że należy pokazać jak i czym da się łatwo statystkę uprawiać czyli liczyć. Przez **czym** oczywiście współcześnie należy rozumieć konkretny program komputerowy i niekoniecznie tym programem musi być arkusz kalkulacyjny.

Sposobem na lepsze zrozumienie statystyki nie może być zasypanie studenta wielopiętrowymi formułami matematycznymi (co dyskwalifikuje część podręczników), ale raczej zastąpienie formalnej ścisłości przykładami i odwołanie się do zdrowego rozsądku, który przecież jest istotą wielu procedur statystycznych, takich jak na przykład testy istotności.

Kierując się tym założeniem całość teorii wyłożono w czterech pierwszych rozdziałach podręcznika. W pierwszym jak to zwykle bywa przedstawiono czym jest statystyka, czym się zajmuje oraz wyjaśniono **podstawowe pojęcia**, którymi statystyka się posługuje, takie jak populacja, próba czy pomiar. W rozdziale drugim przedstawiono metody **opisu statystycznego**. Rozdział trzeci przedstawia ideę **wnioskowania statystycznego**. Najbardziej obszerny rozdział czwarty opisuje metody **analizy zależności**.

Pozostałe rozdziały poświęcono praktyce. Rozdział piąty zawiera przykładowe analizy trzech badań ankietowych a rozdział szósty opisuje jak uprawiać statystykę z wykorzystaniem programu Jamovi.

Podręcznik został napisany dla studentów Wydziału Nauk o Zdrowiu Powiślańskiej Szkoły Wyższej, których uczę od kilku lat statystyki medycznej, ale mam nadzieję że może być przydatny także dla studentów innych kierunków zwłaszcza kierunków medycznych.

Elektroniczne wersje podręcznika w formatach HTML, PDF oraz EPUB znajdują się pod adresem: [hrpunio.github.io/SMI\\_Bookdown/wstęp.html](http://hrpunio.github.io/SMI_Bookdown/wstęp.html). Wykorzystane w omawianych w książce przykładach zbiory danych można znaleźć pod adresem: [github.com/hrpunio/SMBook](https://github.com/hrpunio/SMBook).

## Rozdział 1

# Przedmiot i metody badań statystycznych

## 1.1 Przedmiot statystyki

Wyraz statystyka ma wiele znaczeń: **statystyki zgonów** albo **statystyki alkoholizmu** czyli **dane** dotyczące zgonów lub alkoholizmu. Statystyka to też **diedzina wiedzy**, upraszczając zbiór metod, które służą do tworzenia statystyk w pierwszym znaczeniu tego słowa. Wreszcie statystyka to **pojedyncza metoda** ze zbioru metod opracowanych w dziedzinie, np. średnia to statystyka. Trochę to niefortunne, ale świat nie jest doskonały jak wiemy.

**Statystyka** (obiegowo): dział matematyki, a w związku z tym wiedza absolutnie pewna i obiektywna. Nieprawda choćby z tego powodu, że nie jest działem matematyki. Korzysta z metod matematycznych jak wiele innych dziedzin.

**Statystyka** od strony czysto praktycznej to: **dane + procedury** (zbierania, przechowywania, analizowania, prezentowania *danych*) + **programy**; Jeżeli statystyka kojarzy się komuś ze matematyką, wzorami i liczeniem, to jak widać jest to zaledwie podpunkt procedury → analizowanie.

## 1.2 Podstawowe pojęcia

Celem **badania statystycznego** jest uzyskanie informacji o interesującym zjawisku na podstawie danych. Zjawisko ma charakter masowy czyli dotyczy dużej liczby *obiektów*. Nie interesuje nas jeden zgon (obiekt) tylko zgony wielu ludzi.

**Populacja (zbiorowość statystyczna)** to zbiór obiektów będący przedmiotem badania statystycznego. Na przykład zgony w Polsce w roku 2022.

Każdy **obiekt** w populacji to **obserwacja** (zwana także **jednostką statystyczną** albo **pomiarem**) na jednej lub więcej **zmiennych** (albo **cech**). Jeżeli interesującym zjawiskiem są zgony, obserwacją jest osoba zmarła a zmiennymi (cechami) wiek, płeć, przyczyna zgonu oraz dzień tygodnia (w którym nastąpił zgon) zmarłej osoby.

**Próba** to część **populacji**. Na przykład część zgonów w Polsce w roku 2022.

**Parametr**: wielkość numeryczna obliczona na podstawie populacji.

**Statystyka**: wielkość numeryczna obliczona na podstawie próby.

Populacja powinna być zdefiniowana w taki sposób, aby nie było wątpliwości co tak naprawdę jest badane. Zgony to w oczywisty sposób za mało. *Zgony mieszkańców Kwidzyna w roku 2022*.

Zwróćmy uwagę, że *Zgony w mieście Kwidzyn w roku 2022* to nie to samo (ktoś może być mieszkańcem a umrzeć w Polsce i/lub ktoś może nie być mieszkańcem i umrzeć w Kwidzynie).

**Generalizacja**: ocena całości na podstawie części. Badamy zjawisko wypalenia zawodowego pielęgniarek i pielęgniarzy w Polsce (populacja). Wobec zaporowych kosztów mierzenia wszystkich decydujemy się na przeprowadzenie ankiety wśród studentów pielęgniarstwa PSW (próba). Czy możemy twierdzić na podstawie próby, że wyniki badania dla całej Polski są identyczne? Raczej nie.

Próba, która pozwala na generalizację nazywa się próbą **reprezentatywną**. Naj-

lepszym sposobem na uzyskanie próby reprezentatywnej jest losowanie.

W oczywisty sposób badanie na podstawie próby jest tańsze niż badanie całości, co nie oznacza że jest tanie. Kontynuując przykład: musielibyśmy mieć listę wszystkich pielęgniarek i pielęgniarzy w Polsce. Z tej listy wylosować próbę a następnie skontaktować się z wybranymi osobami (jak?). Dlatego też badania w oparciu o próbę nielosową są całkiem popularne (bo są tanie); należy jednakże mieć świadomość ich ograniczeń, w tym a zwłaszcza uogólnienia uzyskanych wyników.

## 1.3 Pomiar

Potocznie kojarzy się z linijką i wagą, ale w statystyce używany jest w szerszym znaczeniu. Ustalenie płci albo przyczyny zgonu to też pomiar.

**Pomiar** to przyporządkowanie wariantom **zmiennej** liczb lub symboli z pewnej **skali pomiarowej**. Przykładowo jeżeli jednostką statystyczną jest zgon a zmiennymi wiek, płeć, przyczyna zgonu oraz dzień tygodnia to pomiar będzie polegał na ustaleniu (przyporządkowaniu) wieku w latach, płci ('K'/'M'), przyczyny (identyfikatora z katalogu ICD10 zapewne) oraz numeru dnia tygodnia (lub nazwy dnia tygodnia). Wiek oraz numer dnia są liczbami, płeć i przyczyna symbolem.

Wyróżnia się następujące **typy skal pomiarowych**:

- **nominalna** (*nominal scale*), klasyfikuje: płeć zmarłego;
- **porządkowa** (*ordinal scale*), klasyfikuje i porządkuje: dzień tygodnia w którym nastąpił zgon (po poniedziałku jest wtorek);
- **liczbowa**, mierzy w potocznym tego słowa znaczeniu: wiek zmarłego w latach.

Mówimy **zmienna mierzalna** albo **zmienna ilościowa** dla zmiennych mierzonych za pomocą skali liczbowej. Mówimy **zmienna niemierzalna** albo **zmienna jakościowa** dla zmiennych mierzonych za pomocą skali nominalnej lub porządkowej.

Zmienne mierzalne dzielą się na **skokowe** oraz **ciągłe**. Skokowe są to zmienne, które przyjmują skończoną (albo przeliczalną) liczbę wartości. Matematycznym odpowiednikiem zmiennej skokowej jest zbiór liczb całkowitych. Przykładem zmiennej skokowej jest liczba dzieci w gospodarstwie domowym albo liczba dni do pierwszego nawrotu choroby po zakończeniu leczenia. Zmienna ciągła to taka zmienna, która może przyjąć nieskończoną/nieprzeliczalną liczbę wartości. Matematycznym odpowiednikiem zmiennej ciągłej jest zbiór liczb rzeczywistych. Przykładem może być ciśnienie krwi lub waga noworodka.

**Rodzaje danych**

- Przekrojowe (zmarli w Kwidzynie);
- Czasowe: każda obserwacja ma przypisany czas (liczba zmarłych w Polsce w latach 2000–2022);
- Przestrzenne: każda obserwacja ma przypisane miejsce na kuli ziemskiej (współrzędne geograficzne).

## 1.4 Rodzaje i sposoby analizy danych

Rodzaje **analizy statystycznej** zależą od rodzaju danych (jakie mamy dane takie możemy stosować metody):

- jedna zmienna/dane przekrojowe: analiza struktury;

- jedna zmienna/dane czasowe: analiza dynamiki zjawiska;
- co najmniej dwie zmienne: analiza współzależności (nadwaga powoduje cukrzycę).

Sposoby analizy danych zależą od sposobu pomiaru (populacja/próba/generalizacja):

**Opis statystyczny** – (proste) przedstawienie badanych zbiorowości/zmiennych tabel, wykresów lub parametrów (np. średnia, mediana) ; Opis statystyczny może dotyczyć: – struktury zbiorowości; – współzależności; – zmian zjawiska w czasie.

**Wnioskowanie statystyczne:** wnioskowanie na temat całości na podstawie próby; wykorzystuje metody analizy matematycznej.

**Opisujemy** populację lub próbę. **Wnioskujemy** na podstawie próby o całości.

## 1.5 Sposoby pomiaru danych i organizacja badania

Sposób pomiaru/organizacja badania ma zasadnicze znaczenie dla interpretacji wyników. Są dwa fundamentalne rodzaje pomiaru (sposobu zebrania danych) **eksperyment** oraz **obserwacja**.

Mówimy w związku z tym **dane eksperymentalne** albo **dane obserwacyjne**.

### Picie kawy a lepsza ocena

Chcemy ustalić czy spożywanie kawy w czasie sesji egzaminacyjnej skutkuje uzyskaniem lepszej oceny. W celu oceny prawdziwości takiej tezy przeprowadzono badanie wśród studentów pytając ich o to ile kawy pili w czasie sesji i zestawiając te dane z wynikami egzaminów. Średnie wyniki w grupie studentów pijących dużo kawy były wyższe w grupie pijącej mało kawy. Czy można powiedzieć, że udowodniono iż picie dużej ilości kawy poprawia wynik egzaminu?

Raczej nie: można sobie wyobrazić, że studenci którzy poświęcili więcej czasu na naukę pili w tym czasie kawę (na przykład żeby nie zasnąć). Prawdziwą przyczyną jest czas poświęcony na przygotowanie a nie to ile ktoś wypił lub nie wypił kawy. Inaczej mówiąc gdyby ktoś pił dużo kawy, bo uwierzył, że to poprawi mu wyniki i się nie uczył, to pewnie by się rozczarował.

Rodzaje badań: **eksperymentalne** vs **obserwacyjne**.

**Eksperyment kontrolowany** (zrandomizowany lub nie): służy do weryfikacji związku **przyczyna-skutek**. Skutek może być rezultatem działania wielu **czynników** (zmiennych). Eksperymentator manipuluje wielkością przyczyn (zmiennych **niezależnych**) oraz mierzy wielkość skutku (zmiennej **zależnej**); Wszystkie pozostałe czynniki (zmienne **ukryte**) są **kontrolowane** (w tym sensie, że ich wpływ na skutek jest ustalony).

Pomiarowi/manipulacji podlega zbiór jednostek podzielonych **losowo** na dwie grupy: grupa **eksperymentalna** (**experimental group**) oraz **grupa kontrolna** (**control group**).

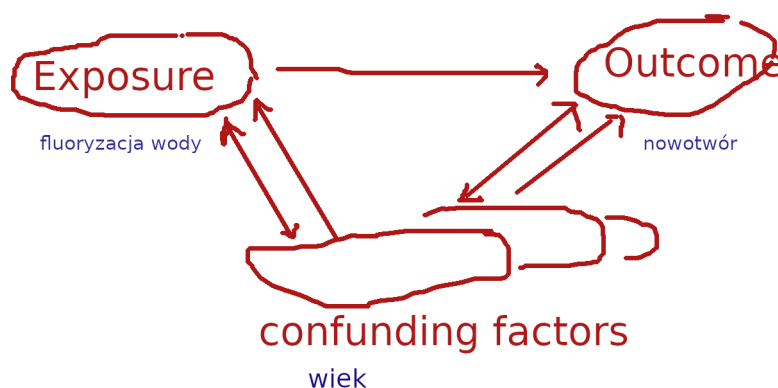
W medycynie używa się terminu **badania kliniczne** czyli badania które dotyczą ludzi. Badania kliniczne także dzielą się na eksperymentalne oraz obserwacyjne. Eksperyment nazywa się RCT (*randomized clinical trial*/randomizowane kontrolowane badania kliniczne.) Manipulacja określana jest jako ekspozycja (**exposure**) albo lecze-



nie (**treatment**) Zmienne ukryte określa się mianem **confunding factors** (czynniki zakłócające).

Rysunek 1.1 przedstawia zależność pomiędzy wynikiem (*outcome*), przyczyną oraz czynnikami zakłócającymi na przykładzie zależności dotyczącej domniemanego wpływu fluoryzowania wody na zwiększenie ryzyka zgonów z powodu nowotworów. W badaniu, którego autor uważał, że udowodnił związek fluoryzowanie→nowotwór porównał on współczynniki zgonów z miast fluoryzujących oraz nie fluoryzujących wodę. Okazało się, że przeciętnie współczynnik ten był wyższy w grupie miast fluoryzujących wodę. Czy to świadczy, że fluoryzowanie wody powoduje raka? Nie.

W innym badaniu tych samych miast okazało się, że w grupie miast fluoryzujących wodę przeciętnie mieszkają starsi ludzie. A ponieważ współczynniki zgonów rosną wraz z wiekiem, to nie można wykluczyć, że prawdziwą przyczyną obserwowanego zwiększenia wartości współczynników zgonów jest wiek a nie fluoryzacja.



Rysunek 1.1: Fluoryzacja, wiek a nowotwór

**Efekt przyczynowy** to ilościowe określenie wpływu ekspozycji na wynik poprzez porównanie wielkości wyniku dla różnych wielkości ekspozycji.

Są dwa typy **efektu przyczynowego**: indywidualny efekt interwencji (*individual treatment effect*) oraz średni efekt interwencji (*average treatment effect*).

#### Individual Treatment Effect (ITE)

Indywidualny efekt interwencji (ITE) określa ilościowo wpływ interwencji dla konkretnej osoby poprzez porównanie wyników dla różnych wartości interwencji.

Mogę pić kawę lub nie pić kawy a wynikiem będzie ocena. Oczywiście nie mogą zrobić tych dwóch rzeczy na raz.

#### Average Treatment Effect (ATE)

Średni efekt interwencji określa ilościowo wpływ interwencji dla grupy osób.

W grupie studentów jedni pili kawę inni nie.

Jeżeli grupa (populacja) została uprzednio podzielona (losowo) na grupę **eksperymentalną** oraz **grupę kontrolną** możemy policzyć ATE oddzielnie dla obu grup. Wtedy efekt przyczynowy można zdefiniować jako:

$ATT - ATC$  (albo  $ATT/ATC$ )

gdzie: ATT oznacza ATE w grupie eksperymentalnej a ATC oznacza ATE w grupie kontrolnej.

**Kawa a ocena (kontynuacja)**

Można przypuszczać, że oprócz kawy na wynik egzaminu ma wpływ np. wrodzone predyspozycje w dziedzinie intelektualnej oraz czas poświęcony na naukę. Aby kontrolować ten czynnik można podzielić losowo grupę studentów; dzięki czemu średnia wielkość predyspozycji oraz czasu w obu grupach będzie podobna. Następnie zalecamy studentom w **grupie eksperymentalnej** picie 1l kawy dziennie a studentom w **grupie kontrolnej** podajemy 1l brązowej wody o smaku i zapachu kawy :-). Średnie wyniki w grupie studentów pijących 1l kawy okazały się wyższe niż w grupie pijącej kolorową wodę. Czy można powiedzieć że udowodniono iż picie dużej ilości kawy poprawia wynik egzaminu? Raczej tak.

**Badania obserwacyjne** można z kolei podzielić na **analityczne** i **opisowe**.

W badaniach **analitycznych** porównuje się grupę kontrolną z grupą poddaną ekspozycji/leczeniu; w badaniach przekrojowych nie ma grupy kontrolnej.

Badania analityczne dzielimy dalej na:

- **kohortowe**,
- **kliniczno-kontrolne**,
- **przekrojowe**.

Badanie **kohortowe** (*cohort study*): wieloletnie badania na dużej grupie jednostek. Pomiar zaczynamy od ekspozycji kończymy na wyniku/chorobie/zgonie (takie badanie nazywamy **prospektywnym**. Problem: koszty (np. choroby rzadkie wymagają ogromnych kohort).

Badanie **kliniczno-kontrolne** (*case-control study*): **retrospektywna** ocena ekspozycji dla jednostek, u których stwierdzono wynik (chorobę). Grupę kontrolną tworzą **dopasowane** jednostki u których wyniku nie stwierdzono (dopasowane w sensie, że są podobne podobne.) W praktyce badanie kliniczno-kontrolne to badanie chorych, którzy zgłosili się do przychodni; grupą kontrolną są podobni chorzy (wiek, płeć) z innej przychodni.

Problem1: błąd pamięci (*recall bias*) pacjenci – zwłaszcza zdrowi – słabo pamiętają fakty które miały miejsce lata temu. Problem2: trudności z **dopasowaniem** grupy kontrolnej (łatwiej powiedzieć niż zrobić).

Badania **prospektywne**: od przyczyny do skutku (*cohort*); badanie **retrospektywne**: od skutku do przyczyny (*case-control*).

Badanie przekrojowe (*cross-sectional study*): badanie związku między wynikiem a ekspozycją bez podziału na grupę eksperymentalną i kontrolną.

Problem: nie da się określić związku przyczyna-skutek w taki sposób jaki się stosuje w badaniach analitycznych, ale można do tego celu zastosować **model** regresji liniowej.

**Palenie a nowotwór**

Badamy grupę pacjentów przychodni onkologicznej. Stwierdzamy że 90% z nich paliło papierosy. Czy z tego wynika że palenie powoduje raka? Niekoniecznie. Możemy **dopasować** pacjentów o podobnym profilu demograficzno-społecznym z innej przychodni (którzy nie chorują na raka) i stwierdzić że 20% z nich paliło. To już jest konkretny argument – ale takie badanie nie jest już **przekrojowe** tylko **kliniczno-kontrolne**.

**Kawa a ocena (kontynuacja)**

Można oprócz pytania studentów o ilość kawy i wynik pytać ich jeszcze o czas poświęcony na naukę oraz o średnią ze studiów (wrodzone predyspozycje w dziedzinie intelektualnej). Za pomocą metody regresji możemy ustalić czy i jak bardzo kawa, czas i predyspozycje wpływają na ocenę. Teoretycznie zamiast eksperymentu można używać regresji, ale jest to w większości przypadków trudne – albo zmienne nie da się zmierzyć (czy średnia ze studiów jest dobrą miarą predyspozycji?) albo jakąś ważną zmienną pominiemy. Więcej na temat regresji w rozdziale 4.

Badanie **ekologiczne**: badanie (przekrojowe) zależności pomiędzy wielkościami **zagregowanymi** a **indywidualnymi**. Przykładowo wielkości zagregowane, to zależność pomiędzy przeciętną wielkością dochodu narodowego, a przeciętną oczekiwaną długością życia np. na poziomie kraju.

Problem: błąd ekologiczmu (**ecological fallacy**.) Zależności na poziomie indywidualnym oraz zagregowanym mogą być różne. Można oczekiwać, że im większy dochód tym osoba dłużej żyje (poziom indywidualny.) Jeżeli w kraju występują duże różnice w dochodach (na przykład w USA), to przeciętnie dochód jest wysoki, ale jest dużo osób o niskich dochodach, o ograniczonym dostępie do służby zdrowia, i krótszej oczekiwanej długości życia. Przeciętna oczekiwana długość życia na poziomie całego kraju jest niższa (bo jest sumą wysokiej dla bogatych + niskiej dla biednych). W rezultacie zależność na poziomie zagregowanym może się znacząco różnić od tej na poziomie indywidualnym.

**1.5.1 Przykłady badań**

Jest ustalony szablon artykułu naukowego, który powinien być podzielony na następujące części:

1. **Wprowadzenie**: określenie problemu badawczego, celu badania.
2. **Materiał i metoda**: Opis danych i zastosowanych metod statystycznych.
3. **Wyniki**: Rezultaty analiz.
4. **Dyskusja**: Znaczenie uzyskanych wyników; jeżeli we wstępie postawiono hipotezy to tutaj należy przedstawić wyniki ich weryfikacji.

Żeby się zorientować jakie dane (jakie zmienne i jak mierzone) oraz jakie metody statystyczne zostały wykorzystane w pracy wystarczy zapoznać się z treścią punktu **materiał i metoda**. W szczególności powinien tam być określony **rodzaj badania**: eksperyment, badanie kohortowe, kliniczno-kontrolne, przekrojowe lub inne.

**Czy konsumpcja soli kuchennej szkodzi? (eksperyment)**

Neal B. i inni zastosowali eksperyment kontrolowany do zbadania wpływu substytucji chlorku sodu chlorkiem potasu na choroby sercowo-naczyniowe (*Effect of Salt Substitution on Cardiovascular Events and Death, New England Journal of Medicine*, <https://doi.org/10.1056/NEJMoa2105675>). W badaniu przeprowadzonym w Chinach, uczestniczyli mieszkańcy 600 wsi, podzieleni losowo na dwie grupy. Uczestnik badania musiał mieć minimum 60 lat oraz nadciśnienie krwi. W badaniu uczestniczyło prawie 21 tysięcy osób. Przez pięć lat trwania eksperymentu

grupa kontrolna używała soli zawierającej 75% chlorku potasu oraz 25% chlorku sodu; grupa badana zaś używała soli tradycyjnej czyli zawierającej wyłącznie chlorek sodu. Obserwowano w okresie pięcioletnim w obu grupach liczbę udarów, incydentów sercowo-naczyniowych oraz zgonów. Wpływ substytucji oceniono porównując współczynniki ryzyka w obu grupach.

#### Konflikt praca-dom w zawodzie pielęgniarstwa (przekrojowe)

Simon i inni badali konflikt Praca-Dom w zawodzie Pielęgniarki/Pielęgniarsza (Work-Home Conflict in the European Nursing Profession Michael Simon 1, Angelika Kümmerling, Hans-Martin Hasselhorn; Next-Study Group Int J Occup Environ Health 2004 Oct-Dec;10(4):384-91. doi: 10.1179/oeh.2004.10.4.384. <https://pubmed.ncbi.nlm.nih.gov/15702752/>).

Konflikt Praca-Dom (WHC) to sytuacja kiedy nie można zająć się zadaniami lub obowiązkami w jednej dziedzinie ze względu na obowiązki w drugiej domenie. Teoria zapożyczona z obszaru Nauk o Zarządzaniu zapewne. Ten konflikt jest mierzony odpowiednią skalą pomiarową składającą się z pięciu pytań. Czynniki które WHC mają powodować są: czas pracy, grafik (w sensie rodzaj etatu/zmianowość), nacisk-na-nadgodziny (występuje lub nie), intensywność pracy, obciążenie emocjonalne oraz jakość zarządzania. (ostatnie trzy mierzone odpowiednimi skalami pomiarowymi, czytaj: serią pytań w ankiecie). Badano 27,603 osoby. Podstawowym narzędziem badawczym jak się łatwo domyśleć była ankieta, a przyczyny WHC ustalono za pomocą metody regresji wielorakiej.

Teraz porównajmy koszty badania #1, w którym jedynie starano się ustalić że sól szkodzi (lub nie) z badaniem #2, w którym starano się ustalić przyczyny stanów psychicznych badanych.

## 1.6 Miary częstości chorób

**Populacja narażona** (*population at risk*): grupa osób podatnych na zdarzenie (chorobę); rak szyjki macicy dotyczy kobiet a nie wszystkich.

**Współczynnik chorobowości** (*prevalence rate*): liczba chorych w określonym czasie (dzień, tydzień, rok) podzielona przez wielkość populacji narażonej. Ponieważ są to zwykle bardzo małe liczby, mnoży się wynik przez  $10^n$  dla ułatwienia interpretacji. Czyli jeżeli chorych w populacji narażonej o wielkości 1mln jest 20 osób, to współczynnik wynosi  $20/1\text{mln} = 0,000002$  co trudno skomentować po polsku. Jeżeli pomnożymy owe 0,000002 przez 100 tys ( $n = 5$ ), to współczynnik będzie równy 2, co interpretujemy jako dwa przypadki na 100 tys. (albo 0,2 na 10 tys, jeżeli  $n = 4$ , co już jednak brzmi trochę gorzej).

**Współczynnik zapadalności** (*incidence rate*): liczba nowych chorych w określonym czasie (dzień, tydzień, rok) podzielona przez wielkość populacji narażonej. Też zwykle pomnożona przez  $10^n$ .

**Współczynnik śmiertelności** (*case fatality rate*): liczba zgonów z powodu X w określonym czasie (dzień, tydzień, rok) podzielona przez liczbę chorych na X w tym samym czasie. Śmiertelność jest miarą ciężkości choroby X.

**Współczynnik zgonów** (*death rate*): liczba zgonów w określonym czasie przez średnią liczbę ludności w tym czasie (pomnożone przez  $10^n$ ).

Jeżeli współczynnik zgonów nie uwzględnia wieku, nazywany jest surowym (*crude*); grupy różniące się strukturą wieku nie powinny być porównywane za pomocą współczynników surowych tylko standaryzowanych (*age-standardized* albo *age-adjusted*). Przykładowo jeżeli porównamy współczynnik zgonów USA i Nigerii to okaże się że w USA jest wyższy a to z tego powodu że społeczeństwo amerykańskie jest znacznie starsze (a umierają zwykle ludzie starzy).

**Współczynnik zgonów** standaryzowany według wieku to ważona średnia współczynników w poszczególnych grupach wiekowych, gdzie wagami są udziały tychże grup wiekowych w pewnej **standardowej populacji**.

## 1.7 Oprogramowanie

Nie da się praktykować statystyki bez korzystania z programów komputerowych i mamy w tym zakresie trzy możliwości:

1. Arkusz kalkulacyjny. Przydatny na etapie zbierania danych i ich wstępnej analizy, później już niekoniecznie. Policzenie niektórych rzeczy jest niemożliwe (brak stosownych procedur) lub czasochłonne (w porównaniu do 2–3).
2. Oprogramowanie specjalistyczne komercyjne takie jak programy STATA czy SPSS. Wady: cena i czas niezbędny na ich poznanie.
3. Oprogramowanie specjalistyczne darmowe: Jamovi oraz R. Same zalety.

W większości podręczników opisuje się **procedury** oraz **program**, w którym te procedury można zastosować jednocześnie. My zdecydowaliśmy się oddzielnie przedstawić teorię statystyki (rozdziały 1–5) a oddzielnie opis posługiwania się konkretnym programem (rozdział 6).

## Rozdział 2

# Analiza jednej zmiennej

**Statystyka opisowa** (opis statystyczny) to zbiór metod statystycznych służących do – surprise, surprise – opisu (w sensie przedstawienia sumarycznego) zbioru danych; w zależności od typu danych (przekrojowe, czasowe, przestrzenne) oraz sposobu pomiaru (dane nominalne, porządkowe liczbowe) należy używać różnych metod.

W przypadku **danych przekrojowych** opis statystyczny nazywany jest **analizą struktury** i sprowadza się do opisania danych z wykorzystaniem:

- tablic (statystycznych);
- wykresów;
- parametrów (takich jak średnia czy mediana).

**Rozkład zmiennej** (cechy) to przyporządkowanie wartościom zmiennej odpowiedniej **liczby wystąpień** w postaci **liczebności** albo **częstości** (popularnych procentów).

**Analiza struktury** (dla jednej zmiennej) obejmuje:

- **określenie tendencji centralnej** (**miary położenia**: wartość przeciętna, mediana, dominanta);
- **zróżnicowanie wartości** (rozproszenie: odchylenie standardowe, rozstęp ćwiartkowy);
- **asymetrię** (rozłożenie wartości zmiennej wokół średniej).

## 2.1 Tablice statystyczne

**Tablica statystyczna** to (w podstawowej formie) dwukolumnowa tabela zawierająca wartości zmiennej oraz odpowiadające tym wartościom liczebności (i/lub częstości).

**Tablica dla zmiennej niemierzalnej (nominalnej albo porządkowej).**

### Absolwenci studiów pielęgniarских w ośmiu największych krajach UE

Tablica: Absolwenci studiów pielęgniarских w ośmiu największych krajach UE w roku 2018

kraj	liczba absolwentów
Belgium	7203
Germany	35742
Spain	9936
France	25757
Italy	11207
Netherlands	9920
Poland	9070
Romania	18664

Źródło: Eurostat, tablica Health graduates (HLTH\_RS\_GRD)

W przykładzie **jednostką statystyczną** jest absolwent studiów pielęgniarских w

roku 2018, badaną **zmienną** zaś **kraj w którym ukończył studia**.

#### Tablica dla zmiennej mierzalnej liczbowej skokowej

Przypomnijmy, że zmienna skokowa to taka zmienna, która może przyjąć skończoną (przeliczalną) liczbę wartości.

Jeżeli tych wartości jest mało, to tablica zawiera wyliczenie wartości zmiennej i odpowiadających im liczebności. Jeżeli liczba wariantów zmiennej jest duża, to tablica zawiera klasy wartości (przedziały wartości) oraz odpowiadające im liczebności.

Liczba przedziałów jest dobierana metodą prób i błędów, tak aby:

- przedziały wartości były jednakowej rozpiętości;
- Na zasadzie wyjątku dopuszcza się aby pierwszy i ostatni przedział były **otwarte**, tj. nie miały dolnej (pierwszy) lub górnej (ostatni) **granicy**;
- nie było przedziałów z zerową liczebnością;
- przedziałów nie było za dużo ani za mało (typowo 5–15);
- większość populacji nie znajdowała się w jednym albo dwóch przedziałach.

#### Gospodarstwa domowe wg liczby osób

Tablica: Gospodarstwa domowe w mieście Kwidzyn wg liczby osób w roku 2021

liczba osób	liczba gospodarstw	%
1	2790	21,64972
2	3420	26,53837
3	2618	20,31505
4	2246	17,42842
5 i więcej	1813	14,06844
razem	12887	100,00000

Źródło: Bank danych lokalnych GUS, podgrupa P4287/Gospodarstwa domowe według liczby osób

W powyższym przykładzie druga kolumna tablicy zawiera liczebności a trzecia częstości (udziały procentowe). W 1813 gospodarstwach domowych mieszkało 5 i więcej osób co stanowiło 14,1% wszystkich gospodarstw domowych w mieście Kwidzyn.

#### Tablica dla zmiennej mierzalnej liczbowej ciągłej

Przypomnijmy, że zmienna ciągła to taka zmienna, która może przyjąć nieskończoną i nieprzeliczalną liczbę wartości.

Tablica zawiera klasy (przedziały) wartości oraz odpowiadające im liczebności.

Liczba przedziałów jest dobierana metodą prób i błędów, tak aby:

- przedziały wartości były jednakowej rozpiętości;
- Na zasadzie wyjątku dopuszcza się aby pierwszy i ostatni przedział były **otwarte**, tj. nie miały dolnej (pierwszy) lub górnej (ostatni) **granicy**;
- nie było przedziałów z zerową liczebnością;
- przedziałów nie było za dużo ani za mało (typowo 8–15);
- większość populacji nie znajdowała się w jednej czy dwóch przedziałach;
- zwykle przyjmuje się za końce przedziałów **okrągłe liczby** bo dziwnie by wyglądało gdyby koniec przedziału np. był równy 1,015 zamiast 1,0.

### Dzietność kobiet na świecie

Współczynnik dzietności (*fertility ratio* albo FR) – przeciętna liczba urodzonych dzieci przypadająca na jedną kobietę w wieku rozrodczym (15–49 lat). Przyjmuje się, iż FR między 2,10–2,15 zapewnia zastępowalność pokoleń.

Dane dotyczące dzietności dla wszystkich krajów świata pobrano ze strony <https://ourworldindata.org/grapher/fertility-rate-complete-gapminder>.

Zbudujemy tablicę przedstawiającą rozkład współczynników dzietności w roku 2018. Wszystkich krajów jest 201. Wartość minimalna współczynnika wynosi 1,22, a wartość maksymalna to 7,13. Decydujemy się na rozpiętość przedziału równą 0,5; dolny koniec pierwszego przedziału przyjmujemy jako 1,0.

Tablica: Kraje świata według współczynnika dzietności (2018)

Wsp. dzietności	liczba krajów
(1,1,5]	24
(1,5,2]	61
(2,2,5]	40
(2,5,3]	17
(3,3,5]	8
(3,5,4]	15
(4,4,5]	11
(4,5,5]	12
(5,5,5]	6
(5,5,6]	5
(6,6,5]	1
(7,7,5]	1

Źródło: <https://ourworldindata.org/grapher/fertility-rate-complete-gapminder>

Zapis  $(1, 1.5]$  oznacza przedział od 1,0 do 1,5 przy czym dolny koniec nie należy do przedziału a górny 1,5 należy. Który koniec „wchodzi”, a który nie powinien być jasno oznaczony. Zwykle jest tak jak w przykładzie: górny „wchodzi”, dolny „nie wchodzi”.

Każda tablica statystyczna **musi** mieć:

1. Część liczbowa (kolumny i wiersze);
  - żadna rubryka w części liczbowej nie może być pusta (żelazna zasada); w szczególności brak danych należy zaznaczyć umownym symbolem.
2. Część opisową:
  - tytuł tablicy;
  - nazwy (opisy zawartości) wierszy;
  - nazwy (opisy zawartości) kolumn;
  - wskazanie źródła danych;
  - ewentualne uwagi odnoszące się do danych liczb.

Pominięcie czegokolwiek z powyższego jest **ciężkim błędem**. Jeżeli nie ma danych (a często – z różnych powodów – nie ma, należy to zaznaczyć a nie pozostawiać pustą rubrykę).



## 2.2 Wykresy

**Wykresy statystyczne** są graficzną formą prezentacji materiału statystycznego, są mniej precyzyjne i szczegółowe niż tablice, natomiast bardziej sugestywne.

Celem jest pokazanie rozkładu wartości zmiennej w populacji: jakie wartości występują często a jakie rzadko, jak bardzo wartości różnią się między sobą. Jak różnią się rozkłady dla różnych, ale logicznie powiązanych populacji (np rozkład czegoś-tam w kraju A i B albo w roku X, Y i Z).

Do powyższego celu stosuje się:

- **wykres słupkowy** (skala nominalna/porządkowa);
- **wykres kołowy** (skala nominalna/porządkowa);
- **histogram** (albo wykres słupkowy dla skal nominalnych).

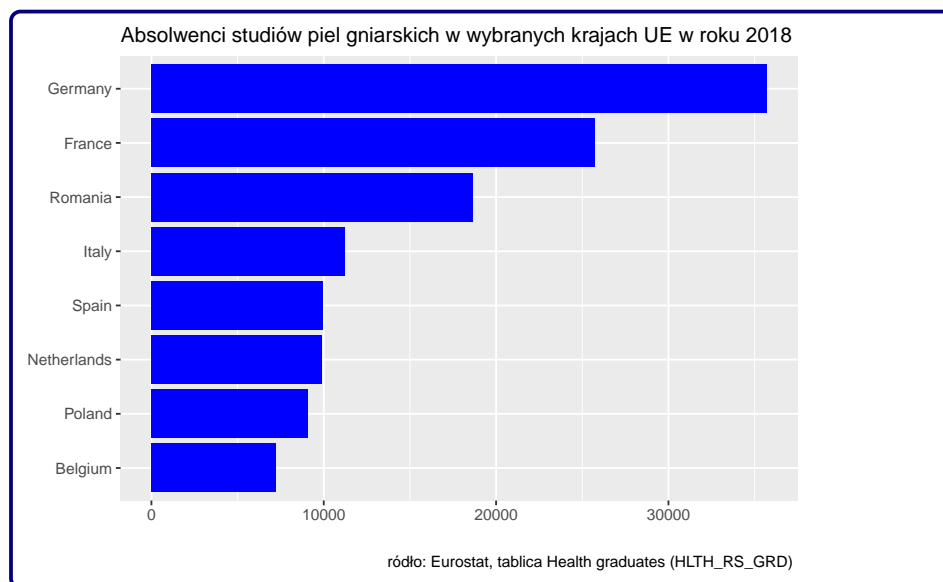
Uwaga: **wykres kołowy** jest zdecydowanie gorszy od wykresu słupkowego i nie jest zalecany. **Każdy** wykres kołowy można wykreślić jako słupkowy i w takiej postaci będzie on bardziej zrozumiały i łatwiejszy w interpretacji.

Podobnie jak tablice, rysunki powinny być opatrzone tytułem oraz zawierać źródło wskazujące na pochodzenie danych (zobacz przedstawione przykłady).

### 2.2.1 Skala nominalna i porządkowa

#### Wykres słupkowy (*bar chart*)

Na wykresie słupkowym długość każdego prostokąta (**słupka**) jest proporcjonalna do liczebności, którą reprezentuje. Wartości zmiennej (etykiety) są umieszczane pod lub obok słupka. Słupki można rysować pionowo lub poziomo (jako na powyższym przykładzie). Jeżeli etykiety są długie, to należy słupki rysować poziomo, bo wtedy można zmieścić etykiety bez potrzeby ich obracania o 90° czy skracania.

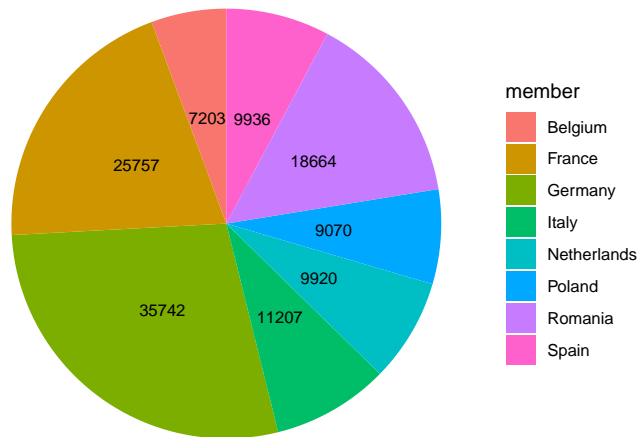


#### Wykres kołowy (*pie chart*)

Na wykresie kołowym długość łuku każdego wycinka (a także kąt środkowy oraz

pole wycinka) jest proporcjonalna do liczebności, którą reprezentuje. Łącznie wszystkie wycinki tworzą pełne koło.

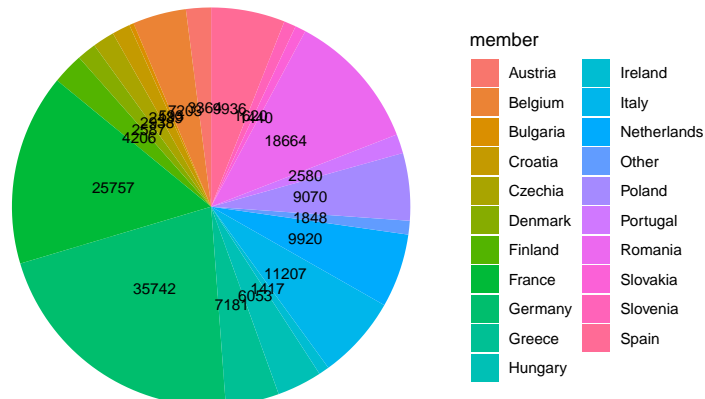
Absolwenci studiów piel gniarskich w wybranych krajach UE w roku 2018  
ródło: Eurostat, tablica Health graduates (HLTH\_RS\_GRD)



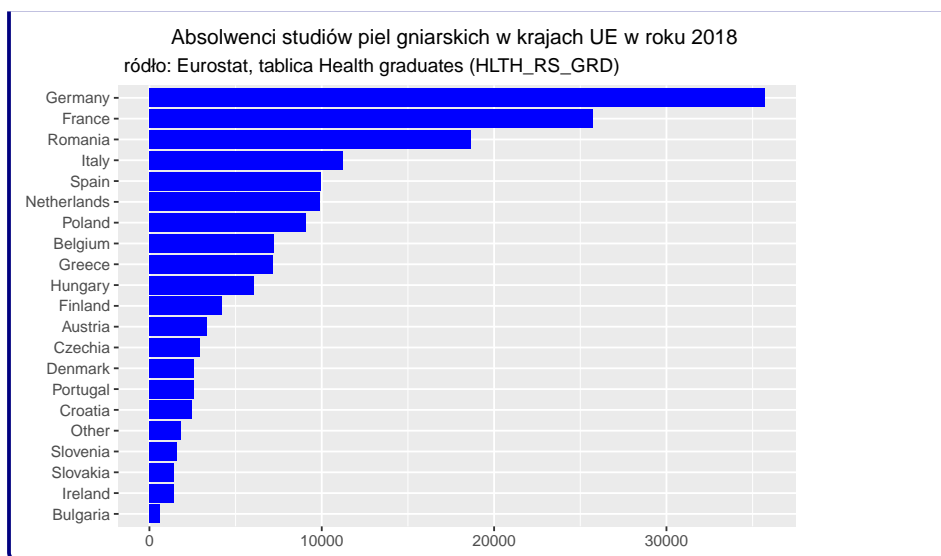
Wykres słupkowy i kołowy przedstawiają dokładnie to samo, zatem który wybrać?

Wykres kołowy wygląda zapewne efektowniej (z uwagi na paletę kolorów), ale jest mniej efektywny. Wymaga dodania legendy, która utrudnia interpretację treści. Jeżeli zwiększymy liczbę krajów, to wykres kołowy staje się zupełnie nieczytelny, bo brakuje rozróżnialnych kolorów, a wycinki koła są zbyt wąskie żeby cokolwiek wyróżniały:

Absolwenci studiów piel gniarskich w krajach UE w roku 2018  
ródło: Eurostat, tablica Health graduates (HLTH\_RS\_GRD)

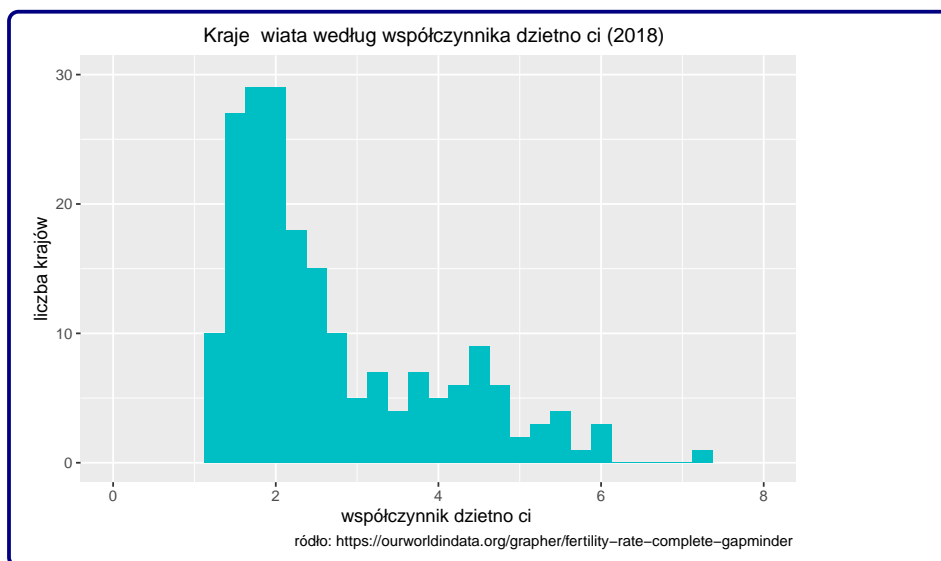


Wykres słupkowy dalej jest natomiast OK:



### 2.2.2 Skala liczbowa

Histogram to coś w rodzaju wykresu słupkowego tylko na osi OX zamiast wariantów zmiennej są przedziały wartości.

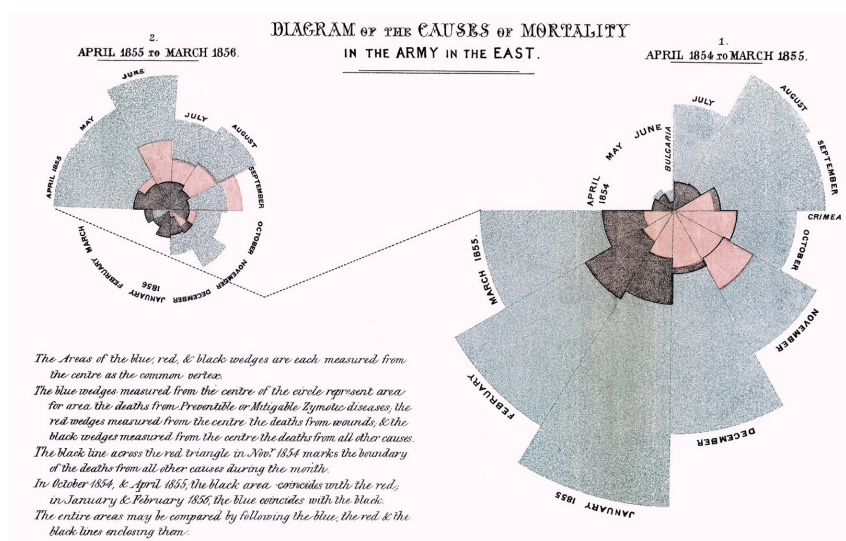


Im więcej przedziałów (mniejsza rozpiętość przedziału), tym histogram jest bardziej szczegółowy co niekoniecznie jest pożądane bo zaciemnia ogólny obraz. Nie ma złotych recept na to ile powinno być przedziałów, a ich liczba determinuje kształt oraz optyczną wielkość histogramu. Im mniej przedziałów tym histogram będzie optycznie większy.

## 2.3 Statystyczka Florence Nightingale

Nie każdy, kto wie kim była Florence Nightingale, wie że była ona także statystykiem. W czasie wojny krymskiej nie tylko zorganizowała opiekę nad rannymi żołnierzami, ale również – aby przekonać swoich przełożonych do zwiększenia nakładów na szpitale polowe – prowadziła staranną ewidencję szpitalną oraz zgromadzone dane potrafiła analizować, używając wykresów własnego projektu.

W szczególności słynny jest diagram Nightingale zwane także różą Nightingale (rys. 2.1), które wprawdzie (podobno) nie okazały się szczególnie użyteczny, no ale nie każdy nowy pomysł jest od razu genialny:



Rysunek 2.1: Róża Nightingale

Jest to coś w rodzaju wykresu słupkowego tyle że zamiast słupków są wycinki koła. Wycinków jest dwanaście tyle ile miesięcy. Długość promienia a co za tym idzie wielkość pola wycinka zależy od wielkości zjawiska, który reprezentuje (przyczyna śmierci: rany/choroby/inne).

Wpisując Florence+Nightingale można znaleźć dużo informacji na temat, w tym: <http://www.matematyka.wroc.pl/ciekawieomatematyce/pielgniarka-statystyczna>.

W 1859 roku Nightingale została wybrana jako pierwsza kobieta na członka Royal Statistical Society (Królewskie Stowarzyszenie Statystyczne) oraz została honorowym członkiem American Statistical Association (Amerykańskiego Stowarzyszenia Statystycznego).

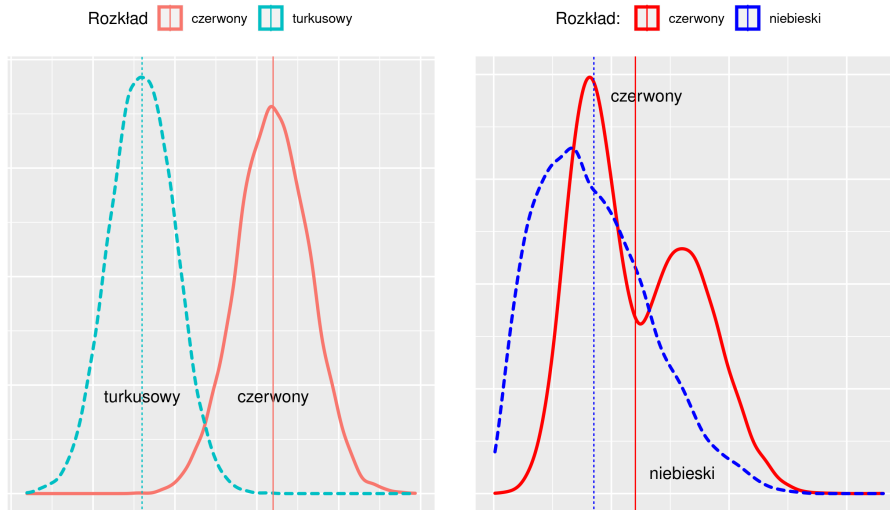
Więc szanowi czytelnicy wnioski są oczywiste.

## 2.4 Analiza parametryczna

Analiza parametryczna z oczywistych względów dotyczy tylko zmiennych mierzonych na skali liczbowej.

### 2.4.1 Miary położenia

Miary przeciętne (**położenia**) charakteryzują średni lub typowy poziom wartości zmiennej. Są to więc takie wartości, wokół których skupiają się wszystkie pozostałe wartości analizowanej zmiennej.



Rysunek 2.2: Rozkłady zmiennej a miary średnie

Na rysunku 2.2 po lewej mamy dwa rozkłady różniące się poziomem przeciętnym. Rozkład czerwony ma przeciętnie większe wartości niż turkusowy. Są to rozkłady **jednomodalne**, czyli takie, w których rozkład zmiennej skupia się wokół jednej wartości. Dla takich rozkładów ma sens obliczanie średniej arytmetycznej. Te średnie wartości są zaznaczone na rysunku linią pionową.

Na rysunku po prawej mamy rozkłady **nietypowe: wielomodalne** (czerwony) lub **niesymetryczne** (niebieski). W rozkładzie niesymetrycznym wartości skupiają się nie centralnie, ale po prawej/lewej od środka przedziału zmienności/wartości średniej).

W świecie rzeczywistym zdecydowana większość rozkładów jest jednomodalna. Rzadkie przypadki rozkładów wielomodalnych zwykle wynikają z łącznego analizowania dwóch różniących się wartości średnią zbiorów danych. Oczywistym zaleceniem w takiej sytuacji jest analiza każdego zbioru oddzielnie.

Rodzaje miar położenia:

- klasyczne:
  - **średnia arytmetyczna**,
- pozycyjne:
  - **mediana**,
  - **dominanta**,
  - **kwartyle**,
  - ewentualnie kwantyle, decyle, centyle (rzadziej używane).

**Średnia arytmetyczna** (*mean, arithmetic mean*), to łączna suma wartości podzielona przez liczbę sumowanych jednostek. Jeżeli wartość  $i$ -tej jednostki w zbiorowości o liczebności  $N$  oznaczmy jako  $x_i$  (gdzie:  $i = 1, \dots, N$ ), to średnią można zapisać jako:

$$\bar{x} = (x_1 + x_2 + \dots + x_N)/N$$

Uwaga: we wzorach statystycznych zmienne zwykle oznacza się małymi literami a średnią dla zmiennej przez umieszczenie nad nią kreski poziomej czyli  $\bar{x}$  to średnia wartość zmiennej  $x$ .

**Mediana** (*median*, kwartył drugi) dzieli **uporządkowaną** zbiorowość na dwie równe części; połowa jednostek ma wartości zmiennej mniejsze lub równe medianie, a połowa wartości zmiennej równe lub większe od mediany. Stąd też mediana bywa nazywana wartością środkową. Mediana jest oznaczana symbolem Me.

Własności mediany: odporna na wartości nietypowe (w przeciwieństwie do średniej).

**Dominanta** (*mode*), wartość najczęściej występująca. Jeżeli rozkład jest wielomodalny to dominanta jest nieokreślona. W szczególności w zbiorowościach o małej liczebności mogą być problemy z ustaleniem dominanty. Dominanta jest oznaczana symbolem D lub Mo.

**Kwartyle** (*quartile*): coś jak mediana tylko bardziej szczegółowo. Kwartyli jest trzy i dzielą one zbiorowość na 4 równe części, każda zawierająca 25% całości. Kwartyle oznaczane są symbolami  $Q_1, Q_2, Q_3, Q_4$ .

Pierwszy kwartył dzieli **uporządkowaną** zbiorowość w proporcji 25%–75%. Trzeci dzieli **uporządkowaną** zbiorowość w proporcji 75%–25%. Drugi kwartył to mediana.

**Kwantyle** (D, wartości dziesiętne), podobnie jak kwartyle, tyle że dzielą na 10 części.

**Centyle** (P, wartości setne), podobnie jak kwantyle tyle że dzielą na 100 części. Przykładowo wartość 99 centyla i mniejszą ma 99% jednostek w populacji.

#### Współczynnik dzietności na świecie w roku 2018

Średnia: 2,68. Interpretacja: średnia wartość współczynnika dzietności wyniosła 2,68 dziecka. Mediana: 2,2. Interpretacja mediany: współczynnik dzietności w połowie krajów na świecie wynosił 2,2 dziecka i mniej. Dominanta: 1,98. Interpretacja dominanty: najwięcej krajów wykazuje współczynnik dzietności równy 1,98. Wartość cokolwiek przypadkowa, bowiem krajów wykazujących współczynnik równy 1,98, jest raptem 4.

Uwaga: średnia dzietność na świecie **nie wynosi** 2,68 dziecka (bo po pierwsze uśredniamy kraje a nie kobiety a po drugie kraje różnią się liczbą ludności). Podobnie dzietność połowy kobiet na świecie wyniosła 2,2 dziecka i mniej jest niepoprawną interpretacją mediany (z tych samych względów jak w przypadku średniej).

**Generalna uwaga:** interpretacja średniej-średnich często jest nieoczywista i należy uważać. (a współczynnik dzietności jest średnią: średnia liczba dzieci urodzonych przez kobietę w wieku rozrodczym. Jeżeli liczymy średnią dla 202 krajów, to mamy *średnią-średnich*). Inny przykład: odsetek ludności w wieku poprodukcyjnym wg powiatów (średnia z czegoś takiego nie da nam odsetka ludności w wieku poproduk-

cyjnym w Polsce, bo powiaty różnią się liczbą ludności).

**Współczynnik dzietności (kontynuacja):**

Pierwszy kwartyl: 1,75; trzeci kwartyl 3,56 co oznacza że 25% krajów miało wartość współczynnika dzietności nie większą niż 1,75 dziecka a 75% krajów miało wartość współczynnika dzietności nie większą niż 3,56 dziecka.

### 2.4.2 Miary zmienności

Miary zmienności określają zmienność (dyspersję albo rozproszenie) w zbiorowości.

Rodzaje miar zmienności:

- Klasyczne:
  - wariancja i odchylenie standardowe,
- Pozycyjne:
  - rozstęp,
  - rozstęp ćwiartkowy.

**Wariancja** (*variance*) jest to średnia arytmetyczna kwadratów odchyleń poszczególnych wartości zmiennej od średniej arytmetycznej zbiorowości. Co można zapisać jako ( $\bar{x}$  oznacza średnią,  $N$  liczebność zbiorowości, a  $x_i$  wartość  $i$ -tej jednostki):

$$s^2 = \frac{1}{N} ((x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 + \dots + (x_N - \bar{x})^2)$$

Przy czym często zamiast dzielenia przez  $N$  dzielimy przez  $N - 1$ .

**Odchylenie standardowe** (*standard deviation*) jest pierwiastkiem kwadratowym z wariancji. Parametr ten określa przeciętną różnicę wartości zmiennej od średniej arytmetycznej. Odchylenie standardowe jest oznaczane symbolem  $s$ .

**Rozstęp** jest to różnica pomiędzy największą a najmniejszą wartością zmiennej w zbiorowości:  $R = x_{\max} - x_{\min}$ . Rzadko stosowana miara, bo jej wartość zależy za bardzo od wartości nietypowych.

**Rozstęp ćwiartkowy** (*interquartile range, IQR*) ma banalnie prostą definicję:  $R_Q = Q_3 - Q_1$ , gdzie:  $Q_1$ ,  $Q_3$  oznaczają odpowiednio pierwszy oraz trzeci kwartyl.

**Współczynnik dzietności (kontynuacja)**

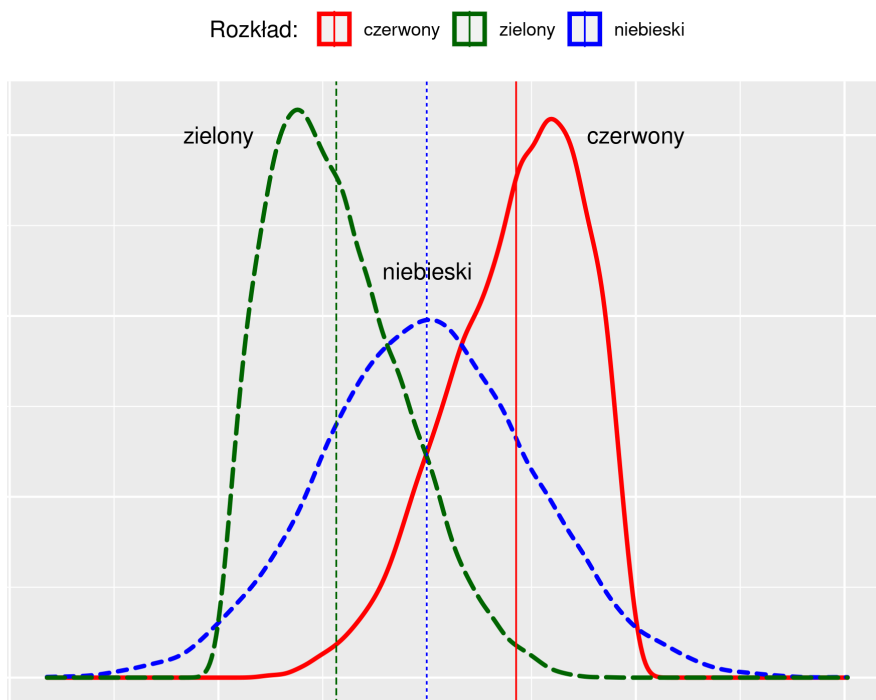
Średnie odchylenie od średniej wartości współczynnika wynosi 1,26 dziecka. Wartość rozstępu ćwiartkowego wynosi 1,81 dziecka. Wartość rozstępu wynosi 5,91 dziecka.

**Uwaga:** odchylenie standardowe/ćwiartkowe są miarami mianowanymi. Zawsze należy podać jednostkę miary.

### 2.4.3 Miary asymetrii

Asymetria albo skośność (*skewness*), to odwrotność symetrii. Szereg jest symetryczny jeżeli jednostki są rozłożone „równomiernie” wokół wartości średniej. W szeregu symetrycznym wartości średniej i mediany są sobie równe. Skośność może być dodatnia (*positive skew*) lub ujemna (*negative skew*). W przypadku asymetrii prawostronnej

większa część zbiorowości przyjmuje wartości poniżej średniej. W przypadku asymetrii lewostronnej jest odwrotnie. Rysunek 2.3 przedstawia rozkład symetryczny oraz rozkłady skośne.



Rysunek 2.3: Rozkłady symetryczne i asymetryczne

Miary asymetrii:

- klasyczny współczynnik asymetrii ( $g$ ):
  - Przyjmuje wartości ujemne dla asymetrii lewostronnej; a dodatnie dla prawostronnej. Teoretycznie może przyjąć dowolnie dużą wartość, ale w praktyce rzadko przekracza 3 co do wartości bezwzględnej.
  - Wartości większe od 2 świadczą o dużej, a większe od 3 o bardzo dużej asymetrii.
- współczynniki asymetrii Pearsona ( $W_s$ ):
  - Wykorzystuje różnice między średnią a dominantą:  $W_s = (\bar{x} - D)/s$ .
  - Wykorzystuje różnice między średnią a medianą:  $W_s = 3(\bar{x} - Me)/s$ .
- Współczynnik asymetrii oparty na odległościach między kwartylami:
  - Obliczany jest według następującej formuły:  $W_{sq} = ((Q_3 - Q_2) - (Q_2 - Q_1))/(Q_3 - Q_1)$ .

#### Współczynnik dzietności (kontynuacja)

Współczynnik asymetrii  $g$  wynosi 1,1; Współczynnik Pearsona wykorzystujący



dominantę wynosi 0,55 a wykorzystujący medianę 1,14. Wreszcie wartość współczynnika opartego o kwartle wynosi 0,5. Wszystkie współczynniki wskazują prawostronną (dodatnią) asymetrię. Współczynniki oparte o dominantę i kwartle wskazują słabą asymetrię, a oparte o medianę oraz klasyczny umiarkowaną. Jest całkowicie normalne, że różne miary wykazują różne wartości asymetrii.

## 2.5 Porównanie wielu rozkładów

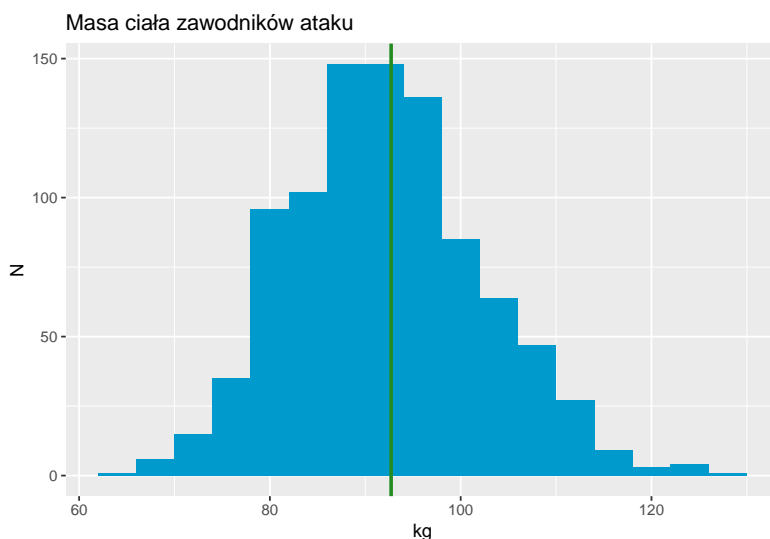
Często strukturę jednego rozkładu należy porównać z innym. Albo trzeba porównać strukturę wielu rozkładów. Pokażemy jak to zrobić na przykładzie.

### Masa ciała uczestników Pucharu Świata w Rugby

W grze w rugby drużyna jest podzielona na dwie **formacje**: ataku i młyna. Należy porównać rozkład masy ciała zawodników obu formacji uczestniczących w turniejach o puchar świata w Rugby w latach 2015, 2019 i 2023.

#### Zawodnicy ataku

Histogram przy przyjęciu długości przedziału równej 4kg (pionowa linia zielona oznacza poziom średniej):



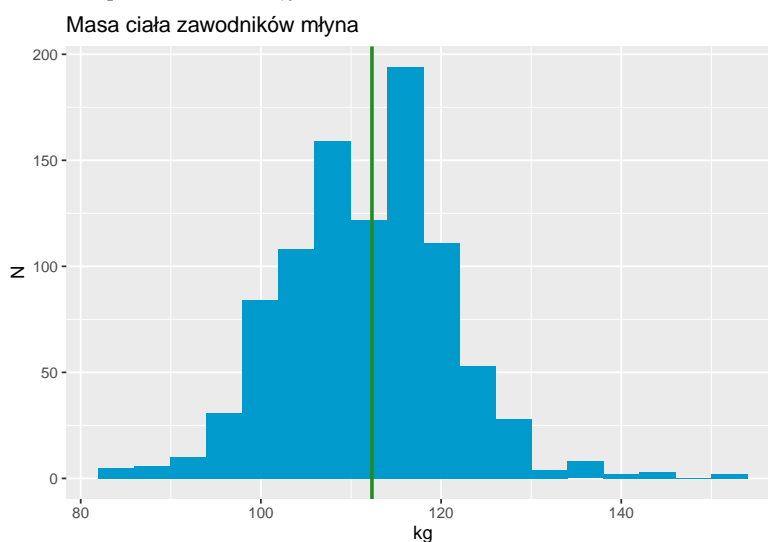
Liczba zawodników ataku wyniosła 936. Przeciętnie zawodnik ataku ważył 92,7 kg. Wartość mediany wyniosła 92 kg (połowa zawodników ataku ważyła 92 kg i mniej). Wartości pierwszego i trzeciego kwartyła wyniosły odpowiednio 85,5 oraz 99 kg (1/4 zawodników ataku ważyła 85,5 kg i mniej; 1/4 zawodników ataku ważyła 99 kg i więcej).

Odchylenie standardowe jest równe 10,1 kg (przeciętnie odchylenie od średniej arytmetycznej wynosi 10,1 kg). Rozstęp ćwiartkowy wynosi 13,5 kg (rozstęp 50% środkowych wartości wynosi 13,5 kg).

Wartość klasycznego współczynnika skośności jest równa 0,34. Wartość współczynnika skośności opartego o kwartle wynosi 0,04, a współczynnika skośności Pearsona wykorzystującego medianę wynosi 0,21.

### Zawodnicy młyna

Histogram przy przyjęciu długości przedziału równej 4kg (pionowa linia zielona oznacza poziom średniej):



Liczba zawodników młyna wyniosła 943. Średnio zawodnik młyna ważył 112,3 kg. Wartość mediany wyniosła 112 kg (połowa zawodników młyna ważyło 112 kg i mniej). Wartości pierwszego i trzeciego kwartyła wyniosły odpowiednio 106 oraz 118 kg (1/4 zawodników młyna ważyło 106 kg i mniej; 1/4 zawodników młyna ważyło 118 kg i więcej).

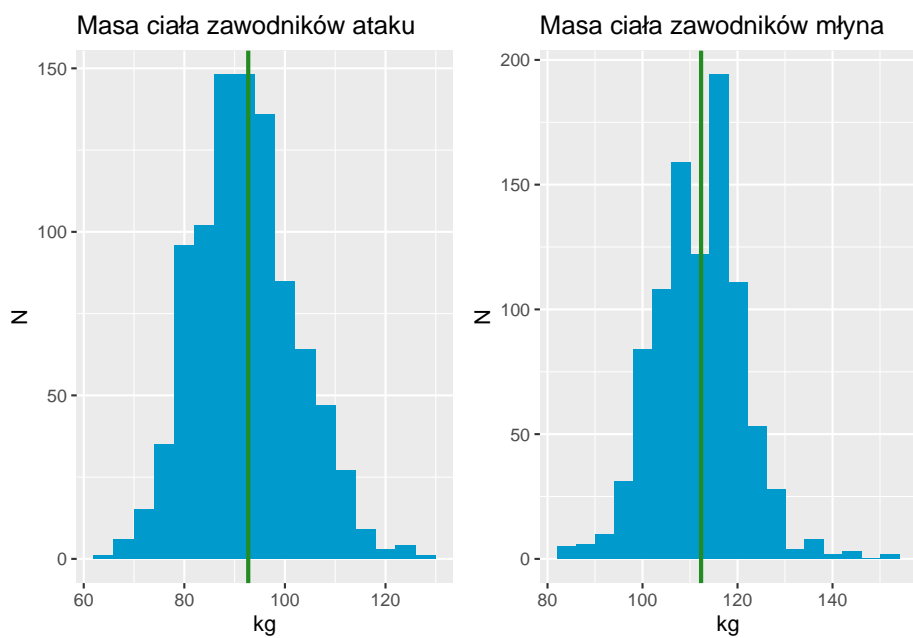
Odchylenie standardowe jest równe 9,2 kg (przeciętnie odchylenie od średniej arytmetycznej wynosi 9,2 kg). Rozstęp ćwiartkowy wynosi 12 kg (rozstęp 50% środkowych wartości wynosi 12 kg).

Wartość klasycznego współczynnika skośności jest równa 0,17. Wartość współczynnika skośności opartego o kwartle wynosi 0, a współczynnika skośności Pearsona wykorzystującego medianę wynosi 0,11.

### Porównanie atak vs młyn

miara	atak	młyn
średnia	92,71	112,33
mediana	92,00	112,00
odchyl.st	10,07	9,24
iqr	13,50	12,00
skośność	0,34	0,17

Średnio zawodnik młyna ważył prawie 20 kg więcej od zawodnika ataku (w przypadku mediany jest to dokładnie 20 kg więcej). Zmienność mierzona wielkością odchylenia standardowego oraz IQR jest w obu grupach podobna. Oba rozkłady są zbliżone do rozkładu symetrycznego.



### 2.5.1 Wykres pudełkowy

Do porównania wielu rozkładów szczególnie użyteczny jest wykres zwany pudełkowym (*box-plot*).

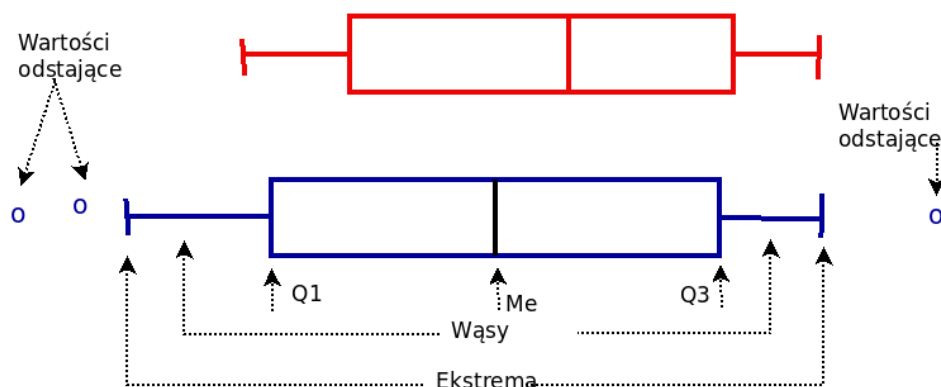
Pudełka na wykresie pudełkowym są rysowane według następujących zasad (por rysunek 2.4):

- lewy i prawy bok pudełka jest równy kwartylom;
- linia pionowa w środku pudełka jest równa medianie;
- linie poziome (zwane wąsami) mają długość równą  $Q_1 - 1,5IQR$  oraz  $Q_3 + 1,5IQR$  (dla przypomnienia:  $Q_1$ ,  $Q_3$  to kwartyle, zaś IQR to rozstęp ćwiartkowy);
- kropki przed oraz za wąsami to wartości zmiennej większe od  $Q_3 + 1,5IQR$  lub mniejsze od  $Q_1 - 1,5IQR$ .

Interpretacja pudełek:

- linia pozioma w środku pudełka określa przeciętny poziom zjawiska;
- długość pudełka oraz wąsów określa zmienność (im większe wąsy/długość pudełka tym większa zmienność);
- kropki przed oraz za wąsami to **obserwacje nietypowe** (albo **wartości odstające**).

Zatem dolny rozkład z rysunku 2.4 ma mniejszą wartość średnią oraz większą zmienność od rozkładu górnego. Dolny rozkład posiada też wartości odstające, a górny nie.



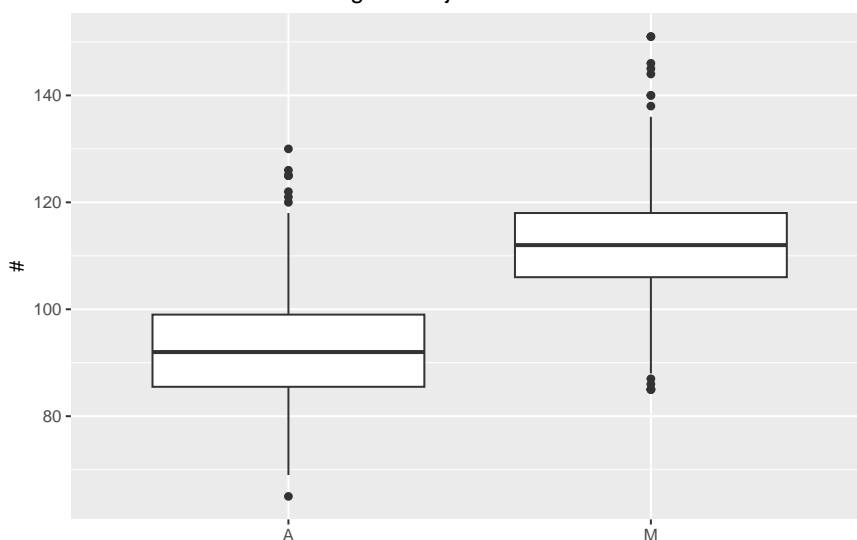
Rysunek 2.4: Wykres pudełkowy

Zwróć uwagę na następującą sztuczkę. Wartości nietypowe nie są definiowane jako na przykład górne/dolne 1% wszystkich wartości, bo wtedy **każdy rozkład** miałby wartości nietypowe, ale jako wartości mniejsze lub większe od  $Q_{1,3} \pm 1,5 \cdot IQR$ . Wszystkie wartości rozkładów o umiarkowanej zmienności mieszczą się wewnątrz tak zdefiniowanego przedziału.

Typowo wykres zawiera wiele pudełek, a każde pudełko wizualizuje jeden rozkład. Pudełka mogą być umieszczone jedno pod drugim, tak jak na rysunku 2.4 lub jedno obok drugiego jak na przykładach poniżej.

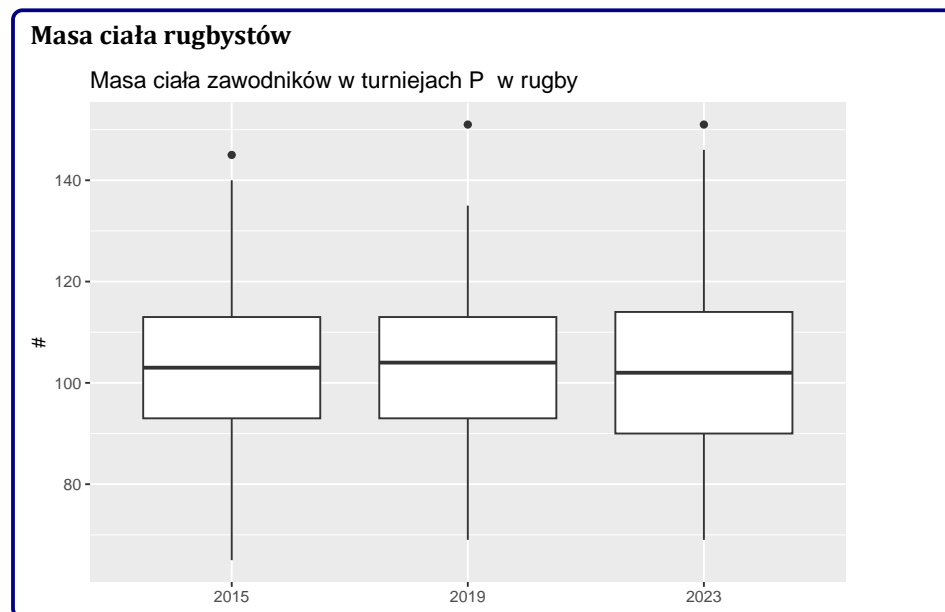
### Masa ciała rugbyistów

Masa ciała zawodników wg formacji



Z wykresu od razu widać, który rozkład ma wyższą średnią (M), który większe rozproszenie (A), oraz w którym występują wartości nietypowe.

Pudełek może być więcej niż dwa oczywiście. Następny przykład pokazuje porównanie rozkładów masy ciała zawodników rugby na poszczególnych turniejach.



Od razu widać, że przeciętnie najciężsi zawodnicy byli na turnieju w roku 2019; największe zróżnicowanie masy ciała występowało na turnieju w roku 2023.

## 2.6 Zestawienie metod opisu statystycznego

W rozdziale przedstawiono osiem sposobów opisanie rozkładu zmiennej:

1. Tablice statystyczne.
2. Wykres słupkowy.
3. Wykres kołowy (niezalecany).
4. Histogram.
5. Wykres pudełkowy.
6. Miary tendencji centralnej: średnia, mediana, kwartyle.
7. Miary rozproszenia: odchylenie standardowe, rozstęp ćwiartkowy.
8. Miary asymetrii.

## Rozdział 3

# Wprowadzenie do wnioskowania statystycznego

**Chcemy się dowiedzieć czegoś na temat populacji (całości) na podstawie próby (części tej całości).**

Przykładowo chcemy ocenić ile wynosi średnia waga główki kapusty na 100 h polu. Można ścinać wszystkie i zważyć, ale można też ścinać trochę (pobrać próbę się mówi uczenie) zważyć i poznać średnią na całym polu z dobrą dokładnością.

## 3.1 Masa ciała uczestników PŚ w rugby

W turnieju o Puchar Świata w rugby w 2015 roku uczestniczyło 623 rugbystów. Znamy szczegółowe dane odnośnie wzrostu i wagi każdego uczestnika turnieju. Obliczamy (prawdziwą) średnią, odchylenie standardowe i współczynnik zmienności masy ciała:

##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	65,0	93,0	103,0	102,8	113,0	145,0

Czyli średnio rugbysta na turnieju ważył 102,8 kg (Mean na wydruku powyżej) a odchylenie standardowe ( $s$ ) wyniosło 12,92 kg.

Rozkład, pokazany na rysunku 3.1, jest dwumodalny, bo w rugby są dwie grupy zawodników i wcale nie wszyscy ważą ponad 110 kilogramów.

**Szacujemy średnią na podstawie 2 zawodników pobranych losowo.**

Powtarzamy eksperyment 1000 razy (dwóch bo dla jednego nie obliczymy wariancji).

##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	80,0	96,5	102,5	102,6	109,0	132,0

Średnia (średnich z próby) ma wartość 102,61 kilogramów a odchylenie standardowe 9,06 kilogramów. Wartość  $s/\sqrt{2}$  (odchylenie standardowe podzielone przez pierwiastek kwadratowy z liczebności próby) jest równa 9,14. Zauważmy, że ta wartość jest zbliżona do wartości odchylenia standardowego uzyskanego w eksperymencie (9,06 vs 9,14).

Zauważmy też, że wartość najniższej średniej wyniosła 80 kilogramów zaś najwyższej 132 kilogramów. Gdybyśmy mieli pecha i wylosowali te skrajnie nieprawdziwe wartości to mylimy się o 22,61 kilogramów na minus lub 29,39 kilogramów na plus.

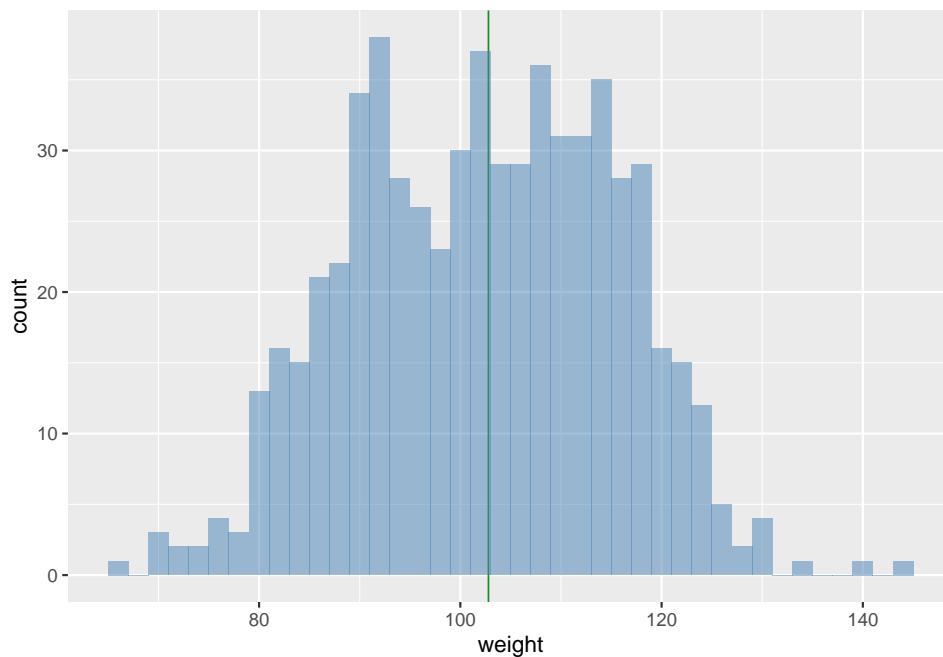
**Szacujemy średnią na podstawie 10 zawodników pobranych losowo.**

Powtarzamy eksperyment 1000 razy.

##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	91,8	99,9	102,6	102,6	105,3	116,9

Średnia ma wartość 102,63 kilogramów, a odchylenie standardowe 3,84 kilogramów. Wartość  $s/\sqrt{10}$  jest równa 4,09.

**Szacujemy średnią na podstawie 40 zawodników pobranych losowo.**



Rysunek 3.1: Rozkład wagi zawodników

Uwaga: 40 zawodników to 6,4% całości. Powtarzamy eksperyment 1000 razy.

##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	95,75	101,42	102,78	102,75	104,13	108,53

Średnia jest równa 102,75 kilogramów, a odchylenie standardowe 1,98 kilogramów. Wartość  $s/\sqrt{40}$  jest równa 2,04.

Zauważmy też, że wartość najniższej średniej wyniosła 95,75 kilogramów zaś najwyższej 108,525 kilogramów. Gdybyśmy mieli pecha i wylosowali te skrajnie nieprawdziwe wartości to mylimy się o 7 kilogramów na minus lub 5,78 kilogramów na plus. Niewątpliwie wynik znacznie lepszy niż dla próby dwuelementowej.

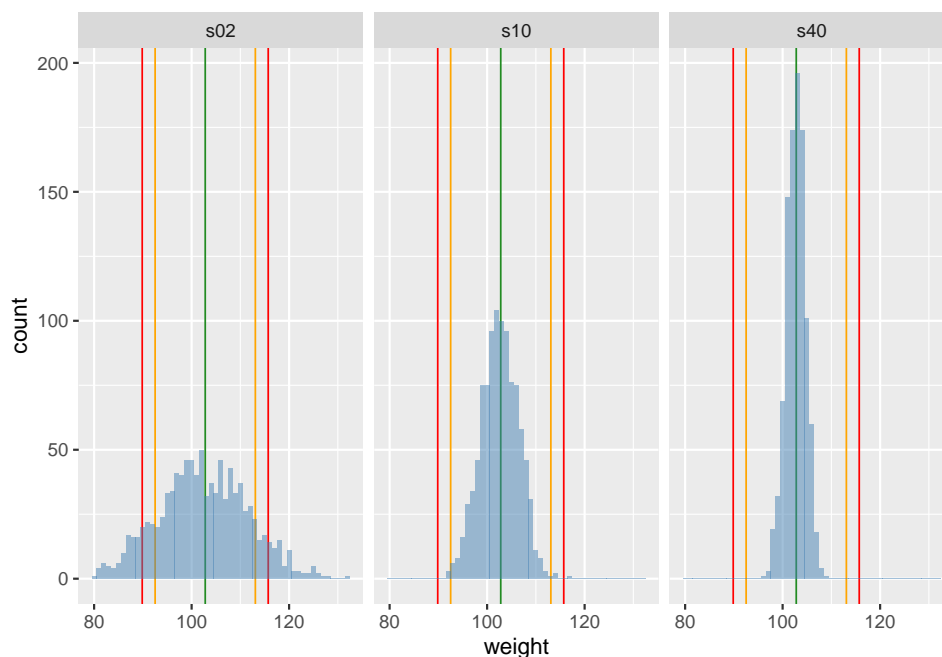
Podsumujmy eksperyment wykresem rozkładu wartości średnich (por. rysunek 3.2).

#### Wnioski z eksperymentu:

Wartość średnią wyznaczamy na podstawie jakiejś konkretnej **metody**. Wydaje się na podstawie powyższych eksperymentów, że z dobrym skutkiem możemy jako metodę wykorzystać **średnią-z-próby**.

W ogólności taką metodą, która formalnie jest funkcją elementów z próby, nazywa się w statystyce **estymatorem**. Warto to pojęcie zapamiętać. Wnioskujemy o wartości nieznanego parametru w populacji posługując się estymatorem.

Kontynuując wnioski z eksperymentu należy zauważyć, że wszystkie średnie-ze-średnich (bez względu na liczebność próby) są zbliżone do wartości prawdziwej (to się nazywa **nieobciążoność** estymatora); Mówiąc innymi słowy jeżeli będziemy oceniać wartość prawdziwej średniej na podstawie próby, a naszą ocenę powtórzymy



Rysunek 3.2: Rozkład średniej wagi rugbyistów w zależności od wielkości próby

wielokrotnie, to średnia będzie zbliżona do wartości prawdziwej (a nie np. niższa czy wyższa). Ta cecha jest niezależna od wielkości próby.

Jeżeli rośnie liczebność próby to zmienność wartości średniej-w-próbie maleje, co za tym idzie prawdopodobieństwo, że wartość oceniona na podstawie średniej z próby będzie zbliżona do wartości szacowanego parametru rośnie (to się nazywa **zgodność**). Co więcej, dobrym przybliżeniem zmienności średniej-w-próbie jest prosta formuła  $s/\sqrt{n}$  gdzie  $n$  jest liczebnością próby a  $s$  jest odchyleniem standardowym w populacji z której pobrano próbę.

Jeżeli mamy dwa różne estymatory służące do oszacowania parametru i oba są **nieobciążone** oraz **zgodne**, to który wybrać? Ten która ma **mniejszą wariancję**. Taki estymator nazywa się **efektywny**.

Estymator zatem powinien być **nieobciążony**, **zgodny** oraz **efektywny** (czyli mieć małą wariancję). Można matematycznie udowodnić, że pewien estymator ma tak małą wariancję, że niemożliwe jest wynalezienie czegoś jeszcze bardziej efektywnego. Takim estymatorem średniej w populacji jest średnia z próby.

Konkretną wartość estymatora dla konkretnych wartości próby nazywamy **oceną** (parametru).

## 3.2 Wiek kandydatów na radnych

W wyborach samorządowych w Polsce w roku 2018 o mandat radnego sejmików wojewódzkich ubiegało się 7076 kandydatów. Znamy szczegółowe dane odnośnie

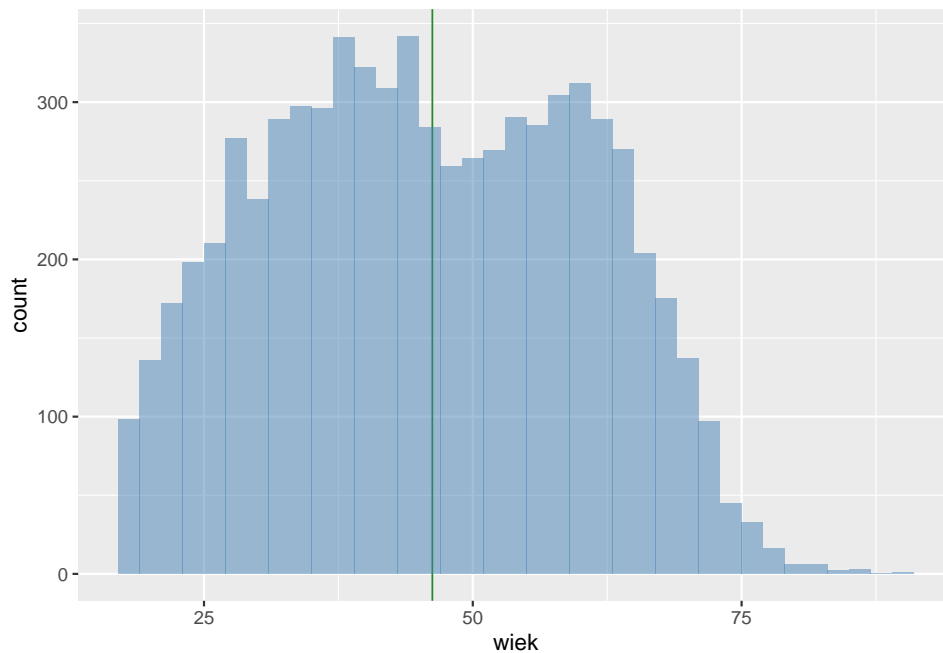


wieku każdego kandydata, bo to zostało publicznie podane przez Państwową Komisję Wyborczą. Obliczamy (prawdziwą) średnią, odchylenie standardowe i współczynnik zmienności wieku kandydatów:

##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	18,00	34,00	46,00	46,24	58,00	91,00

Czyli średnio kandydat miał 46,24 lat a odchylenie standardowe wieku wyniosło 14,61 lat.

Rozkład znowu jest dwumodalny z jakiś powodów (por. rysunek 3.3).



Rysunek 3.3: Rozkład wieku kandydatów na radnych

#### Szacujemy średnią na podstawie 2 kandydatów pobranych losowo.

Powtarzamy eksperyment 1000 razy.

##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	18,50	39,00	46,00	46,15	53,62	76,50

Średnia średnich z próby ma wartość 46,15 lat. Odchylenie standardowe wyniosło 10,55. Wartość  $s/\sqrt{2}$  jest równa 10,33.

Wartość najniższej średniej wyniosła 18,5 lat zaś najwyższej 76,5 lat. Gdybyśmy mieli pecha i wylosowali te skrajnie nieprawdziwe wartości to mylimy się o 27,65 lat na minus lub 30,35 lat na plus.

#### Szacujemy średnią na podstawie 10 kandydatów pobranych losowo.

Powtarzamy eksperyment 1000 razy.

##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	33,90	43,00	46,20	46,16	49,10	59,20

Średnia średnich z próby ma wartość 46,16 lat. Odchylenie standardowe wyniosło 4,48. Wartość  $s/\sqrt{10}$  jest równa 4,62.

**Szacujemy średnią na podstawie 40 kandydatów pobranych losowo**

Uwaga: 40 kandydatów to ok 0.6% całości. Powtarzamy eksperyment 1000 razy.

##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	38,85	44,58	46,15	46,20	47,90	53,05

Średnia średnich z próby ma wartość 46,2 lat. Odchylenie standardowe wyniosło 2,34. Wartość  $s/\sqrt{40}$  jest równa 2,31.

**Szacujemy średnią na podstawie 70 kandydatów pobranych losowo.**

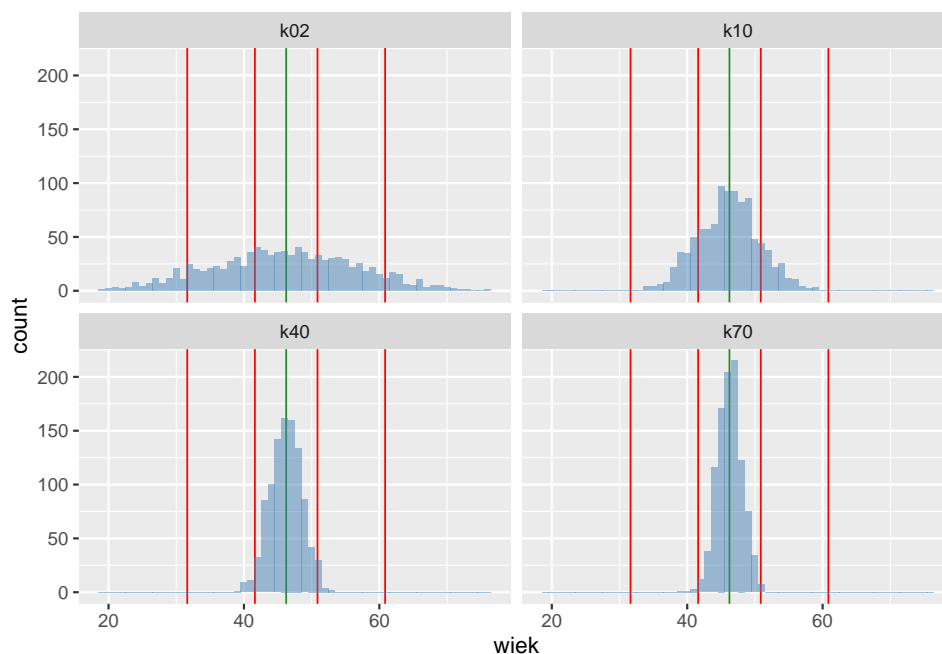
Uwaga: 70 kandydatów to około ok 1% całości (1000 powtórzeń).

##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	38,79	45,07	46,30	46,29	47,47	51,31

Średnia średnich z próby ma wartość 46,29 lat. Odchylenie standardowe wyniosło 1,82. Wartość  $s/\sqrt{70}$  jest równa 1,75.

Wartość najniższej średniej wyniosła 38,79 lat zaś najwyższej 51,31 lat. Gdybyśmy mieli pecha i wylosowali te skrajnie nieprawdziwe wartości to mylimy się już tylko o 7,5 lat na minus lub 5,02 lat na plus.

Podsumujmy eksperyment wykresem rozkładu wartości średnich (rysunek 3.4).



Rysunek 3.4: Rozkład średniej wieku kandydatów w zależności od wielkości próby

Obserwujemy to samo zjawisko, co w przypadku wagi rugbyistów: im większa próba, tym dokładniejsza wartość średniej wieku. Bez względu na wielkość próby przeciętnie otrzymujemy prawdziwą wartość średniej.

Wnioski:

- Precyzja wnioskowania zwiększa się wraz z liczebnością próby;
- Precyzja wnioskowania zwiększa się tym szybciej im rozproszenie w populacji generalnej jest mniejsze;
- Żeby z dużą dokładnością wnioskować o średniej dla dużej populacji wcale nie trzeba pobierać dużej próby (w ostatnim przykładzie było to 1% całości).

### 3.3 Rozkład normalny

**Rozkład empiryczny** zmiennej to przyporządkowanie kolejnym wartościom zmiennej odpowiadających im liczebności.

Założmy, że istnieje zapotrzebowanie społeczne na wiedzę na temat wieku kandydatów na radnych. Możemy to jak widać łatwo liczyć, ale jednocześnie jest to kłopotliwe. Należy do tego mieć zbiór ponad 7 tys liczb. **Rozkład teoretyczny** to matematyczne uogólnienie **rozkładu empirycznego**. Jest to model matematyczny operujący pojęciem (ściśle sformalizowanym) **prawdopodobieństwa** (zamiast liczebności).

**Rozkład teoretyczny** jest:

- zbliżony do empirycznego jeżeli chodzi o wyniki (jest przybliżeniem empirycznego);
- jest zdefiniowany za pomocą kilku liczb; nie ma potrzeby korzystania z liczebności.

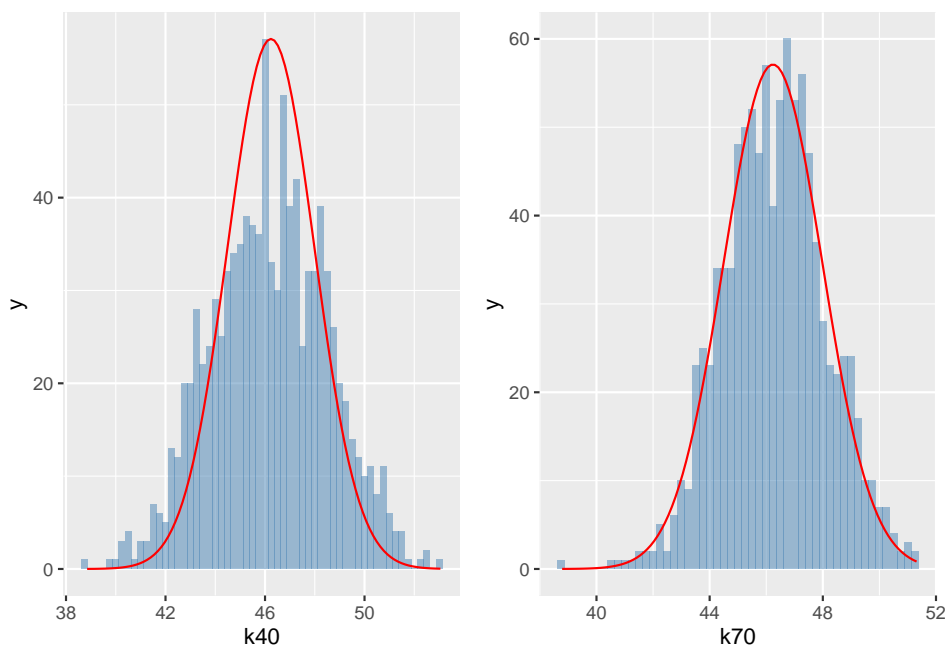
Okazuje się, że istnieje jeden **rozkład teoretyczny**, który z dobrą dokładnością opisuje rozkłady empiryczne będące wynikiem powyższej zabawy. Ten rozkład (zwany **normalnym**) zależy tylko od dwóch parametrów: średniej i odchylenia standardowego, gdzie średnia będzie równa (prawdziwej) średniej w populacji a odchylenie standardowe równe odchyleniu standardowemu w populacji podzielonemu przez pierwiastek z wielkości próby.

Przybliżenie za pomocą rozkładu normalnego średniego rozkładu wieku kandydatów na radnych dla próby 40- oraz 70-elementowej pokazuje rysunek 3.5.

Prawda, że wynik jest całkiem dobry? Teoretyczność czerwonej (w kolorowej wersji podręcznika) krzywej polega na tym, że ona zawsze będzie identyczna, podczas gdy histogram będzie różny. Gdybyśmy powtórzyli nasz eksperyment (generowania 1000 losowych prób przypominam), to zapewne trochę by się różnił, bo byśmy wylosowali inne wartości do prób. Ta **teoretyczna abstrakcja** określana jest **prawdopodobieństwem**. Rzucając monetą 1000 razy spodziewamy się po 500 orłów i reszek, co w modelu matematycznym będzie opisane jak: prawdopodobieństwo wyrzucenia orła wynosi 0,5. Rzucanie monetą to bardzo prosty eksperyment; nasz z liczeniem średniej wieku jest bardziej skomplikowany więc miło jest się dowiedzieć, że używając czerwonej krzywej można łatwo obliczyć jak bardzo prawdopodobne jest na przykład popełnienie błędu większego niż 10% średniej, albo większego niż 0,1 lat. Albo jak duża powinna być próba żeby ten błąd był nie większy niż 0,1 lat.

Interpretacja wartości rozkładu empirycznego zwykle jest w kategoriach ryzyka/szansy czy prawdopodobieństwa. Przykładowo interesuje nas prawdopodobieństwo, że kandydat ma mniej niż 30 lat. Takich kandydatów jest 1091 a wszystkich kandydatów dla przypomnienia jest 7076. Iloraz tych wartości będzie interpretowany jako ryzyko/szansa/prawdopodobieństwo (wynosi ono 15,42%).

Podobnie można obliczyć prawdopodobieństwo, że wiek kandydata będzie się



Rysunek 3.5: Rozkład normalny

zawierał w przedziale 50–60 lat. Ponieważ kandydatów w wieku 50–60 lat jest 1570, to szukane prawdopodobieństwo jest równe: 22,19%).

Jeżeli zamiast rozkładu empirycznego będziemy używać rozkład normalnego, który jak widzimy jest jego dobrym przybliżeniem, to nie musimy liczyć empirycznych liczebności. Wystarczy że znamy średnią i odchylenie standardowe a potrafimy obliczyć każde prawdopodobieństwo dla każdego przedziału wartości zmiennej.

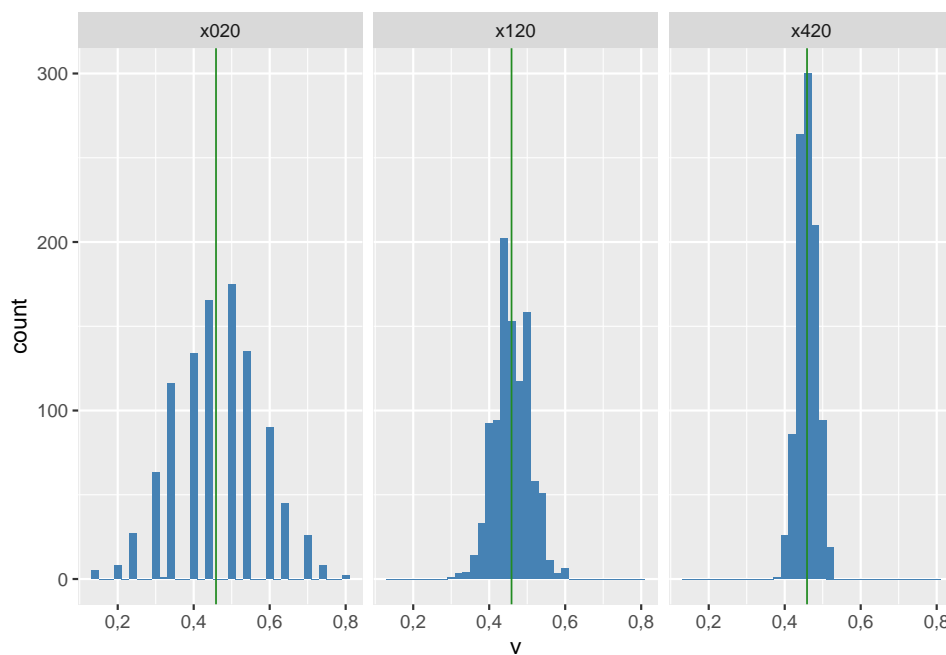
W szczególności dla rozkładu normalnego prawdopodobieństwo przyjęcia wartości z przedziału  $m \pm s$  (średnia plus/minus odchylenie standardowe) wynosi około 0,68 prawdopodobieństwo przyjęcia wartości z przedziału  $m \pm 2 \times s$  wynosi około 0,95 a przyjęcia wartości z przedziału  $m \pm 3 \times s$  około 0,997. Czyli w przedziale  $[-3s < m, m + 3s]$  znajdują się praktycznie wszystkie wartości rozkładu. Albo innymi słowy przyjęcie wartości spoza przedziału średnia plus/minus trzykrotność odchylenia standardowego jest bardzo mało prawdopodobne.

Za pomocą rozkładu normalnego można opisać rozkład wagi rugbyistów, wieku posłów, wagę noworodków i miliony innych rozkładów. Uogólnieniem teoretycznym pojęcia **zmiennej statystycznej**, które do tej pory używaliśmy jest **zmienna losowa**, tj. zmienna, która przyjmuje wartości liczbowe z określonym prawdopodobieństwem np. określonym przez rozkład normalny.

### 3.4 Odsetek kobiet wśród kandydatów na radnych

Dane dotyczące kandydatów na radnych do sejmików wojewódzkich zawierają także płeć kandydata. Ktoś może być ciekaw jaki był odsetek kobiet w tej grupie. Taki

parametr nazywa się proporcją albo ryzykiem, a potocznie i niefachowo procentem. Matematycznym modelem jest **zmienna dwuwartościowa**, która z określonym prawdopodobieństwem przyjmuje wartość kobieta. Obliczmy empiryczną wartość tego prawdopodobieństwa jako liczbę kobiet do liczby wszystkich kandydatów. Wartość tego parametru wynosi 0,4587 (albo 45,87%). Potraktujmy to jako prawdziwą wartość prawdopodobieństwa ( $p$ ), że kandydat jest kobietą i empirycznie sprawdzimy czy możemy szacować o prawdziwej wartości tego parametru używając (jako estymatora żeby się przyzwyczajać do nowych terminów) proporcji z próby. Tradycyjnie powtarzamy eksperyment 1000 razy dla trzech różnych wielkości próby. Rozkład otrzymanych wartości przedstawia rysunek 3.6.



Rysunek 3.6: Rozkład wielkości  $p$  dla różnych wielkości próby

#### Wnioski:

- Dla próby 20 elementowej rozkład nie przypomina rozkładu normalnego.
- Dla prób 120 i 420 elementowej rozkład jest podobny do normalnego.
- Zmienność estymatora maleje wraz ze wzrostem liczebności próby; każe nam to przypuszczać (i tak jest w istocie) że jest on zgodny.
- W każdym przypadku średnia z 1000 eksperymentów jest zbliżona do wartości prawdziwej; każe nam to przypuszczać (i tak jest w istocie) że estymator jest nieobciążony.

Rozkład normalny jest tak magiczny że nawet jeżeli zmienna, której parametr szacujemy nie ma rozkładu zbliżonego do normalnego (jak w przypadku zmiennej, która przyjmuje tylko dwie wartości) to i tak estymator tego parametru będzie normalny. Co najwyżej będziemy potrzebowali większej próby żeby „znormalniał” (jak w opisywanym przykładzie).

## 3.5 Wnioskowanie statystyczne

Celem analizy danych z próby jest **uogólnienie** uzyskanych wyników na całą populację. To uogólnienie nazywa się wnioskowaniem (*inference*). Przypominamy, że **wnioskujemy** o wartości parametru w populacji posługując się **estymatorem**. W przypadku wnioskowania o średniej estymatorem jest średnia-z-próby. Dobrze by było wiedzieć jak bardzo wiarygodna jest ta wartość (zwana oceną parametru) uzyskana na podstawie konkretnego estymatora, inaczej mówiąc jak dużo mogliśmy się pomylić.

Do oceny tej wiarygodności można użyć wariancji-średniej-z-próby, która nazywa się **wariancją błędu** (*error variance*). Jeżeli wariancja błędu jest duża, to w pojedynczej próbie mogą wystąpić wartości znacznie różniące się od prawdziwej średniej; jeżeli jest mała to wartości bardzo różniące się od prawdziwej średniej mają małe szanse na zaistnienie. W przypadku rozkładu normalnego wiemy, że wariancja błędu jest równa  $s^2/n$  (gdzie  $s^2$  jest wariancją w populacji, a  $n$  wielkością próby).

W ramach wnioskowania stosowane są trzy metody (podejścia):

- estymacja punktowa,
- estymacja przedziałowa,
- testowanie hipotez.

### 3.5.1 Estymacja punktowa

Szacujemy średnią (albo inny parametr) i tę wartość uznajemy za wartość prawdziwą; dokładność szacunku jest nieokreślona. Inaczej mówiąc wartość **estymatora** dla konkretnej próby przyjmujemy za ocenę parametru.

Estymatorem punktowym średniej jest średnia z próby a estymatorem punktowym proporcji/ryzyka jest proporcja/ryzyko z próby.

### 3.5.2 Estymacja przedziałowa

Nie można ustalić prawdopodobieństwa popełnienia błędu dla dokładnej wartości parametru (co wynika z właściwości matematycznych modelu), ale można dla dowolnego przedziału od-do.

Czyli nie można ustalić, że z prawdopodobieństwem 95% oszacujemy wartość średnią czegoś jako 5,000000, ale można z prawdopodobieństwem 95% oszacować **przedział**, w którym znajdzie się średnia (przykładowo, że będzie to przedział 4,9–5,1).

Estymacja przedziałowa to oszacowanie przedziału wartości od-do, który z zadanym z góry prawdopodobieństwem zawiera prawdziwą wartość parametru.

Z góry wyznaczone prawdopodobieństwo nazywa się **poziomem ufności** (określa jak często mamy się **NIE pomylić**).

### 3.5.3 Testowanie hipotez

Większość analiz statystycznych polega na porównaniu. W wyniku tego porównania otrzymujemy liczbę. Załóżmy, że mamy dwie próby dotyczące wieku kandydatów na radnych do sejmików wojewódzkich z roku 2018 (średnia 46,1) oraz z roku 2014 (47,2). Różnica wynosi 1,1 lat i może być spowodowana błędem przypadkowym (tj.

gdybyśmy wylosowali jeszcze raz dwie próby, to wynik byłby zupełnie odmienny np 46,9 vs 46,5) i/lub wynikać z tego, że faktycznie w roku 2014 kandydaci byli starsi.

Formalnie stawiamy **hipotezę**, że różnica średnich wynosi zero. Jest to tzw. **hipoteza zerowa**. Niezbędne jest także postawienie **hipotezy alternatywnej**, którą może być proste zaprzeczenie zerowej. Zapisuje się to następująco ( $m_{14}/m_{18}$  oznacza odpowiednio średnie w latach 2014/2018):

$H_0$ : różnica średnich wieku wynosi zero ( $m_{14} = m_{18}$ )

$H_1$ : różnica średnich wieku jest różna od zera ( $m_{14} \neq m_{18}$ )

Hipotezy sprawdzamy wykorzystując **test statystyczny** czyli zmienną losową, której rozkład prawdopodobieństwa zależy (jest funkcją powiedziała by matematyk) od wartości testowanych parametrów (w tym przypadku  $m_{14}$  oraz  $m_{18}$ ). Tę zmienną losową nazywa się **statystyką testu**.

Nie jest chyba wielkim zaskoczeniem, że **statystyką testu** w teście różnicy średnich jest różnica średnich w próbie (poprawnie mówiąc różnica uwzględniająca liczebność próby oraz zmienność obu populacji). Całkiem **zdroworozsądkowo** możemy przyjąć, że duże wartości **statystyki testu** świadczą na rzecz hipotezy alternatywnej, natomiast małe na rzecz hipotezy zerowej.

Duża różnica pomiędzy **hipotezą** a wynikiem z próby może wynikać z tego, że:

1. Pechowo trafiła nam się nietypowa próba, który zdarza się rzadko (rozkład normalny).
2. Hipoteza jest fałszywa, średnie mają inną wartość niż zakładamy w hipotezie zerowej.

Statystyk zawsze wybierze drugą wersję. Pozostaje tylko ustalić co to jest rzadko (dla statystyka)?

Rzadko, to z prawdopodobieństwem mniejszym niż z góry ustalone małe prawdopodobieństwo. zwane **poziomem istotności**. Określa ono jak często możemy się pomylić **odrzucając hipotezę zerową, która jest prawdziwa**.

Teraz wystarczy obliczyć prawdopodobieństwo wystąpienia różnicy, którą otrzymaliśmy lub jeszcze większej i porównać je z poziomem istotności. Jeżeli to prawdopodobieństwo jest równe lub niższe od poziomu istotności odrzucamy hipotezę zerową (różnica jest istotna statystycznie).

Przyjmijmy przykładowo, że prawdopodobieństwo wystąpienia różnicy 1,1 lat (i większej) oszacowane na podstawie odpowiedniego modelu matematycznego (rozkład normalny) wynosi 0,3 co znaczy że coś takiego zdarza się względnie często – trzy razy na 10 pobranych prób.

Załóżmy z kolei że, ta różnica wyniosła 3,2 lata. Prawdopodobieństwo wystąpienia takiej różnicy (i większej) wynosi 0,009 co znaczy że coś takiego zdarza się względnie rzadko – 9 razy na tysiąc prób.

Przyjmując, że możemy się mylić 5 razy na 100 w pierwszym przypadku statystyk powie, że nie ma podstaw do odrzucenia hipotezy  $H_0$ . Różnica 1,1 lat wynika z przypadku. W drugim wypadku statystyk powie, że hipoteza jest fałszywa, bo zdarzyło się coś co nie powinno się zdarzyć.

Ale jest jeszcze drugi przypadek popełnienia błędu: **przyjmujemy hipotezę zerową, która jest fałszywa**. W testach statystycznych nie określa się prawdopodobieństwa popełnienia tego błędu, a w związku z tym nie można **przyjąć hipotezy zerowej** (bo nie znamy ryzyka popełnienia błędu).

W konsekwencji hipotezę zerową albo się odrzuca albo **nie ma podstaw do od-**

**rzucenia.** Wniosek cokolwiek niekonkluzywny, ale tak jest.

Dlatego też często „opłaca się” tak postawić hipotezę zerową aby ją następnie odrzucić, bo taki rezultat jest bardziej konkretny.

### 3.5.4 Testy nieparametryczne

Można testować hipotezy na temat wartości parametrów, ale można też testować przypuszczenia o charakterze mniej konkretnym. Na przykład, że dwie zmienne są niezależne (co to znaczy wyjaśniono w następnym rozdziale), albo że dwa rozkłady są podobne do siebie (rozkłady nie średnie). Takie hipotezy/testy określa się jako **nieparametryczne**. Przykładami są testy niezależności chi-kwadrat albo normalności Shapiro-Wilka (opisane w następnym rozdziale).

Oczywiście, ale powtórzmy: przypuszczenia o charakterze nieparametrycznym możemy tylko testować (sprawdzać hipotezy); nie obliczamy wtedy ani ocen ani nie wyznaczamy przedziałów ufności.

## 3.6 Statystyk Carl Pearson

W punkcie 2.3 przypomnieliśmy postać Florence Nightingale – matki statystyki i bardzo dobrej kobiety. A kto był ojcem tejże statystyki? Ojców było więcej niż matek oczywiście, a wśród nich Francis Galton (regresja), Carl Pearson (współczynnik korelacji liniowej, test niezależności chi-kwadrat) oraz Ronald Fisher (podstawy wnioskowania). Niestety wszyscy wymienieni byli zadeklarowanymi rasistami oraz wyznawcami społecznego darwinizmu i eugeniki. Pierwszymi zastosowaniami „nowoczesnych” metod statystycznych było naukowe udowodnienie, że biali ludzie są lepsi od innych:

*Przez ile stuleci, ile tysięcy lat Kaffirowie [...] lub Murzyni rządili w Afryce nie niepokojeni przez białych ludzi? Jednak ich walki międzyplemienne nie stworzyły cywilizacji w najmniejszym stopniu porównywalnej z **aryjską** [...] Historia pokazuje jeden i tylko jeden sposób, w jaki powstaje wysoka cywilizacja, a mianowicie **walka rasy** i przetrwanie rasy sprawniejszej fizycznie i psychicznie...*

Powyższe to cytaty z *National Life from the standpoint of science* Carla Pearsona (Londyn 1905).

Naszym zdaniem dobrze jest pamiętać o tym fatalnym starciu „nowoczesnej statystyki”, bo chociaż jest mało prawdopodobne, że zostanie ona znowu wykorzystana do równie odrażających celów, to jest raczej więcej niż pewne, że będzie użyta do innych szwindli. Jeszcze jeden argument żeby nie traktować wyników analiz statystycznych jako wiedzy objawionej, absolutnie pewnej i 100% prawdziwej (por. uwagę w punkcie 1.1).

## 3.7 Słownik terminów, które warto znać

1. Estymacja (punktowa, przedziałowa): szacowanie wartości parametru na podstawie próby.
2. Estymator (nieobciążony, zgodny, efektywny): funkcja na wartościach próby która służy do oszacowania parametru.



3. Hipoteza statystyczna: przypuszczenie dotyczące parametru lub rozkładu zmiennej.
4. Ocena (parametru): konkretna wartość estymatora dla pewnej próby.
5. Poziom istotności (testu; oznaczany jako  $\alpha$ ; zwykle 0,05): prawdopodobieństwo popełnienia błędu.
6. Poziom ufności = prawdopodobieństwo, że przedział ufności zawiera prawdziwą wartość parametru; oznaczany jako  $1 - \alpha$ ; zwykle 0,95.
7. Rozkład (prawdopodobieństwa): przypisanie prawdopodobieństwa wartościom zmiennej losowej.
8. Test statystyczny: metoda weryfikacji hipotezy statystycznej.
9. Wnioskowanie statystyczne: wnioskowanie o całości na podstawie próby.

## Rozdział 4

# Analiza współzależności pomiędzy zmiennymi

Pomiędzy zjawiskami występują związki (zależności.) Nauki formułują te związki w postaci **praw**. Jak takie **prawo naukowe** powstaje? Typowo w dwu etapach, najpierw za pomocą **dedukcji** stawia się **hipotezę**, potem konfrontuje się hipotezę z danymi (podejście hipotetyczno-dedukcyjne). Na tym drugim etapie używa się statystyki (lub matematyki jeżeli prawo ma charakter deterministyczny).

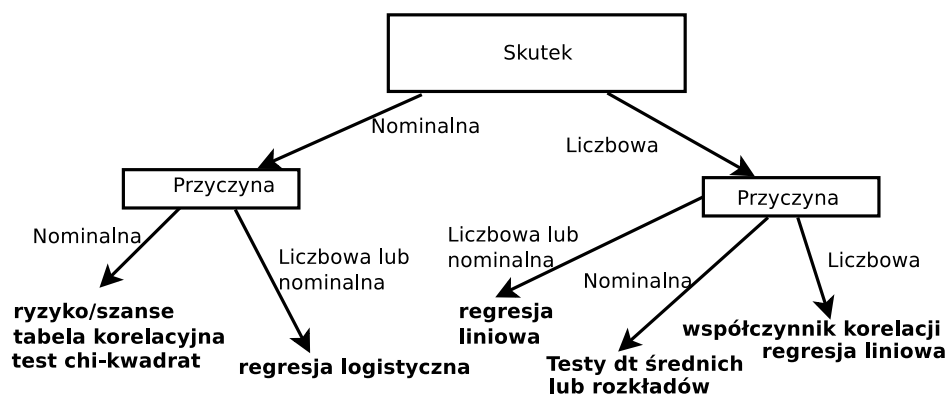
Upraszczając *metoda hypodedukcji* sprowadza się do dedukcyjnego sformułowania hipotezy, która następnie jest empirycznie *falsyfikowana*, tj. próbuje się wykazać, że jest ona nieprawdziwa. Konsekwencje: nie można dowieść prawdziwości żadnej hipotezy, można natomiast wykazać, że hipoteza jest fałszywa.

Związki między zmiennymi mogą być albo **funkcyjne** – wartościom jednej zmiennej odpowiada tylko jedna wartość drugiej zmiennej lub **stochastyczne** – wartościom jednej zmiennej odpowiadają z pewnym przybliżeniem wartości innej zmiennej.

Problem: czy istnieje związek (zależność) pomiędzy cechami? Przykładowo czy istnieje związek pomiędzy paleniem (przyczyna) a chorobą nowotworową (skutek), wiekiem a prawdopodobieństwem zgonu z powodu COVID19 itd.

Jaki jest charakter zależności? Jaka jest siła zależności?

Rodzaj konkretnej metody zastosowanej do empirycznej weryfikacji zależy w szczególności od sposobu pomiaru danych (nominalne, porządkowe, liczbowe), co pokazuje na rysunku 4.1.



Rysunek 4.1: Metody statystycznej weryfikacji zależności pomiędzy zmiennymi

Optymistyczną informacją jest, że metod (oznaczonych krojem pogrubionym na diagramie), które omawiamy dalej w rozdziale, jest raptem siedem, czyli niedużo.

## 4.1 Dwie zmienne nominalne

### 4.1.1 Ryzyko względne oraz iloraz szans

Ryzyko to udział (iloraz) liczby sukcesów do liczby prób (zdarzeń pozytywnych/wyróżnionych do wszystkich). Zwykle podawany w procentach. Warto zauważyć, że jest to empiryczny odpowiednik prawdopodobieństwa.

#### Podawanie witaminy C a przeziębienie/brak przeziębienia

Eksperyment, który przeprowadził Linus Pauling (laureat nagrody Nobla za odkrycie witaminy C), polegał na tym, że podzielił 280 narciarzy na dwie 140 osobowe grupy. Przez 5–7 dni podawał witaminę C jednej grupie oraz placebo drugiej grupie. Obserwował zachorowania na przeziębienie przez następne dwa tygodnie. Jeden narciarz nie dokończył eksperymentu. Historia milczy dlaczego :-)

W grupie 139 narciarzy, którym podano witaminę C (grupa C) zachorowało 17. W grupie 140 narciarzy, którym podano placebo (grupa P) zachorowało 31. Zatem:

- Ryzyko zachorowania w grupie C wyniosło  $17/139 = 12,2\%$ .
- Ryzyko zachorowania w grupie P wyniosło  $31/140 = 22,14\%$ .

Na tzw. chłopski rozum jeżeli witamina C **nie działa**, to powinien zachorować ten sam odsetek narciarzy w obu grupach. A tak nie jest jak widać.

Prostymi miarami oceny siły zależności mogą być:

- różnica ryzyk (**risk difference**),
- ryzyko względne (**relative risk**), oraz
- iloraz szans (**odds ratio**).

Jeżeli  $r_e$  oznacza ryzyko w grupie eksperymentalnej (*test group*, grupa narażona, *exposed group*), a  $r_k$  w grupie kontrolnej (*control group*, grupa nienarażona, *unexposed group*), to **różnica ryzyk** to po prostu  $r_e - r_k$ . W przykładzie będzie to  $22,14 - 12,2 = -9,94\%$ . Ta miara aczkolwiek prosta jest rzadko stosowana.

Znacznie częściej używa się **ryzyka względnego** definiowanego jako  $RR = r_e/r_k$ . Oczywiście jest, że  $RR < 1$  oznacza zmniejszenie ryzyka;  $RR > 1$  oznacza zwiększenia;  $RR = 1$  oznacza brak zależności.

Zamiast ryzyka (czyli ilorazu liczby sukcesów do liczby prób) można używać pojęcia szans/szansy (*odds*) definiowanego jako iloraz sukcesów do porażek.

Jeżeli  $o_e$  oznacza szanse w grupie eksperymentalnej a  $o_k$  w grupie kontrolnej, to **iloraz szans** (*odds ratio*), jest definiowany jako stosunek  $OR = o_e/o_k$ .

Przykładowo jeżeli w dwukrotnym rzucie monetą otrzymano orła i reszkę to ryzyko otrzymania orła wynosi  $1/2 = 0,5$  a szansa otrzymania orła wynosi  $1/1 = 1$ .

#### Narciarze Paulinga (kontynuacja)

Przypomnijmy: ryzyko zachorowania w grupie C wynosi 12,2; ryzyko zachorowania w grupie P wyniosło 22,14. Ryzyko względne wynosi zatem  $12,2/22,14 = 0,55$ . Podanie witaminy C zmniejsza ryzyko zachorowania o prawie połowę.

Szansa, że narciarz z grupy C zachoruje wynosi  $17/122 = 13,9\%$ . Szansa, że narciarz w grupie P zachoruje wynosi 28,44%.

Zatem iloraz szans dla narciarzy wyniesie  $13,9/28,44 = 0,48$ . Podanie witaminy C zmniejsza szansę na zachorowanie o ponad połowę. Albo  $1/0,48 = 2,04$ , co oznacza, że narciarz, który nie brał witaminy C ma ponad dwukrotnie większą szansę na zachorowanie.

Jak widać dla dużych ryzyk (rzut monetą) szansa różni się znacznie od prawdopodobieństwa, ale dla małych ryzyk obie miary mają zbliżoną wartość.

Właściwości ilorazu szans:

- jeżeli równe 1 to sukces/porażka równie prawdopodobne;
- jeżeli większe od 1 to sukces jest bardziej prawdopodobny;
- jeżeli jest mniejsze od 1 to porażka jest bardziej prawdopodobna.

Dane w badaniach wykorzystujących ryzyko/szanse mają często postać następującej tabeli dwudzielnej o wymiarach  $2 \times 2$  (a, b, c i d to liczebności):

Grupa	sukces	porażka
grupa kontrolna	a	b
grupa eksperymentalna	c	d

Dla danych w tej postaci:

$$RR = (a/(a + b))/(c/(c + d))$$

$$OR = (ad)/(bc)$$

#### Narciarze Paulinga (tabela dwudzielna)

Grupa	katar	zdrowy
grupa C	17	122
grupa P	31	109

$$RR = (17/(17 + 122))/(31/(31 + 109)) = 0,55 \text{ oraz}$$

$$OR = (17 * 109)/(31 * 122) = 0,48$$

Otrzymaliśmy oczywiście identyczny wynik jak w poprzednim przykładzie.

### 4.1.2 Przedziały ufności dla ryzyka względnego i ilorazu szans

Ryzyko, ryzyko względne czy iloraz szans to parametry podobne do odsetka kobiet wśród kandydatów na radnych z przykładu w poprzednim rozdziale. Wiemy, że estymatorem punktowym proporcji jest proporcja z próby. Nie będzie wielkim odkryciem, że estymatorem punktowym ryzyka jest ryzyko z próby, ryzyka względnego/ilorazu szans zaś ryzyko względne/iloraz szans z próby.

Standardem jest obliczanie dla ryzyka względnego oraz ilorazu szans oprócz ocen punktowych także przedziałów ufności czyli podawania dwóch wartości, pomiędzy

którymi z zadanym prawdopodobieństwem znajduje się nieznana wartość szacowanego parametru.

#### Narciarze Paulinga (przedziały ufności)

Końce przedziałów ufności dla ilorazu szans (ocena punktowa 0,4899524) wynoszą: [0,2569389; 0,934282] zaś dla ryzyka względnego (ocena punktowa 0,5523323) przedział ufności wynosi [0,3209146; 0,9506298].

**Uwaga:** nie jest specjalnie istotne jaka jest konkretna formuła obliczania przedziałów ufności, przecież obliczenia i tak koniec-konców wykona program komputerowy.

Przedział ufności dla ilorazu szans nie zawiera 1; zatem branie witaminy C zmniejsza szansę na zachorowanie; albo zwiększa na niezachorowanie od  $1/0,25 = 4$  do  $1/0,9 \approx 1,1$ . Żeby to zabrzmiało ładnie i po polsku: zwiększa na niezachorowanie od 10% do 300%.

Dlaczego taka znacząca rozpiętość? Bo próba jest względnie mała. Gdyby Pauling zwerbował nie 280 a 2800 narciarzy mógłby weryfikować działanie swojej witaminy z większą pewnością.

### 4.1.3 Tabele wielodzielne

Łączny rozkład dwóch lub większej liczby zmiennych można przedstawić w tabeli. Taka tabela nazywa się dwudzielna (dla dwóch zmiennych) lub wielodzielna albo wielodzielcza (dla więcej niż dwóch zmiennych). Inne nazwy dla tabel wielodzielnych to krzyżowe albo kontyngencji (*cross-tabulation*, *contingency* albo *two-way tables*).

Ograniczmy się do analizy tabel dwudzielnych.

#### Narciarze Paulinga (kontynuacja)

Eksperyment Paulinga można przedstawić w postaci tablicy dwudzielnej (P/C oznacza czy narciarz zażywał witaminę czy placebo; cold/nocold czy zachorował czy nie zachorował na katar):

	nocold	cold	razem
C	122	17	139
P	109	31	140
Sum	231	48	279

Taka tabela składa się z wierszy i kolumn. Dolny wiersz (Sum czyli Razem po polsku) zawiera łączną liczebność dla wszystkich wierszy w danej kolumnie. Podobnie prawa skrajna kolumna zawiera łączną liczebność dla wszystkich kolumn dla danego wiersza. Dolny wiersz/Prawą kolumnę nazywamy **rozkładami brzegowymi**. Pozostałe kolumny oraz wiersze nazywane są **rozkładami warunkowymi**. Rozkładów warunkowych jest tyle ile wynosi suma  $r + c$  gdzie  $r$  to liczba wariantów jednej cechy a  $c$  to liczba wariantów drugiej cechy.

Przy warunku że narciarz brał witaminę C, 122 takich osób nie zachorowało (**nocold**) a 17 zachorowało (**cold**). Drugi rozkład warunkowy: 109 narciarzy, któ-

rzy brali placebo nie zachorowało, a 31 zachorowało. Są także rozkłady warunkowe dla drugiej cechy. W grupie narciarzy, którzy zachorowali 122 brało witaminę C, a 109 brało placebo. Wreszcie w grupie narciarzy, którzy nie zachorowali 109 brało witaminę C, a 31 brało placebo. Rozkładów warunkowych jest 4 bo obie cechy mają po dwa warianty. Jest to najmniejsza możliwa tabela wielodzielna.

Zamiast liczebności można posługiwać się odsetkami (procentami):

	nocold	cold	razem
C	43,73	6,09	49,82
P	39,07	11,11	50,18
Sum	82,80	17,20	100,00

Narciarze, którzy brali witaminę C oraz nie zachorowali stanowią 43,73% wszystkich narciarzy. Mało przydatne...

Ciekawsze jest obliczenie procentów każdego wiersza osobno, tj. dzielimy liczebności w każdej kolumnie przez liczebności rozkładu brzegowego (wartości ostatniej kolumny):

	nocold	cold	razem
C	87,77	12,23	100
P	77,86	22,14	100
Sum	82,80	17,20	100

Otrzymaliśmy ryzyka zachorowania na katar (lub nie zachorowania). Ryzyko zachorowania dla całej grupy wynosi 17,2% a nie zachorowania 82,8%. Jest przynajmniej całkiem **zdroworozsądkowym założeniem** (uczenie hipotezą statystyczną), że jeżeli przyjmowanie witaminy nie ma związku z zachorowaniem lub nie na katar, to w grupie tych co brali i tych co nie brali powinniśmy mieć identyczne rozkłady warunkowe równe rozkładowi brzegowemu. Czyli powinno przykładowo zachorować 17,2% narciarzy, którzy brali witaminę C a widzimy, że zachorowało jedynie 12,23%.

Na oko księgowego witamina C działa (bo są różnice), ale dla statystyka liczy się czy ta różnica jest na tyle duża, że (z założonym prawdopodobieństwem) można wykluczyć działanie przypadku.

Rozumowanie jest następujące: jeżeli prawdopodobieństwo wystąpienia tak dużej różnicy jest małe, to cechy nie są niezależne. Jest to istota i jedyny wniosek z czegoś co się nazywa testem istotności-chi-kwadrat. Test chi-kwadrat porównuje liczebności tablicy wielodzielnej z idealną-tablicą-wielodzielną, która zakłada niezależność jednej zmiennej od drugiej.

Można udowodnić, że taka idealna tablica powstanie przez przemnożenie dla każdego elementu tablicy odpowiadających mu wartości brzegowych a następnie podzieleniu tego przez łączną liczebność (czyli przykładowo pierwszy element poniższej „idealnej” tablicy to 231 pomnożone przez 139 i podzielone przez 279; proszę sprawdzić, że jest to 115,086):

	N	Y	Sum
C	115,086	23,914	139
P	115,914	24,086	140
Sum	231,000	48,000	279

Proszę zwrócić uwagę że **rozkłady brzegowe** są identyczne, identyczna jest też łączna liczebność. Różnią się tylko rozkłady warunkowe (które nie są liczbami całkowitymi, ale tak ma być i nie jest to błąd).

Za pomocą testu chi-kwadrat obliczamy jakie jest prawdopodobieństwo wystąpienia tak dużych lub większych różnic. Wynosi ono 0,0419. Czyli wystąpienie tak dużych różnic pomiędzy **oczekiwanymi** (przy założeniu o niezależności zmiennych) a obserwowanymi liczebnościami zdarza się około 4 razy na 100.

Jeszcze raz przypominamy ideę testu: jeżeli prawdopodobieństwo zaobserwowania różnic jest małe to zakładamy że:

- albo mamy pecha i pięć razy podzując monetą zawsze nam spadła reszka (prawdopodobieństwo około 0,03), albo
- że założenie co do niezależności jest fałszywe.

Statystyk zawsze wybierze drugie. Pozostaje tylko ustalenie co to znaczy **małe**.

Małe to takie które jest mniejsze od arbitralnie przyjętego przez statystyka. Zwykle jest to 0,05 lub 0,01 (czasami 0,1) co oznacza że odrzucając założenie o braku związku pomiędzy katarą a braniem witaminy C pomylimy się pięć lub raz na 100.

**Uwaga:** proszę zwrócić uwagę że wniosek z testu niezależności jest słabszy niż z porównania ryzyk. Tam mamy informację że zależność istnieje i oszacowaną jej wielkość (np. za pomocą ryzyka względnego) tutaj tylko zweryfikowaliśmy fakt czy obie zmienne są niezależne lub nie.

#### Palenie a status społeczno-ekonomiczny

Dla pewnej grupy osób odnotowujemy ich status-społeczno-ekonomiczny (wysoki/**high**, średni/**middle**, niski/**low**) oraz status-względem-palenia (wartości: pali/**current**, palił-nie-pali/**former**, nigdy-nie-palił/**never**). Obie zmienne są nominalne, obie mają po trzy wartości. Można poklasyfikować wszystkich badanych w następujący sposób:

	High	Low	Middle	Sum
current	51	43	22	116
former	92	28	21	141
never	68	22	9	99
Sum	211	93	52	356

Uwaga: status-społeczno-ekonomiczny to powiedzmy miara prestiżu używana w socjologii (można na Wikipedii doczytać co to dokładnie jest).

Tym razem tabela składa się z 3 wierszy i 3 kolumn (ostatni wiersz/kolumna się nie liczą bo to sumy – rozkłady brzegowe).

Przedstawmy tę tabelę w postaci udziałów procentowych sumujących się dla każdego wiersza osobno do 100% (tj. dzielimy liczebności w każdej kolumnie przez liczebności rozkładu brzegowego (wartości ostatniej kolumny):

	High	Low	Middle	Sum
current	43,96552	37,06897	18,965517	100
former	65,24823	19,85816	14,893617	100
never	68,68687	22,22222	9,090909	100
Sum	59,26966	26,12360	14,606742	100

Rozumowanie jest identyczne jak dla narciarzy Paulinga. Jeżeli nie ma zależności pomiędzy paleniem a statusem to procenty w ostatnim wierszu powinny być identyczne jak w wierszach 1–3 (nagłówka nie liczymy). Tym idealnym procentom odpowiadają następujące liczebności:

	High	Low	Middle	Sum
current	68,75281	30,30337	16,94382	116
former	83,57022	36,83427	20,59551	141
never	58,67697	25,86236	14,46067	99
Sum	211,00000	93,00000	52,00000	356

Wartość prawdopodobieństwa dla testu chi-kwadrat określająca, że przy założeniu niezależności obu zmiennych tak duża różnica między liczebnościami rzeczywistymi a idealnymi (porównaj stosowne tabele wyżej) jest dziełem przypadku wynosi 0,001. Jest to prawdopodobieństwo tak małe, że statystyk odrzuca założenie o niezależności statusu i palenia (myląc się w przybliżeniu  $0,001 \approx$  raz na tysiąc).

## 4.2 Zmienna liczbowa i zmienna nominalna

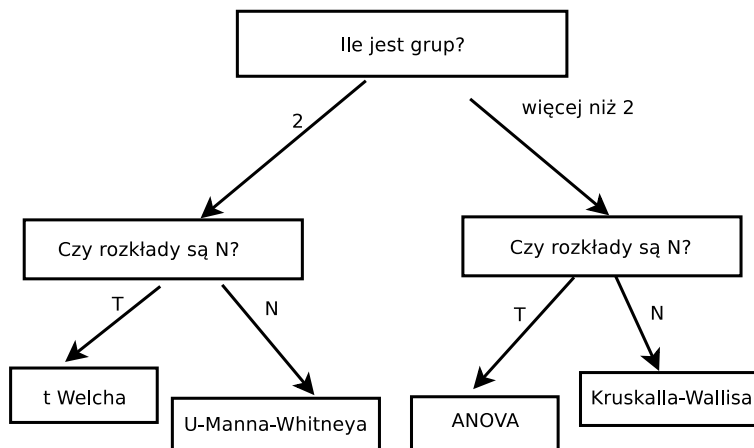
Obserwujemy wartości zmiennej liczbowej **w grupach** określonych przez wartości zmiennej nominalnej, np. wypalenie zawodowe w podziale na miejsce pracy. Grup może być dwie lub więcej.

Stawiamy hipotezę, że wartości zmiennej w każdej grupie są równe, wobec hipotezy alternatywnej, że tak nie jest (że są różne jeżeli grup jest dwie; co najmniej jedna jest różna jeżeli grup jest więcej niż dwie). Stosujemy odpowiedni test statystyczny:

- jeżeli liczba grup wynosi 2 oraz można przyjąć założenie o przybliżonej normalności rozkładów, to stosujemy test *t* Welcha;
- jeżeli liczba grup wynosi 2, ale nie można założyć normalności rozkładów, to stosujemy test U-Manna-Whitneya;
- jeżeli liczba grup jest większa niż dwie oraz można przyjąć założenie o normalności rozkładów, to stosujemy test ANOVA z poprawką Welcha;
- jeżeli liczba grup jest większa od dwóch oraz nie można przyjąć założenia o normalności rozkładów, to stosujemy test Kruskala-Wallisa.

Powyższe w postaci diagramu ze strzałkami przedstawiono na rysunku 4.2.





Rysunek 4.2: Testowanie istotności różnicy pomiędzy średnimi

Każdy z testów jest interpretowany identycznie:

1. Obliczana jest wartość statystyki testu  $t_k$ .
2. Obliczane jest prawdopodobieństwo  $t \geq t_k$  czyli przyjęcia przez statystykę testu  $t$  równej lub większej od  $t_k$  (co do wartości bezwzględnej). To prawdopodobieństwo zwyczajowo oznacza się literą  $p$  albo  $p$ -value (czyli wartość  $p$ ).
3. Jeżeli  $p$  jest mniejsze/równe od przyjętego poziomu istotności to hipotezę zerową odrzucamy; jeżeli  $p$  jest większe od przyjętego poziomu istotności to nie ma podstaw do odrzucenia hipotezy zerowej.

Odrzucenie hipotezy zerowej oznacza, że istnieje związek pomiędzy jedną a drugą zmienną. Jeżeli nie ma podstaw do odrzucenia hipotezy zerowej to oznacza to że takiej zależności nie udało nam się wykazać.

Omawiając wynik należy podać się wartość  $t_k$  oraz  $p$ . Statystyka testu może się różnie nazywać i być oznaczana różnym symbolem, np.:  $t$  (test  $t$  Welcha),  $U$  (test  $U$  Manna-Whitneya).

Testy  $t$  Welcha oraz ANOVA są **parametryczne**, porównujemy średnie w grupach. Testy  $U$ -Manna-Whitneya oraz Kruskala-Wallisa są **nieparametryczne** porównujemy rozkłady wartości zmiennej w grupach.

#### 4.2.1 Test $t$ Welcha

Test stosujemy jeżeli porównujemy dwie grupy oraz można przyjąć założenie, że rozkład wartości w obu grupach jest normalny. Test  $t$  Welcha jest testem parametrycznym. Sprawdzamy czy średnie w grupach są równe.

##### Poziom depresji a miejsce pracy

Studenci pielęgniarstwa i ratownictwa PSW w 2023 roku wypełnili ankietę zawierającą test depresji Becka, mierzący **poziom depresji** (wartość liczbową) oraz pytanie o rodzaj miejsca pracy (skala nominalna). Poniżej zestawiono średnie wartości **poziomu depresji** w podziale na rodzaj miejsca pracy (szpital/przychodnia).

m-pracy	średnia	mediana	n
Przychodnia	7,833333	7	12
Szpital	13,252747	11	91

Kolumna n zawiera liczebności.

Średnie różnią się o 5,42. Pytanie czy to dużo czy mało?

Przyjmijmy (na razie bez sprawdzania), że rozkłady wartości poziomu depresji w obu grupach są normalne. Można zatem zastosować test  $t$  Welcha.

Grupa1	Grupa2	n1	n2	t	p
Przychodnia	Szpital	12	91	-2,895988	0,00978

Kolumna t zawiera wartość statystyki testu  $t_k$ . Kolumna p zawiera oczywiście wartość prawdopodobieństwa  $p$ .

Ponieważ wartość  $p$  równa 0,00978 jest mniejsza od każdego zwyczajowo przyjmowanego poziomu istotności (0,05 na przykład, albo 0,1) hipotezę, że średnie w obu grupach są równe należy odrzucić. W konsekwencji stwierdzamy, że poziom depresji pracujących w Szpitalu był istotnie wyższy od pracujących w Przychodni.

#### 4.2.2 Testowanie normalności

Statystyk nie przyjmuje założeń na słowo honoru. Kiedy zatem można przyjąć założenie o normalności a kiedy nie? Można to ocenić na podstawie wykresu kwantylowego oraz posługując się testem Shapiro-Wilka (bo statystycy na każde pytanie mają zawsze **jakiś** stosowny test).

##### Poziom depresji a miejsce pracy

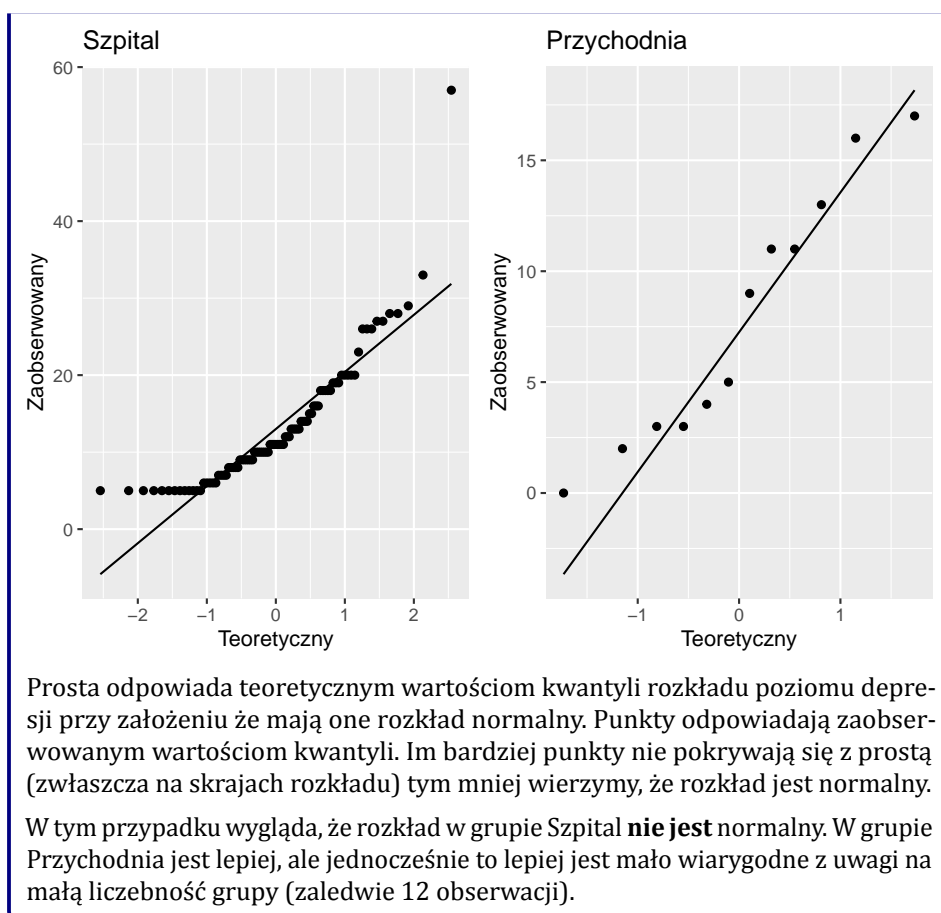
Hipoteza zerowa w teście Shapiro-Wilka (S-W) zakłada, że rozkład cechy jest normalny. Interpretacja tego testu jest „standardowa”, mianowicie małe wartości  $p$  świadczą przeciwko hipotezie zerowej.

m-pracy	S-W	p
Przychodnia	0,9256178	0,3359655
Szpital	0,8191903	0,0000000

Kolumna S-W zawiera wartości statystyki testu S-W oczywiście.

Rozkład w grupie szpital nie jest normalny (o czym świadczy niska wartość  $p$ ). Nasze założenie co do normalności było niepoprawne i należy do weryfikacji hipotezy o równości wartości zmiennej w grupach zamiast testu  $t$  Welcha zastosować test U Manna-Whitneya.

Wykres kwantylowy jest graficzną metodą weryfikacji normalności rozkładu zmiennej. Dla **poziomu depresji** wygląda jak na poniższym rysunku:



### 4.2.3 Test U Manna-Whitneya

Test U Manna-Whitneya jest testem nieparametrycznym. Sprawdzamy czy rozkłady zmiennej w grupach są identyczne.

#### Poziom depresji a miejsce pracy

Ponieważ grup jest dokładnie 2 a rozkład nie jest normalny, stosujemy test U Manna-Whitneya.

Grupa1	Grupa2	n1	n2	U	p
Przychodnia	Szpital	12	91	317	0,0185

Prawdopodobieństwo wystąpienia tak dużej wartości statystyki testu (U) przy założeniu, że rozkłady zmiennej w obu grupach są identyczne wynosi 0,0185. Różnica jest zatem istotna na każdym zwyczajowo przyjmowanym poziomie istotności (0,05 na przykład, albo 0,1); oba rozkłady różnią się. Kolumna U zawiera wartość

statystyki testu U. Przypominamy, że dobry zwyczaj nakazuje podawać tę wartość omawiając wynik testu (więc ją podajemy).

Depresja wśród pracowników szpitali jest wyższa niż wśród pracowników przychodni (w przypadku testu U Manna-Whitneya porównujemy wartości mediany).

#### 4.2.4 Test ANOVA

Jeżeli liczba grup jest większa niż dwie, ale można przyjąć założenie że w każdej grupie zmienna ma rozkład normalny, to stosujemy test ANOVA z poprawką Welcha. Test ANOVA jest testem parametrycznym. Sprawdzamy czy średnie w wszystkich grupach są równe.

##### Poziom depresji a staż pracy

W ankiecie, którą wypełnili Studenci pielęgniarstwa i ratownictwa PSW w 2023 roku było też pytanie o staż pracy. Oryginalną liczbową wartość zmiennej staż zamieniono na zmienną w skali nominalnej o następujących czterech wartościach: <6 (oznacza od 0 do 6 lat stażu pracy), 07-12 (7-12 lat), 13-18 (13-18 lat) oraz >19 (19 i więcej lat).

staż (kategoria)	średnia	mediana	n
<06	12,84615	11,0	39
07-12	12,14286	9,0	7
13-18	10,91667	6,5	12
>19	12,95556	11,0	45

Zakładając, że rozkłady w grupach są normalne, do weryfikacji hipotezy o równości wszystkich średnich możemy zastosować test ANOVA z poprawką Welcha. Na poniższym wydruku kolumna F zawiera wartość statystyki testu ANOVA, a kolumna p jak zwykle wartość prawdopodobieństwa p:

```
##
## One-way analysis of means (not assuming equal variances)
##
## data: P and staz
## F = 0,14844, num df = 3,000, denom df = 20,847, p-
value = 0,9295
```

Wartość p równa 0,9295 świadczy o tym, że nie ma istotnych różnic pomiędzy średnimi, co oznacza, że pomiędzy poziomem depresji a kategoriami stażu pracy nie ma zależności.

Czy zastosowanie testu ANOVA było poprawne? Żeby się o tym przekonać trzeba zastosować (znowu) test Shapiro-Wilka:

m-pracy	S-W	p
<06	0,9127826	0,0052385
07-12	0,8939017	0,2956402
13-18	0,8138678	0,0135286
>19	0,7239373	0,0000001

Wobec takiego wyniku testu do oceny istotności różnic należy zastosować bardziej ogólny test Kruskala-Wallisa.

#### 4.2.5 Test Kruskala-Wallisa

Test Kruskala-Wallisa jest testem nieparametrycznym. Sprawdzamy, czy rozkłady zmiennej w grupach są identyczne.

##### Poziom depresji a staż pracy

Na poniższym wydruku wartość statystyki testu jest oznaczona jako Kruskal-Wallis chi-squared a wartość p symbolem p-value:

```
##
## Kruskal-Wallis rank sum test
##
## data: P by staz
## Kruskal-Wallis chi-squared = 2,6982, df = 3, p-value = 0,4405
```

Prawdopodobieństwo tak dużej wartości statystyki testu przy założeniu, że rozkłady wartości zmiennej we wszystkich grupach są identyczne wynosi 0,4405 (różnice są zatem nieistotne; wszystkie rozkłady są identyczne i nie ma zależności).

### 4.3 Dwie zmienne liczbowe

#### 4.3.1 Korelacyjny wykres rozrzutu

Wykres rozrzutu (*scatter plot*) znany także jako korelogram, albo wykres XY, to prosty wykres kreślony w układzie kartezjańskim, w którym każdej obserwacji (składającej się z dwóch liczb) odpowiada kropka o współrzędnych XY.

O występowaniu związku świadczy układanie się kropek według jakiegoś kształtu (krzywej). O braku związku świadczy chmura punktów niepodobna do żadnej krzywej.

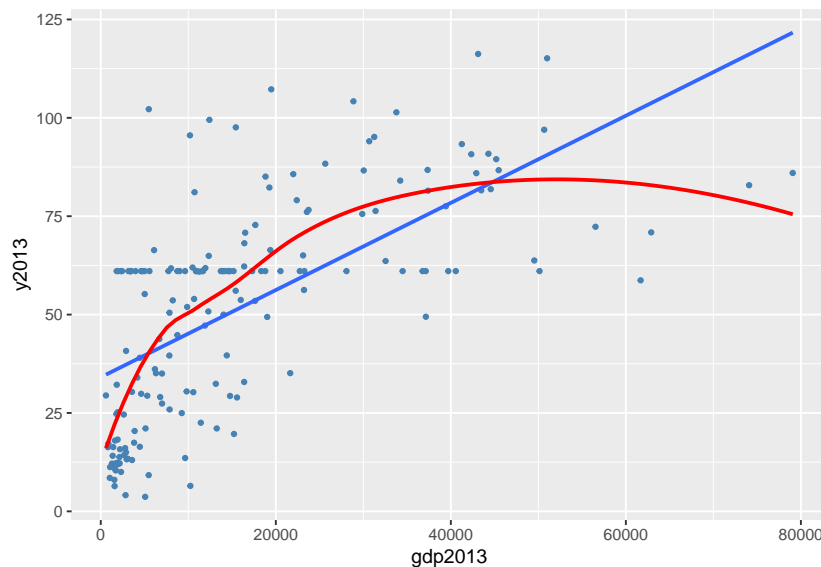
Punkty układające się według prostej świadczą o zależności liniowej (wyjątek: linia pozioma lub pionowa o czym dalej) zaś punkty układające się według krzywej świadczą o zależności nieliniowej.

##### Zamożność a konsumpcja mięsa

Organizacja Narodów Zjednoczonych do spraw Wyżywienia i Rolnictwa znana

jako FAO udostępnia dane dotyczące konsumpcji żywności na świecie (<https://www.fao.org/faostat/en/#home>). Bank światowy udostępnia dane dotyczące dochodu narodowego (<https://data.worldbank.org/>).

Konsumpcja mięsa jest mierzona jako średnia roczna konsumpcja w kilogramach w każdym kraju (*per capita* się mówi). Dochód podobnie jako średnia wielkość dochodu narodowego *per capita*. Dane dotyczą roku 2013.



Przy dużej dozie wyobraźni można dostrzec relację liniową pomiędzy konsumpcją mięsa a GDP co oznaczono na wykresie linią prostą. Można też założyć, że relacja pomiędzy konsumpcją mięsa a GDP ma charakter nieliniowy (linia krzywa). Liniowa czy nieliniowa, relacja jest na pewno mocno przybliżona co jest najbardziej pewnym wnioskiem, który można wysnuć z wykresu rozrzutu.

#### 4.3.2 Pomiar siły zależności: współczynnik korelacji liniowej Pearsona

Kowariancja to średnia arytmetyczna iloczynów odchyłeń wartości zmiennych  $X$ ,  $Y$  od ich wartości średnich. Dla  $n$  obserwacji na zmiennych  $X$  oraz  $Y$  można powyższe zapisać w postaci następującej formuły:

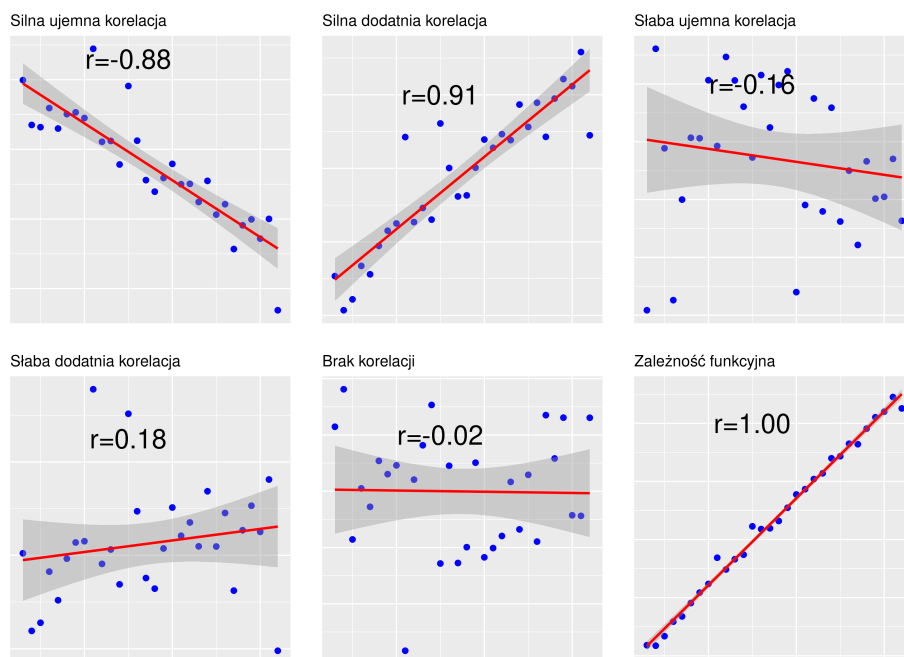
$$\text{cov}(xy) = \frac{1}{n} ((x_1 - \bar{x})(y_1 - \bar{y}) + \dots + (x_n - \bar{x})(y_n - \bar{y}))$$

Kowariancja zależy od rozproszenia (im większe tym większa), ma też dziwną jednostkę (jednostka $X$  · jednostka $Y$ ) oraz zależy od wybranych skal (tony vs gramy na przykład).

Z powyższych powodów do pomiaru związku pomiędzy cechami używa się standaryzowanego współczynnika kowariancji, zwanego **współczynnikiem korelacji liniowej Pearsona**, (*Pearson correlation coefficient*). Standaryzacja polega na podzieleniu wartości kowariancji przez iloczyn odchyłeń standardowych  $s_x$  oraz  $s_y$ .

$$r_{xy} = \frac{\text{cov}(xy)}{s_x \cdot s_y}$$

Współczynnik jest miarą niemianowaną, przyjmującą wartości ze zbioru  $[-1; 1]$ ; Skrajne wartości  $\pm 1$  świadczą o związku funkcyjnym (wszystkie punkty układają się na linii prostej); wartość zero świadczy o braku związku co odpowiada linii poziomej lub pionowej (por. rysunek 4.3).



Rysunek 4.3: Wykresy rozrzutu dla korelacji o różnej sile

Interpretacja opisowa: wartości powyżej 0,9 świadczą o silnej zależności.

#### Zamożność a konsumpcja mięsa (kontynuacja)

Współczynnik korelacji liniowej wynosi 0,6823158 (umiarkowana korelacja).

Czy ta wartość jest istotnie różna od zera? Jest na to stosowny test statystyczny, który sprowadza się do określenia jakie jest prawdopodobieństwo otrzymania  $r = 0,6823158$  przy założeniu że prawdziwa wartość  $r$  wynosi zero. Otóż w naszym przykładzie to prawdopodobieństwo wynosi  $3.850676e-26$  (czyli jest ekstremalnie małe –  $r$  jest istotnie różne od zera).

### 4.3.3 Macierz korelacji

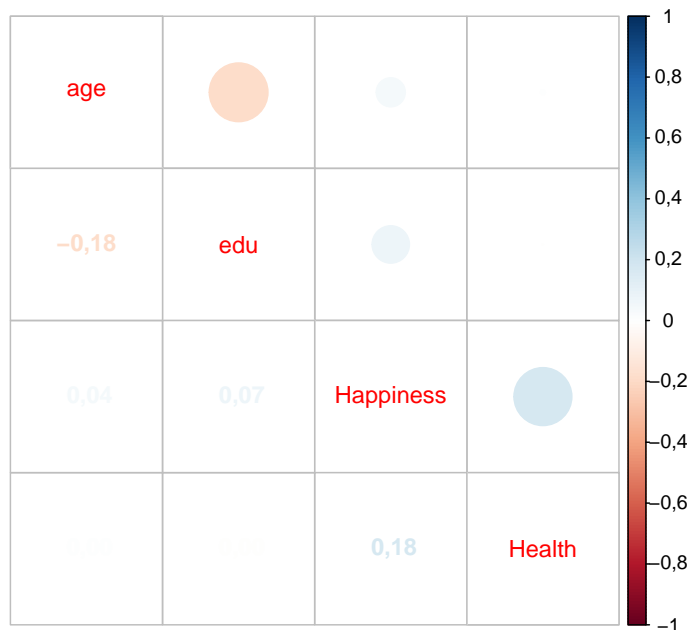
Wstępnym etapem analizy zależności między zmiennymi jest często hurtowa ocena współczynników korelacji w postaci kwadratowej **macierzy korelacji**.

### Korelacja pomiędzy wiekiem, edukacją, szczęściem a stanem zdrowia

Mohammadi S. i inni badali zależność pomiędzy wiekiem, poziomem edukacji, szczęściem a stanem zdrowia. (The relationship between happiness and self-rated health: A population-based study of 19499 Iranian adults; <https://doi.org/10.1371/journal.pone.0265914>).

##	age	edu	Happiness	Health
## age	1,0000	-0,1834	0,0449	0,0013
## edu	-0,1834	1,0000	0,0742	0,0000
## Happiness	0,0449	0,0742	1,0000	0,1786
## Health	0,0013	0,0000	0,1786	1,0000

Albo w bardziej efektownej postaci tekstowo-graficznej:



Ze wszystkich zmiennych analizowanych w badaniu Mohammadiego i innych jedynie zależność pomiędzy wiekiem a wykształceniem (raczej trywialna) oraz szczęściem i zdrowiem (raczej oczywista) okazały się znacząco różne od zera.

#### 4.3.4 Pomiar siły zależności: regresja liniowa

**Regresja liniowa** zakłada, że istnieje związek przyczyna-skutek i ten związek można opisać linią prostą (stąd liniowa). Skutek jest jeden i nazywa się go **zmienną zależną** a przyczyn może być wiele i noszą nazwę **zmiennych niezależnych** (albo **predyktorów**). W przypadku gdy związek dotyczy dwóch zmiennych mówi się o **regresji prostej**. Przykładowo zależność pomiędzy spożywaniem kawy w czasie sesji egzaminacyjnej a wynikiem egzaminu można formalnie zapisać jako:



$$\text{wynik} = b_0 + b_1 \cdot \text{kawa}$$

Współczynnik  $b_1$  określa wpływ spożycia kawy na wynik egzaminu. W szczególności jeżeli  $b_1 = 0$  to nie ma związku między spożywaniem kawy a wynikiem egzaminu.

#### 4.3.5 Regresja prosta

Równanie regresji dla zmiennych  $Y$  (skutek) oraz  $X$  (przyczyna) można zapisać następująco:

$$Y = b_0 + b_1 \cdot X + e$$

$Y = b_0 + b_1 \cdot X$  to **część deterministyczna**, a  $e$  oznacza **składnik losowy**. O tym składniku zakładamy, że średnia jego wartość wynosi zero. Można to sobie wyobrazić następująco: w populacji jest jakaś prawdziwa zależność  $Y = b_0 + b_1 \cdot X$  pomiędzy  $X$  a  $Y$ , która w próbie ujawnia się z błędem o charakterze losowym. Ten błąd może wynikać z pominięcia pewnej ważnej zmiennej (model jest zawsze uproszczeniem rzeczywistości), przybliżonego charakteru linii prostej jako zależności pomiędzy  $X$  a  $Y$  (prosta, ale nie do końca prosta) albo błędu pomiaru.

Współczynnik  $b_1$  (nachylenia prostej) określa wielkość efektu w przypadku regresji, tj. siły zależności pomiędzy zmiennymi.

Współczynnik  $b_1$  ma prostą interpretację: jeżeli wartość zmiennej  $X$  rośnie o jednostkę to wartość zmiennej  $Y$  zmienia się przeciętnie o  $b_1$  jednostek zmiennej  $Y$ . Wyraz wolny zwykle nie ma sensownej interpretacji (formalnie jest to wartość zmiennej  $Y$  dla  $X = 0$ ).

Oznaczmy przez  $y_i$  wartości obserwowane (zwane też empirycznymi) a przez  $\hat{y}_i$  wartości teoretyczne (leżące na prostej linii regresji).

Wartości  $b_0$  oraz  $b_1$  wyznacza się minimalizując sumę kwadratów odchyłeń wartości teoretycznych od wartości empirycznych, tj.:

$$(\hat{y}_1 - y_1)^2 + (\hat{y}_2 - y_2)^2 + \dots + (\hat{y}_n - y_n)^2 \rightarrow \min$$

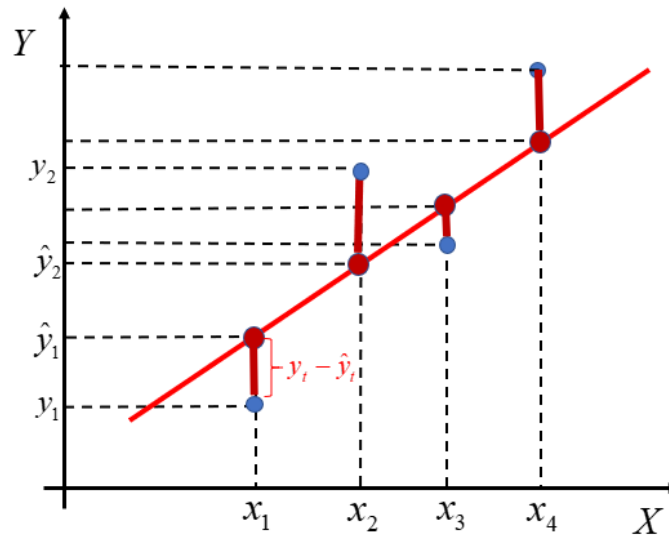
Rozwiązując powyższy **problem minimalizacyjny** otrzymujemy wzory określające wartości parametrów  $b_0$  oraz  $b_1$ . Metoda wyznaczania parametrów linii prostej w oparciu o minimalizację sumy kwadratów odchyłeń nosi nazwę **metoda najmniejszych kwadratów**.

Przypominamy, że **estymatorem** nazywamy metodę oszacowania parametru na podstawie próby. Ponieważ traktujemy  $b_0$  oraz  $b_1$  jako parametry pewnej populacji generalnej to „wzory na  $b_0$  oraz  $b_1$ ” statystyk nazwie estymatorami parametrów  $b_0$  oraz  $b_1$ .

Przypominamy dalej, że wartość średnia **dobrego estymatora** powinna wynosić zero (bo wtedy nie ma błędu systematycznego), oraz że wariancja estymatora powinna maleć wraz ze wzrostem liczebności próby. Można udowodnić że estymatory parametrów  $b_0$  oraz  $b_1$  uzyskane **metodą najmniejszych kwadratów** posiadają obie właściwości.

Graficznie **kryterium minimalizacyjne** przedstawia rysunek 4.4.

Suma podniesionych do kwadratu odległości pomiędzy czerwonymi (leżącymi na linii prostej w wersji czarno-białej) i niebieskimi kropkami ma być minimalna. Kropki



Rysunek 4.4: Metoda najmniejszych kwadratów

niebieskie to wartości empiryczne; kropki czerwone to wartości teoretyczne. Zadanie wyznaczenie parametrów takiej prostej oczywiście realizuje program komputerowy.

Można udowodnić, że bez względu czy punkty na wykresie układają się w przybliżeniu wzdłuż prostej czy nie, zawsze **jakaś prosta** zostanie dopasowana (jeżeli tylko punktów jest więcej niż jeden). Jak ocenić w sposób bardziej konkretny, a nie tylko na oko jakość dopasowania prostej do wartości empirycznych?

**Ocena dopasowania: wariancja resztowa oraz średni błąd szacunku**

Oznaczając *resztę* jako:  $e_i = y_i - \hat{y}_i$ , definiujemy **wariancję resztową** jako:

$$s_e^2 = \frac{e_1^2 + e_2^2 + \dots + e_n^2}{n - k}$$

Gdzie  $n$  oznacza liczbę obserwacji (liczebność próby), a  $k$  liczbę szacowanych parametrów bez wyrazu wolnego czyli jeden w regresji prostej (a więcej niż jeden w regresji wielorakiej o czym dalej).

Pierwiastek kwadratowy z **wariancji resztowej** nazywamy **średnim błędem szacunku** (*mean square error*, MSE).

**Ocena dopasowania: współczynniki zbieżności i determinacji**

Suma kwadratów reszt (albo odchyłeń wartości teoretycznych od wartości empirycznych, albo suma kwadratów błędów vel **resztowa suma kwadratów**):

$$RSK = (y_1 - \hat{y}_1)^2 + (y_2 - \hat{y}_2)^2 + \dots + (y_n - \hat{y}_n)^2$$

Suma kwadratów odchyłeń **wartości empirycznych** od średniej (**ogólna suma kwadratów**):

$$\text{OSK} = (y_1 - \bar{y})^2 + (y_2 - \bar{y})^2 + \dots + (y_n - \bar{y})^2$$

Suma kwadratów odchyleń **wartości teoretycznych** od średniej (**wyjaśniona suma kwadratów**):

$$\text{WSK} = (\hat{y}_1 - \bar{y})^2 + (\hat{y}_2 - \bar{y})^2 + \dots + (\hat{y}_n - \bar{y})^2$$

Można wykazać, że  $\text{OSK} = \text{WSK} + \text{RSK}$  zatem (po podzieleniu obu stron równania przez OSK) otrzymujemy:

$$1 = \text{WSK}/\text{OSK} + \text{RSK}/\text{OSK}$$

**Współczynnik zbieżności** oznaczany jako  $R^2$  to  $\text{WSK}/\text{OSK}$ .

**Współczynnik determinacji** oznaczany jako  $\Phi^2$  (duża grecka litera Fi) to  $\text{RSK}/\text{OSK}$ .

Współczynniki przyjmują wartość z przedziału  $[0, 1]$  lub  $[0, 100]\%$  jeżeli ich wartości zostaną pomnożone przez 100.

Interpretacja współczynnika zbieżności: udział (procent) zmienności wyjaśnianej przez linię regresji. Im  $R^2$  jest bliższe jedności (lub 100% jeżeli jest współczynnik zbieżności jest wyrażony w procentach) tym lepiej.

**Ocena dopasowania: istotność parametru  $b_1$**

Jeżeli:  $Y = 0 \cdot X + b_0$ , to  $Y = b_0$  czyli nie ma zależności pomiędzy  $X$  oraz  $Y$ . Wartości  $b_1$  bliskie zero wskazują na słabą zależność pomiędzy cechami.

Przypominamy, że **estymator** parametru  $b_1$  ma średnią równą prawdziwej wartości  $b_1$ . Dodatkowo zakładamy, że rozkład tego estymatora jest normalny, co pozwala wiarygodnie oszacować jego wariancję. W konsekwencji znamy jego dokładny rozkład, bo przypominamy, że rozkład normalny jest określony przez dwa parametry: średnią oraz wariancję (lub odchylenie standardowe).

Można teraz zadać pytanie jeżeli faktycznie  $b_1 = 0$ , to jakie jest prawdopodobieństwo, że współczynnik  $\hat{b}_1$  oszacowany na podstawie  $n$  obserwacji będzie (co do wartości bezwzględnej) większy niż  $b_e$ . Albo inaczej: otrzymaliśmy  $b_e$ , jakie jest prawdopodobieństwo otrzymania takiej wartości (lub większej co do wartości bezwzględnej) przy założeniu, że istotnie  $b_1 = 0$ .

Jeżeli takie prawdopodobieństwo jest duże, to uznajemy, że  $b_1 = 0$ , a jeżeli małe, to będziemy raczej sądzić, że  $b_1 \neq 0$ . Duże/małe przyjmujemy arbitralnie, zwykle jest to 0, 1, 0, 05 lub 0, 01. To prawdopodobieństwo to oczywiście poziom istotności.

W każdym programie komputerowym na wydruku wyników linii regresji są podane wartości prawdopodobieństwa  $b_1 > b_e$  (co do wartości bezwzględnej). Jeżeli jest ono mniejsze niż ustalony **poziom istotności** to  $b_1$  ma wartość istotnie różną od zera. Oprócz wartości prawdopodobieństwa drukowana jest także wartość błędu standardowego parametru zwykle oznaczana jako SE (*standard error*). Wartość SE nie jest wprawdzie potrzebna do oceny istotności (wystarczy prawdopodobieństwo), ale dobry zwyczaj nakazuje podawać także tę wartość w raporcie.

Testowanie istotności współczynnika regresji jest ważnym kryterium oceny jakości dopasowania. Regresja z **nieistotnym** współczynnikiem nie może być podstawą do interpretowania zależności pomiędzy  $X$  oraz  $Y$ .

### Waga a wzrost rugbyistów

Zależność między wagą (weight) a wzrostem (height):

$$\text{height} = b_0 + b_1 \text{weight}$$

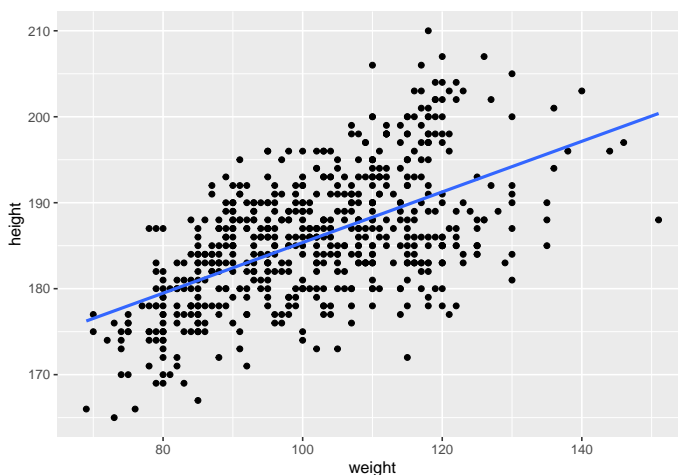
Oszacowanie tego równania na próbie 635 uczestników Pucharu Świata w rugby w 2023 roku daje następujące wyniki:

Zmienna	B	SE	z	p	CI95
(Intercept)	155,926	1,753	88,969	0	152,48 159,37
weight	0,294	0,017	17,305	0	0,26 0,33

Pierwsza kolumna Zmienna zawiera nazwy zmiennych ((Intercept) oznacza wyraz wolny). Druga kolumna oznaczona jako B zawiera oszacowane wartości (oceny) parametrów linii regresji. Kolumna SE zawiera oceny błędu standardowego estymatorów parametrów linii regresji. Kolumna p zawiera prawdopodobieństwo  $b > b_e$ .

Wzrost wagi zawodnika o 1kg skutkuje przeciętnie większym wzrostem o 0,294 cm. Współczynnik determinacji wynosi 32,86%. Współczynnik nachylenia prostej jest istotny ponieważ wartość p (tak mała, że w tabeli oznaczona jako 0) jest grubo poniżej zwyczajowego poziomu istotności ( $p < 0,05$ ).

Kolumna CI95 zawiera 95% przedziały ufności: z 95% prawdopodobieństwem wartość współczynnika nachylenia prostej znajduje się w przedziale 0.260–0.330.

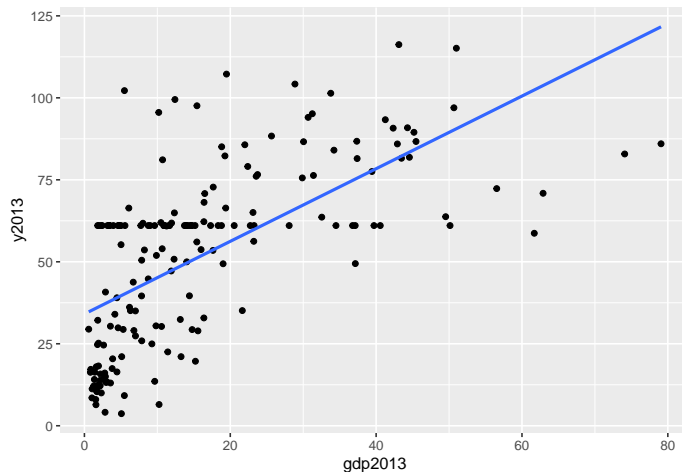


### Zamożność a konsumpcja mięsa

Następujące równanie opisuje zależność pomiędzy dochodem narodowym na głowę (tys USD *per capita*) a konsumpcją mięsa w kilogramach:

$$\text{konsumpcja} = b_0 + b_1 \text{gdp}$$

Model oszacowano dla krajów świata w roku 2013 na podstawie danych pobranych z bazy FAO Food Balance Sheet oraz Banku Światowego, otrzymując następujące wyniki:



Zmienna	B	SE	z	p	CI95
(Intercept)	34,085	2,232	15,268	0	29,68 38,5
gdp2013	1,108	0,100	11,124	0	0,91 1,3

Każdy tysiąc USD *per capita* więcej dochodu narodowego (GDP) oznacza przeciętny wzrost spożycia mięsa o 1,108 kg. Przeciętna różnica wartości teoretycznych od empirycznych wynosi 21,04 kg (średni błąd szacunku). Współczynnik zbieżności wynosi 40,88%. Współczynnik nachylenia prostej (którego wartość wynosi 1,108) jest statystycznie istotny.

Nie ma przykładów zastosowania regresji prostej w literaturze przedmiotu, bo jest ona zbyt dużym uproszczeniem rzeczywistości. Jest to jednak dobry punkt startu do bardziej skomplikowanego modelu **regresji wielorakiej**.

## 4.4 Zmienna liczbowa i zmienne liczbowe lub nominalne

### 4.4.1 Regresja wieloraka

Jeżeli zmiennych niezależnych jest więcej niż jedna, to mówimy o **regresji wielorakiej**. Przykładowo zależność pomiędzy wynikiem egzaminu, spożyciem kawy czasem nauki oraz predyspozycjami opisuje następujący model regresji:

$$\text{wynik} = b_0 + b_1 \cdot \text{kawa} + b_2 \cdot \text{czas} + b_3 \cdot \text{predyspozycje}$$

Współczynnik  $b_1$  określa wpływ spożycia kawy,  $b_2$  czasu poświęconego na naukę, a  $b_3$  predyspozycji (intelektualnych, mierzonych np. średnią oceną ze studiów). Ogólnie model regresji wielorakiej zapisać można jako:

$$Y = b_0 + b_1 \cdot X_1 + b_2 \cdot X_2 + \dots + b_k \cdot X_k$$

Wpływ każdej zmiennej  $X_i$  na zmienną zależną  $Y$  jest określony przez odpowiedni współczynnik  $b_i$ . Zmienne  $X_i$  mogą być zmiennymi liczbowymi lub nominalnymi.

Podobnie jak w przypadku regresji prostej do oceny stopnia dopasowania modelu do danych wykorzystuje się: średni błąd szacunku, współczynnik zbieżności  $R^2$  oraz weryfikuje się istotność współczynników  $b_i$ .

#### Standaryzacja współczynników regresji

Ponieważ współczynniki regresji  $b_1, \dots, b_k$  mogą być wyrażone w różnych jednostkach miary, bezpośrednie porównanie jest niemożliwe; mały współczynnik może w rzeczywistości być ważniejszy niż większy. Jeżeli chcemy porównywać wielkości współczynników to trzeba je **zestandaryzować**.

Standaryzowany współczynnik regresji dla  $i$ -tej zmiennej obliczony jest poprzez pomnożenie współczynnika regresji  $b_i$  przez  $s_{xi}$  i podzielenie przez  $s_y$ :

$$\beta_i = b_i \frac{s_{xi}}{s_y}$$

Dla przypomnienia  $s_{xi}$  to odchylenie standardowe zmiennej  $X_i$ , a  $s_y$  to odchylenie standardowe zmiennej  $Y$ . Interpretacja współczynnika standaryzowanego jest cokolwiek dziwna: zmiana zmiennej  $X_i$  o jedno odchylenie standardowe ( $s_{xi}$ ) skutkuje zmianą zmiennej  $Y$  o  $b_i$  jej odchylenia standardowego  $s_y$ . Na szczęście współczynniki regresji standaryzuje się nie w celu lepszej interpretacji, tylko w celu umożliwienia porównania ich względnej wielkości (*wielkości efektu*). W publikacjach medycznych zwykle używa się litery  $b$  na oznaczenie współczynników niestandaryzowanych a litery  $\beta$  na oznaczenie współczynników standaryzowanych.

#### Wielkość efektu

Współczynniki regresji to miara wielkości efektu, która wskazuje na siłę zależności między zmiennymi. Standaryzacja pozwala na porównanie wielkości efektu zmiennych mierzonych w różnych jednostkach miary. Standaryzacja przydaje się także w przypadku posługiwania się skalami pomiarowymi mierzącymi przekonania i postawy, które z definicji są bezjednostkowe.

#### Wybór zmiennych objaśniających

Zwykle jest tak, że do objaśniania kształtowania się wartości zmiennej  $Y$  kandyduje wiele potencjalnych predyktorów  $X_k$ . Model zawierający wszystkie  $X_k$  predyktory niekoniecznie będzie najlepszy. Nie wdając się w omawianie szczegółowych zasad porzeczamy na dwóch kryteriach:

1. Model prostszy jest lepszy od modelu bardziej skomplikowanego jeżeli adekwatnie objaśnia zmienność  $Y$  (zasada brzytwy Ockhama).
2. Model powinien zawierać tylko zmienne o współczynnikach, których wartości są statystycznie różne od zera.

Regresja krokowa (*stepwise regression*) jest metodą wyboru najlepszych predyktorów spośród większego zbioru zmiennych. Występuje w dwóch wariantach **dołączania** i **eliminacji**. Ponieważ **eliminacja** wydaje się prostsza omówimy tylko ten wariant.

W metodzie eliminacji początkowym modelem jest model zawierający wszystkie potencjalne  $X_k$  predyktory. Następnie testujemy istotność wszystkich współczynników regresji i usuwamy ze zbioru predyktorów ten, który jest „najbardziej nieistotny” (ma największą wartość  $p$ ). Procedurę powtarzamy dla modelu bez usuniętej zmiennej. Procedurę przerywamy gdy wszystkie współczynniki regresji są statystycznie istotne.

#### Zależność pomiędzy ciśnieniem skurczowym, BMI oraz wiekiem

$$\text{ciśnienie} = b_0 + b_1 \text{BMI} + b_2 \text{wiek}$$

Dane pochodzą z badania: Zależność pomiędzy BMI i wiekiem a występowaniem cukrzycy wśród dorosłych osób w Chinach. Badanie kohortowe (Chen i inni, *Association of body mass index and age with incident diabetes in Chinese adults: a population-based cohort study*. BMJ Open. 2018 Sep 28;8(9):e021768. doi: 10.1136/bmjopen-2018-021768. PMID: 30269064; PMCID: PMC6169758).

Oryginalny zbiór danych liczy 60 tysięcy obserwacji. Dla celów przykładu losowo wybrano 90, 490 oraz 4490 obserwacji. Zobaczmy jaki ma wpływ wielkość próby na wynik szacowania modelu.

Oszacowanie równania dla próby o wielkości 90 obserwacji daje następujące wyniki:

Zmienna	B	SE	z	p	Beta	CI95
(Intercept)	59,698	11,965	4,990	0,000	NA	35,92 83,48
BMI	1,742	0,486	3,583	0,001	0,33	0,78 2,71
age	0,484	0,124	3,906	0,000	0,36	0,24 0,73

Współczynnik zbieżności wynosi 26,24%. Kolumna Beta zawiera standaryzowane oceny parametrów regresji. Tej kolumny na poprzednich wydrukach (punkt 4.3.5) nie było, bo w przypadku regresji prostej standaryzacja jest zabiegiem raczej zbędnym. Dla wyrazu wolnego nie ma wartości standaryzowanej (co oznaczono jako NA czyli *not available*), ale to żadna strata – oceny tego parametru nie są interpretowane. Wpływ BMI na wielkość ciśnienia jest nieco niższy niż age.

Oszacowanie równania dla próby o wielkości 490 obserwacji daje następujące wyniki:

Zmienna	B	SE	z	p	Beta	CI95
(Intercept)	79,061	4,378	18,057	0	NA	70,46 87,66
BMI	1,213	0,183	6,637	0	0,28	0,85 1,57
age	0,259	0,053	4,856	0	0,21	0,15 0,36

Współczynnik zbieżności wynosi 14,97%. Wpływ BMI na wielkość ciśnienia jest teraz wyższy niż age. Przedziały ufności są węższe co wynika z większej liczebności próby.

Oszacowanie równania dla próby o wielkości 4490 obserwacji daje następujące wyniki:

Zmienna	B	SE	z	p	Beta	CI
(Intercept)	74,011	1,530	48,358	0	NA	71,01 77,01
BMI	1,375	0,064	21,404	0	0,30	1,25 1,50
age	0,320	0,018	18,270	0	0,25	0,29 0,35

Współczynnik zbieżności wynosi 18,54%. Przedziały ufności są jeszcze węższe. Ocena age z 95% prawdopodobieństwem znajduje się w przedziale [0.290 0.350] a w pierwszym oszacowaniu dla znacznie mniejszej próby było to [0.240 0.730]. Przedział jest ponad 8 razy węższy.

#### 4.4.2 Zmienne zero-jedynkowe

Zamiast (celem wykazania związku między zmienną liczbową a nominalną) porównywać średnie w grupach możemy wykorzystać metodę regresji wielorakiej. Zmienna nominalna jest zamieniana na jedną lub więcej zmiennych binarnych, które przyjmują tylko dwie wartości 0 lub 1.

Przykładowo rodzaj miejsca pracy (skala nominalna; dwie wartości: szpital, przychodnia) można zamienić na zmienną binarną praca przypisując 1 = szpital, oraz 0 = przychodnia (lub odwrotnie). Załóżmy że poziom stresu zależy od stażu pracy, satysfakcji (obie mierzone na skali liczbowej) i rodzaju miejsca pracy. Możemy to zapisać jako następujące równanie regresji:

$$\text{stres} = b_0 + b_1 \text{staż} + b_2 \text{satysfakcja} + b_3 \text{praca}$$

Jaka jest interpretacja współczynnika  $b_3$ ? Zakładając że 0 = przychodnia,  $b_3$  oznacza przeciętną zmianę wielkości stresu spowodowaną pracą w szpitalu w porównaniu do pracy w przychodni. Jeżeli ten współczynnik jest istotny statystycznie, to istnieje zależność pomiędzy stresem a miejscem pracy. Czyli zamiast stosować test  $t$  Welcha i porównywać średnie w grupach, możemy oszacować model regresji z wykorzystaniem stosownej zmiennej zero-jedynkowej a następnie sprawdzić czy współczynnik stojący przy tej zmiennej jest istotny.

Jeżeli zmienna nominalna ma  $n$  wartości należy ją zamienić na  $n - 1$  zmiennych zero-jedynkowych. Załóżmy że stress zależy także od wykształcenia, mierzonego w skali nominalnej (średnie, licencjat, magisterskie.) Tworzymy dwie zmienne: magister (jeden jeżeli respondent ma wykształcenie magisterskie lub 0 jeżeli nie ma) oraz licencjat (jeden jeżeli respondent ma licencjat lub 0 jeżeli nie ma). Równanie regresji ma postać:

$$\text{stres} = b_0 + b_1 \text{staż} + b_2 \text{satysfakcja} + b_3 \text{praca} + b_4 \text{magister} + b_5 \text{licencjat}$$

Jeżeli magister = 0 oraz licencjat = 0 to osoba ma wykształcenie średnie.

Interpretacja:  $b_4$  (jeżeli istotne) oznacza przeciętną zmianę wielkości stresu osoby z wykształceniem magisterskim w porównaniu do osoby z wykształceniem średnim.



Podobnie  $b_5$  oznacza przeciętną zmianę wielkości stresu osoby z wykształceniem licencjackim w porównaniu do osoby z wykształceniem średnim.

#### Zależność pomiędzy ciśnieniem skurczowym, BMI, wiekiem, płcią, paleniem i piciem

Poprzednio rozważany model rozszerzymy o trzy zmienne: płeć (kobieta/mężczyzna), status względem picia alkoholu (pije, pił, nigdy nie pił) oraz status względem palenia (palił, pali, nigdy nie palił). Zwróćmy uwagę że zmienne mierzące status względem palenia/picia mają nie dwie a trzy wartości. Należy każdą zamienić na dwie zmienne binarne, wg schematu:

CS (*current smoker*/pali) = 1 jeżeli pali, 0 w przeciwnym przypadku

PS (*past smoker*/palił) = 1 jeżeli palił ale nie pali, 0 w przeciwnym przypadku

CD (*current drinker*/pije) oraz PD (*past drinker*/pił) *per analogiam* do CS/CD

Zmienna płeć  $genderF$  = 1 jeżeli kobieta, lub 0 jeżeli mężczyzna. Zauważmy, że nazwa zmiennej dwuwartościowej wskazuje, która wartość jest zakodowana jako 1. Przykładowo  $genderF$  (*female* żeby się trzymać języka angielskiego) wskazuje że jedynką jest kobieta. Taka konwencja ułatwia interpretację. Gdybyśmy zamiast  $genderF$  nazwali zmienną  $gender$  to na pierwszy rzut oka nie było by wiadomo co zakodowano jako jeden. A tak wiadomo od razu jak interpretować parametr stojący przy tej zmiennej: zmiana wielkości ciśnienia u kobiet w porównaniu do mężczyzn.

Rozważany model ma postać:

$$SBP = b_0 + b_1BMI + b_2age + b_3genderF + b_4CS + b_5PS + b_6CD + b_7PD$$

Oszacowanie dla próby o wielkości 90 obserwacji daje następujące wyniki:

Zmienna	B	SE	z	p	Beta	CI
(Intercept)	90,332	15,745	5,737	0,000	NA	59,01 121,65
BMI	0,778	0,592	1,314	0,193	0,15	-0,40 1,96
age	0,441	0,121	3,658	0,000	0,33	0,20 0,68
genderF	-13,820	4,399	-3,141	0,002	-0,40	-22,57 -5,07
CS	-6,890	3,972	-1,735	0,087	-0,18	-14,79 1,01
PS	7,626	6,882	1,108	0,271	0,11	-6,07 21,32
CD	-3,959	8,523	-0,465	0,643	-0,04	-20,91 13,00
PD	-4,001	4,575	-0,875	0,384	-0,08	-13,10 5,10

Współczynnik zbieżności wynosi 38,11%. Tylko dwie na siedem zmiennych są istotne. Zwróćmy uwagę że nieistotnie zmienne mają przedziały ufności zawierające zero. W konsekwencji z 95% prawdopodobieństwem wartości tych współczynników mogą być raz ujemne raz dodatnie – nie mamy nawet pewności co do kierunku zależności między zmienną objaśnianą a ciśnieniem. Zmienne, które okazały się

istotne jednocześnie mają największą wielkość efektu (kolumna Beta) i nie jest to przypadek.

Wyniki oszacowania równania dla próby o wielkości 4490 obserwacji:

Zmienna	B	SE	z	p	Beta	CI
(Intercept)	80,089	1,623	49,334	0,000	NA	76,91 83,27
BMI	1,192	0,066	18,009	0,000	0,26	1,06 1,32
age	0,336	0,018	19,087	0,000	0,26	0,30 0,37
genderF	-5,329	0,508	-10,496	0,000	-0,16	-6,32 -4,33
CS	-2,752	0,583	-4,717	0,000	-0,07	-3,90 -1,61
PS	-2,021	1,045	-1,933	0,053	-0,03	-4,07 0,03
CD	3,621	1,535	2,360	0,018	0,03	0,61 6,63
PD	0,193	0,623	0,310	0,757	0,00	-1,03 1,41

Współczynnik zbieżności wynosi 20,72%. Zwiększenie liczebności próby z 90 do 4490 obserwacji spowodowało, że tylko dwie z siedmiu zmiennych mają nieistotne wartości. Analizując wartości standaryzowane możemy ustalić które zmienne mają największy wpływ na wielkość ciśnienia krwi.

Ktoś mógłby dojść do wniosku że wszystko da się **uistotnić** wystarczy zwiększyć wielkość próby. Teoretycznie tak, praktycznie nie. W praktyce nie interesuje nas niewielka wielkość efektu (znikomy wpływ czegoś na coś). Dodatkowo zebranie dużej próby może być kosztowne czyli w praktyce niemożliwe – nie mamy dość dużo pieniędzy. Można teoretycznie określić jaką wielkość próby pozwoli na ocenę jakiej wielkości efektu. Sposób postępowania jest wtedy następujący: określamy jaką wielkość efektu ma **znaczenie praktyczne** i na tej podstawie określamy liczebność próby. Takie zaawansowane podejście wykracza poza ramy tego podręcznika.

### Regresja krokowa

W modelu zależność pomiędzy ciśnienie skurczowym, BMI, wiekiem, płcią, paleniem i piciem (próba 4490) zmienne PD oraz PS są nieistotne przy czym współczynnik przy zmiennej PD ma wartość  $p$  równą 0,309 zaś przy zmiennej PS ma wartość 0,05324. Usuwamy zmienną PD (bo wartość  $p$  jest większa) i szacujemy równanie regresji dla sześciu pozostałych zmiennych. Otrzymujemy:

Zmienna	B	SE	z	p	Beta	CI
(Intercept)	80,108	1,622	49,387	0,000	NA	76,93 83,29
BMI	1,193	0,066	18,056	0,000	0,26	1,06 1,32
age	0,335	0,018	19,097	0,000	0,26	0,30 0,37
genderF	-5,358	0,499	-10,740	0,000	-0,16	-6,34 -4,38
CS	-2,740	0,582	-4,708	0,000	-0,07	-3,88 -1,60
PS	-1,980	1,037	-1,910	0,056	-0,03	-4,01 0,05
CD	3,579	1,528	2,342	0,019	0,03	0,58 6,58

Współczynnik przy zmiennej PS dalej uparcie jest istotny. Usuamy teraz tę zmienną. Otrzymujemy:

Zmienna	B	SE	z	p	Beta	CI
(Intercept)	79,865	1,618	49,375	0,00	NA	76,69 83,04
BMI	1,195	0,066	18,075	0,00	0,26	1,07 1,32
age	0,336	0,018	19,100	0,00	0,26	0,30 0,37
genderF	-5,155	0,488	-10,572	0,00	-0,16	-6,11 -4,20
CS	-2,540	0,573	-4,435	0,00	-0,06	-3,66 -1,42
CD	3,551	1,529	2,323	0,02	0,03	0,55 6,55

Wszystkie współczynniki mają istotnie różne od zera wartości. Wartość współczynnika zbieżności ostatecznego modelu wynosi 20,66%. Usuając nieistotne zmienne z modelu obniżyliśmy wartość współczynnika zmienności o 20,72% - 20,66% = 0,07%, czyli tyle co nic.

## 4.5 Przypadek specjalny: regresja logistyczna

Jeżeli zmienna  $Y$  jest zmienną **dwuwartościową**, czyli taką która przyjmuje tylko dwie wartości (np. chory/zdrowy), to metoda regresji nie może być zastosowana. Przykładowo jeżeli zakodujemy te wartości jako chory=0 i zdrowy=1, to zastosowanie regresji doprowadzi do obliczenia (teoretycznych) wartości  $Y$  różnych od 0 i 1. Taki wynik nie ma sensownej interpretacji.

Ale zamiast szacować regresję  $Y$  względem  $(X/X\text{-ów})$  można szacować regresję względem ryzyka dla  $Y$  (czyli prawdopodobieństwa że  $Y$  przyjmie wartość 1). Tutaj znowu pojawia się jednak trudność, bo ryzyko może przyjąć tylko wartości z przedziału  $[0, 1]$ . Nie wchodząc w matematyczne zawiłości, model zapisuje się jako (ln oznacza logarytm naturalny):

$$\ln\left(\frac{p}{1-p}\right) = b_0 + b_1 \cdot x_1 + \dots + b_k \cdot x_k$$

Zauważmy, że  $o = \frac{p}{1-p}$  to nic innego jak szansa (*odds*, por. punkt 4.1.1). Parametr  $b_i$  jest miarą wpływu zmiennej  $X_i$  na zmienną  $Y$ . Jeżeli  $X_i$  wzrośnie o jednostkę, to logarytm ilorazu szans wzrośnie o  $\ln(o)$  (przy założeniu, że pozostałe zmienne  $X$  mają pewne ustalone wartości a zmienia się tylko  $X_i$ ). Jeżeli  $X_i$  jest zmienną **dwuwartościową** to interpretacja jest jeszcze prostsza: jest to logarytm ilorazu szans dla wartości  $X_i = 1$  względem  $X_i = 0$ .

Zwykle zamiast **logarytmu ilorazu szans** wolimy interpretować zmianę w kategoriach **ilorazu szans**. Aby otrzymać ów iloraz należy wykonać następujące przekształcenie (exp oznacza podstawę logarytmu naturalnego):

$$o = \exp^{\ln(o)}$$

Zwykle iloraz szans wyraża się w procentach, czyli mnoży przez 100. Jeżeli ta liczba jest większa od 100 oznacza to wzrost szansy, a jeżeli mniejsza od 100, spadek szansy.

### Ocena dopasowania

Nie ma w przypadku regresji logistycznej możliwości obliczenia sumy kwadratów reszt (*residual sum of squares*) oraz współczynnika zbieżności. Model ocenia się używając jako kryterium dewiancję (*deviance*). Dewiancja to miara, której wielkość zależy od proporcji pomiędzy liczbą sukcesów obliczonych z modelu a liczbą sukcesów zaobserwowanych (jak dokładnie dewiancja jest liczona nie jest dla nas istotne).

Wyjaśnijmy to na przykładzie prostego modelu pomiędzy wystąpieniem osteoporozy a płcią. Model ma postać:

$$\ln(o) = b_0 + b_1 \text{płeć}$$

Po oszacowaniu  $b_0$  oraz  $b_1$  możemy łatwo obliczyć  $\ln(o)$ . Wiedząc że  $\ln(o) = \frac{p}{1-p}$  możemy stąd obliczyć prawdopodobieństwo, które jak widać będzie różne dla kobiet i mężczyzn. Po pomnożeniu tych prawdopodobieństw przez liczebności dostajemy (teoretyczne) liczebności sukcesów (tj. wystąpienia osteoporozy). Dewiancja będzie tym większa im różnica między tymi teoretycznymi liczebnościami a liczebnościami empirycznymi będzie większa.

Jako minimum porównuje się wielkość dewiancji szacowanego modelu z modelem zerowym (*null model*), tj. modelem w którym po prawej stronie równania występuje tylko stała:

$$\ln(o) = b_0$$

W tym modelu prawdopodobieństwo osteoporozy jest identyczne dla kobiet i mężczyzn, zatem w oczywisty sposób dewiancja tego modelu będzie większa. Pytanie jest czy różnica jest istotna statystycznie. Jeżeli jest większa to przyjmuje się, że szacowany model jest lepszy od modelu trywialnego (warunek minimum przydatności).

Jeżeli model zawiera wiele zmiennych w tym zmienne liczbowe, idea liczenia dewiancji jest podobna, ale oczywiście szczegóły są już bardziej skomplikowane. Szczegóły te nie są wszakże dla nas istotne bo zajmuje się tym program komputerowy.

**Minimalne kryteria oceny przydatności modelu regresji logistycznej:** istotnie mniejsza od modelu zerowego dewiancja oraz istotnie różne od zera parametry przy zmiennych niezależnych (predyktorach).

### Ocena skuteczności klasyfikacji

Model regresji logistycznej nie oblicza wartości zmiennej prognozowanej, bo ta nie jest liczbą, tylko **klasyfikuje**, tj. ustala (albo prognozuje) wartość zmiennej nominalnej w kategoriach „sukces”/„porażka”. Ważnym kryterium oceny jakości modelu jest ocena jakości klasyfikacji, to jest ocena na ile model poprawnie przypisuje przypadkom kategorii zmiennej prognozowanej. Im mniejsza rozbieżność pomiędzy wartościami rzeczywistymi, a prognozowanymi tym oczywiście lepiej.

Tę jakość klasyfikacji ocenia się za pomocą dwóch wskaźników: czułość (*sensitivity*) oraz swoistość (*specifity*).

1. Odsetek sukcesów zaklasyfikowanych jako „sukces” (**Czułość**); określane także jako TPR (*true-positive-rate*).
2. Odsetek porażek zaklasyfikowanych jako „porażka” (**Swoistość**); określane także jako TNR (*true-negative-rate*).

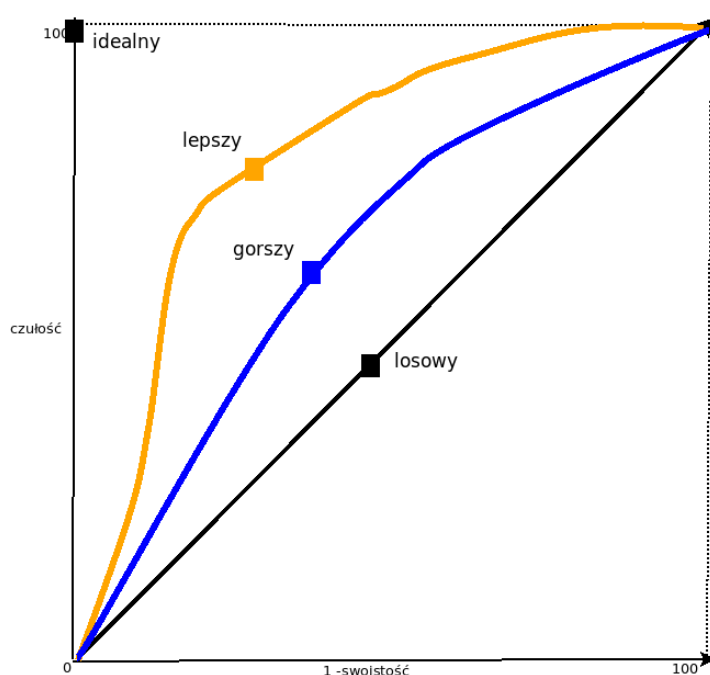
Klasyfikacja w modelu regresji logistycznej wygląda następująco. Jeżeli prawdopodobieństwo obliczone z modelu jest wyższe-lub-równe niż założona **wartość granic-**

na ( $p_g$ ), to zakładamy „sukces”, jeżeli tak nie jest, to zakładamy „porażkę”. Wartość graniczna jest ustalana albo arbitralnie albo na podstawie jakiejś dodatkowej (pozastatystycznej) informacji. Domyślnie za wartość graniczną przyjmuje się zwykle  $p_g = 0,5$ , co oznacza że wartości  $p \geq 0,5$  zostaną zamienione na „sukces” a wartości  $p < 0,5$  zostaną zamienione na „porażkę”.

#### Ocena dopasowania: krzywa ROC

Czułości oraz swoistości zależą od prawdopodobieństwa granicznego. Im wyższa jest wartość prawdopodobieństwa granicznego tym mniej będzie „sukcesów”.

Krzywa ROC przedstawia w układzie współrzędnych XY wartości czułości oraz swoistości dla różnych wartości granicznych. Współczynnik AUC (*area under curve*) to wielkość pola pod krzywą wyrażona w procentach pola kwadratu o boku 100%. AUC zawiera się w przedziale 50–100. Im większa wartość współczynnika tym lepiej. Model który klasyfikuje czysto losowo ma wartość AUC równą 50% (por. rysunek 4.5).



Rysunek 4.5: Krzywa ROC

#### Osteoporoza i witamina D

Al Zarooni A.A.R i inni badali wpływ różnych czynników na ryzyko wystąpienia osteoporozy (Risk factors for vitamin D deficiency in Abu Dhabi Emirati population; <https://doi.org/10.1371/journal.pone.0264064>), takich jak deficyt witaminy D, wiek oraz płeć w grupie 392 osób.

Zacznijmy od modelu zerowego tj. takiego w którym ryzyko/prawdopodobieństwo/szansa wystąpienia osteoporozy jest takie same bez względu na wielkości

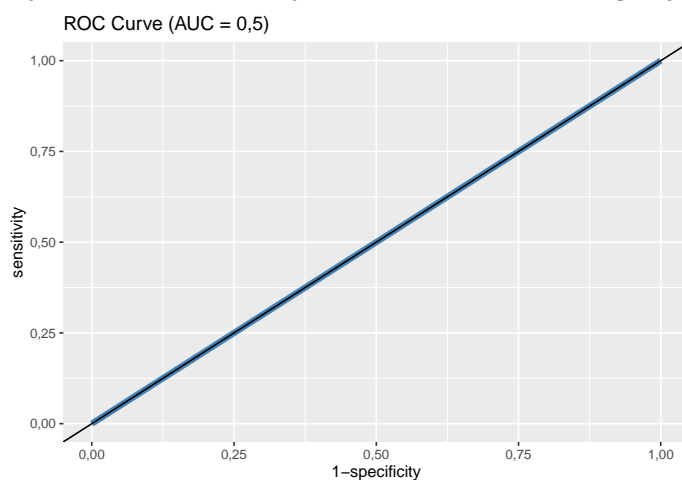
innych zmiennych. Odpowiada to następującemu równaniu:

$$\ln(o) = b_0$$

W tabeli zestawiono wartości parametrów oszacowanego modelu, ilorazy szans, przedziału ufności oraz prawdopodobieństwo.

Parametr	Ocena	SE	z	p
(Intercept)	-2,644537	0,2029618	-13,02973	0

Można obliczyć że (teoretyczne) prawdopodobieństwo wystąpienia osteoporozy wyniosło 0,0663265. Krzywa ROC dla modelu zerowego wygląda następująco:



Model zerowy jak sama nazwa wskazuje może tylko służyć do porównania z bardziej skomplikowanymi modelami.

Takim bardziej skomplikowanym modelem będzie przykładowo zależność pomiędzy wystąpieniem osteoporozy a płcią, którą można opisać następującym równaniem regresji:

$$\ln(o) = b_0 + b_1 \text{kobieta}$$

Zmienna *kobieta* przyjmuje wartość 1 jeżeli osoba była kobietą oraz zero w przypadku jeżeli była mężczyzną. Dla przypomnienia *o* jest szansą wystąpienia osteoporozy.

W tabeli zestawiono wartości parametrów oszacowanego modelu, ilorazy szans, przedziału ufności oraz prawdopodobieństwo.

Parametr	Ocena	SE	z	p	OR	CI
(Intercept)	-3,367	0,455	-7,403	0,000	0,03	0,01 0,08
genderF	1,014	0,509	1,992	0,046	2,76	1,09 8,40

Znając wartości współczynników równania można obliczyć wartości  $\ln(o)$ .

Dewiancja modelu jest istotnie mniejsza od modelu zerowego (wartość  $p$  wynosi bowiem 0,0303521).

Zależność pomiędzy wystąpieniem osteoporozy a płcią, wiekiem oraz poziomem witaminy D można opisać następującym równaniem regresji:

$$\ln(o) = b_0 + b_1 \text{kobieta} + b_2 \text{wiek} + b_3 \text{poziomD}$$

W tabeli zestawiono wartości parametrów oszacowanego modelu, ilorazy szans, przedziału ufności oraz prawdopodobieństwo.

Parametr	Ocena	SE	z	p	OR	CI
(Intercept)	-12,183	1,766	-6,898	0,000	0,00	0,00 0,00
d	0,005	0,009	0,536	0,592	1,00	0,99 1,02
age	0,156	0,026	5,930	0,000	1,17	1,12 1,24
genderF	2,463	0,662	3,722	0,000	11,74	3,54 48,76

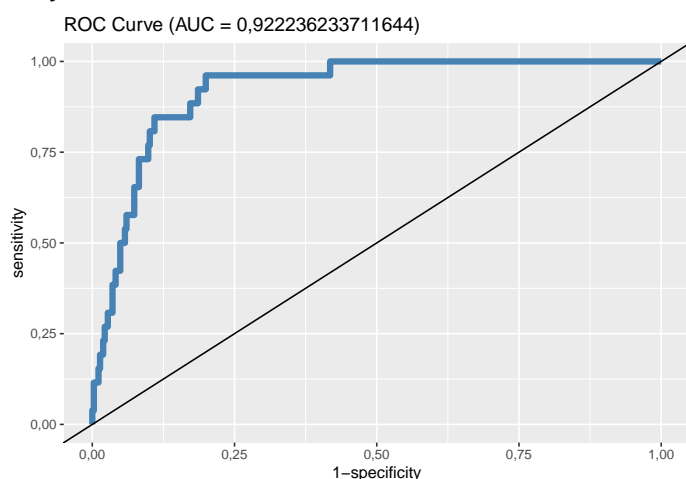
Macierz pomyłek (*confussion matrix*)

```
##          Osteoporoza
## Prognoza  0    1
##          0 362  22
##          1   4   4
```

Stąd: czułość 0,1538462; swoistość 0,989071.

Istotność modelu: dewiancja jest istotnie mniejsza od dewiancji modelu zerowego ( $p = 0$ ).

Krzywa ROC



## 4.6 Przypadek specjalny: co najmniej dwie zmienne porządkowe

### 4.6.1 Pomiar siły zależności: współczynnik korelacji rang

Współczynnik korelacji rang Spearmana (*Spearman's Rank-Order Correlation*) może być stosowany w przypadku gdy cechy są mierzone w skali porządkowej (lub lepszej czyli liczbowej).

Obliczenie współczynnika Spearmana dla  $N$  obserwacji na zmiennych  $XY$  polega na zamianie wartości zmiennych  $X$  oraz  $Y$  na **rangi** (numery porządkowe od 1 do  $N$ ). Następnie stosowana jest formuła współczynnika korelacji liniowej Pearsona ( $\tau_x$  oraz  $\tau_y$  oznaczają **rangi**):

$$\rho_{xy} = \frac{\text{cov}(\tau_x, \tau_y)}{s_{\tau_x} s_{\tau_y}}$$

Współczynnik  $\rho_{xy}$  to – podobnie jak **oryginalny** współczynnik korelacji liniowej Pearsona – miara niemianowana, o wartościach ze zbioru  $[-1;1]$ .

#### Spożycie mięsa

Współczynnik Pearsona i Spearmana dla zależności między spożyciem mięsa w 1980 a spożyciem mięsa w 2013 roku (zmienna objaśniana):

- współczynnik Pearsona: 0,68;
- współczynnik Spearmana: 0,68.

Nie ma sensu liczenia współczynnika korelacji rang w przypadku kiedy obie cechy są liczbami, bo wtedy należy użyć normalnego współczynnika Pearsona. Ale nie jest to też błędem więc w powyższym przykładzie go liczymy.

Współczynnik korelacji liniowej Spearmana wynosi 0,68 (umiarkowana korelacja).

Czy ta wartość jest istotnie różna od zera? Jest na to stosowny test statystyczny, który sprowadza się do określenia jakie jest prawdopodobieństwo otrzymania  $r_s = 0,68$  przy założeniu że prawdziwa wartość  $r_s$  wynosi zero. Otóż w naszym przykładzie to prawdopodobieństwo wynosi  $2.302116e-26$  (czyli jest ekstremalnie małe – współczynnik jest istotnie różny od zera).

## 4.7 Podsumowanie

Przedstawiono 7 następujących metod ustalania zależności między zmiennymi:

1. Wykres rozrzutu.
2. Tablica wielodzielna i test chi-kwadrat.
3. Współczynnik korelacji liniowej Pearsona.
4. Współczynnik korelacji Spearmana.
5. Regresja liniowa.
6. Regresja logistyczna.
7. testy  $t$  Welcha, U Manna-Whitneya, ANOVA albo test Kruskala-Wallisa.



## Rozdział 5

# Przykłady badań ankietowych

Uwaga: ankieta nie jest kolejną metodą statystyczną tylko techniką zbierania danych. Wszystkie metody już zostały przedstawione i żadnej nowej nie będzie.

## 5.1 Jak zacząć badanie?

Każde badanie, w tym ankietowe.

Należy zastanowić się nad trzema sprawami:

1. Co chcemy ustalić?
2. Jakie dane są nam potrzebne, żeby ustalić to co chcemy ustalić?
3. Jak te dane zebrać (czyli co i w jaki sposób zmierzyć)?

### Co chcemy ustalić?

Najlepiej jakąś zależność. Na przykład: stress a wypalenie zawodowe; satysfakcja zawodowa a retencja; determinanty satysfakcji zawodowej.

Może być od biedy opis czegoś lub porównanie czegoś z czymś. Przykłady: nadwaga wśród studentów wydziału zdrowia PSW; analiza porównawcza wypalenia zawodowego pielęgniarek pracujących w różnych systemach opieki.

### Co i jak mierzyć?

Jeżeli mamy zamiar badać nadwagę, to powinniśmy zmierzyć masę ciała. Jeżeli celem jest ustalenie zależności pomiędzy stresem a wypaleniem zawodowym to niewątpliwie powinniśmy zmierzyć stress i wypalenia. Jak dotąd banalnie prosto. Problem zaczyna się w momencie odpowiedzi na pytanie **jak**?

### 5.1.1 Mierzenie twardych faktów vs mierzenia przekonań

Możemy pytać w ankiecie o dwie rzeczy:

- **fakty** (wiek, staż, zawód, tętno, przebyte choroby);
- **przekonania, wartości, postawy; uczucia** (strach / radość) albo **zamiary** (w języku attitudes/emotions/intentions).

Mierzenie **faktów** nie wymaga dodatkowych objaśnień. Problem jest z mierzeniem **przekonań**.

**Przekonanie** to idea, którą jednostka uważa za prawdziwą. **Wartości** to trwałe przekonania o tym, co jest ważne dla jednostki. Stają się standardami, według których jednostki dokonują wyborów. **Postawy** to mentalne dyspozycje/nastawienie przed podjęciem decyzji, które skutkują określonym zachowaniem (zrobię to a nie tamto). Postawy kształtowane są wartościami i przekonaniami.

### 5.1.2 Pomiar przekonań, wartości i postaw

Postawy/uczucia/zamiary są to pojęcia abstrakcyjne. Często (albo zawsze) definiowane w obszarze psychologii, nauk o zarządzaniu itp.

Pomiar *przekonań* jest dokonywany w specyficzny sposób. **Definicja konceptualna** definiuje pojęcie (zaufanie do kogoś/czegoś to **przekonanie**, że *działania*

*tego kogoś/czegoś okazać się zgodne z naszymi oczekiwaniami; satysfakcja to uczucie przyjemności, zadowolenia z czegoś; samoskuteczność to przekonanie, iż jest się w stanie zrealizować określone działanie lub osiągnąć wyznaczone cele).* **Definicja operacyjna** określa jak zmierzyć pojęcie (jak zmierzyć satysfakcję) Przejście od definicji konceptualnej do definicji operacyjnej bywa czasami mocno, hmm... arbitralne.

### 5.1.3 Skala Likerta

Przykładowo chcemy się dowiedzieć czy i jak bardzo respondenci boją się COVID19.

W najprostszej wersji się po prostu pytamy: **Czy pan/pani boi się COVID19?** i dajemy respondentowi trzy możliwe warianty odpowiedzi: Tak/Nie/Nie wiem. Taki pomiar jest mocno zgrubny: ktoś się może bać panicznie a ktoś inny dużo mniej.

Subtelniejszy pomiar to wybór spośród pięciu wariantów: bardzo się boję – boję się – trudno powiedzieć – nie boję się – zupełnie się nie boję.

Taką skalę pomiarową określamy jak wiemy jako **porządkową**. Pomiar nie są liczbami, ale są uporządkowane. Rangi wartości są już liczbami (np. 1–5 w przykładowej skali pięciowariantowej), można je np. uśredniać. Tego typu skala pomiarowa, typowa dla ankiet, nosi nazwę skali **Likerta**. Można sobie wymyślać skalę Likerta 7-punktową i więcej.

Naszym zdaniem powyżej 7 wariantów normalny respondent będzie miał problem czy się bardziej-bardziej czy jednak bardziej-bardziej-bardziej boi.

### 5.1.4 Skala pomiarowa czyli inwentarz albo kwestionariusz

Ponieważ skala Likerta, mimo że lepsza od pytań tak/nie, jest ciągle zgrubna, to uważa się powszechnie, że lepszy wynik da pomiar wielokrotny. W naukach podstawowych mierzymy (np. liniijką) parę razy, a wynik uśredniamy co daje pomiar bardziej precyzyjny. Tutaj pytamy się parę razy o to samo co ma dać podobny efekt (mniejszy średni błąd pomiaru). Taka seria pytań nosi też nazwę skali albo **inwentarza**. Nie pytamy się zatem **Czy pan/pani boi się COVID19?** tylko zadajemy serię pytań o strach względem COVID19:

#### Strach przed COVID19

The Fear of COVID-19 Scale: Development and Initial Validation. International Journal of Mental Health and Addiction, 1–9. (<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7100496/>).

Lęk przed koronawirusem COVID-19 i lęk przed śmiercią – polskie adaptacje narzędzi (<https://www.termedia.pl/Fear-of-COVID-19-and-death-anxiety-Polish-adaptations-of-scales,116,44937,1,1.html>).

1. I am most afraid of Corona.
2. It makes me uncomfortable to think about Corona.
3. My hands become clammy when I think about Corona.
4. I am afraid of losing my life because of Corona.
5. When I watch news and stories about Corona on social media, I become

nervous or anxious.

6. I cannot sleep because I'm worrying about getting Corona.
7. My heart races or palpitates when I think about getting Corona.

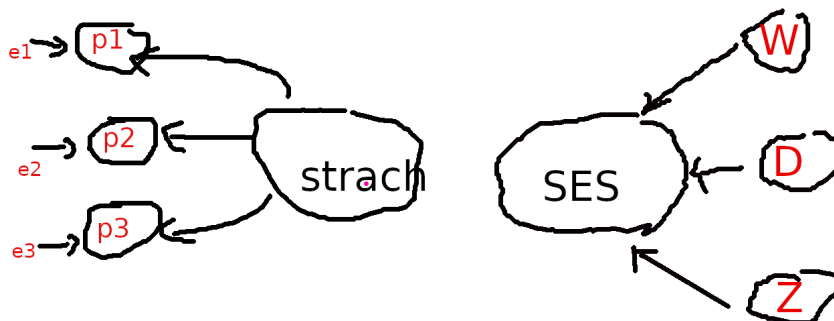
albo:

1. Boję się koronawirusa.
2. Czuję dyskomfort, gdy myślę o koronawirusie.
3. Pocą mi się dłonie, gdy myślę o koronawirusie.
4. Boję się, że mogę stracić życie z powodu koronawirusa.
5. Gdy oglądam wiadomości i czytam o koronawirusie w mediach społecznościowych, robię się nerwowy i niespokojny.
6. Nie mogę spać, ponieważ martwię się, że ja lub moi bliscy zarażą się.
7. Dostaję palpitacji serca, gdy myślę o tym, że mógłbym się zarazić.

Odpowiadający ma do wyboru pięć wariantów odpowiedzi: **zdecydowanie nie/nie/nie mam zdania/tak/zdecydowanie tak**.

### 5.1.5 Model pomiaru

**Ukryty czynnik** (strach) kształtuje wartości **indykatorów** (odpowiedzi na pytania) Taki sposób pomiaru **ukrytego czynnika** (*latent* w języku angielskim) określa się mianem refleksyjnego (co jest kalką od *reflexive*). Na rysunku 5.1 kierunek strzałki obrazuje zależność (czynnik → indyktor).



Rysunek 5.1: Modele pomiaru czynnika ukrytego

Alternatywny sposób definiowania ukrytego (w pewnym sensie, raczej złożonego) czynnika nosi nazwę **formatywnego** (albo indeksu): czynnik jest sumą indykatorów. Przykładem może być SES: status socjo-ekonomiczny będący agregatem wykształcenia (W), dochodu (D) oraz zawodu (Z).

W założeniu indykatory są jednakowo dobrymi miarami czynnika refleksyjnego i jako takie powinny być mocno skorelowane (mierzą to samo). Natomiast składniki czynnika formatywnego nie powinny być skorelowane, raczej każdy powinien mierzyć

**inny aspekt** czynnika. Ktoś może być profesorem za przeproszeniem filozofii, nie mieć pracy i kiepskie dochody. Tylko jeden z trzech aspektów podwyższa mu SES; albo świetnie zarabiająca prostytutka bez matury.

Jeżeli w czynniku refleksyjnym pominiemy jeden z trzech indyktorów, to nic się nie stanie oprócz tego, że pomiar będzie mniej precyzyjny. Jeżeli w czynniku formatywnym pominiemy indyktor, to popełniamy grubo błąd, bo pomijamy istotny składnik całości.

Dobłą wiadomością jest, że najprostszy sposób pomiaru traktuje czynniki refleksyjne i formatywne jednakowo: wartością czynnika jest suma albo ewentualnie średnia wartość indyktorów. Jeżeli indykatory są mierzone za pomocą skali Likerta suma (albo średnia) wartości rang po prostu. W skali strachu przed COVID ten kto się najbardziej boi powinien odpowiedzieć 7 razy **zdecydowanie tak**, co odpowiada sumie 35 rang (jeżeli rangujemy od 1 do 5). Ten który się wcale nie boi zaś będzie miał 7.

Małym utrudnieniem mogą być **pytania odwrócone**. Jeżeli pytamy o strach przed COVID i w każdym pytaniu jak bardzo ktoś się boi, albo jak bardzo mu serce bije, ale w jednym z pytań zapytamy **nie boję się COVID**, to ranga 5 odpowiada uczuciu **braku strachu**. Rangi w pytaniach odwróconych należy przeliczyć (odwrócić): 1 zamienić na 5, 2 na 4 itd... Jeżeli używamy cudzych skal to w opisie powinno być wskazane, które pytania są odwrócone.

**Zalecany schemat postępowania jeżeli w ankiecie mają być mierzone przekonania** (strach, samoskuteczność, wypalenie zawodowe, stress czy satysfakcja):

- Doksztalcamy się nieco z psychologii mimo wszystko;
- Robimy przegląd literatury i znajdujemy skalę, którą ktoś już wymyślił żeby zmierzyć to co my chcemy zmierzyć, bo **raczej nie należy wymyślać własnych skal**;
- Robimy ankietę (w Internecie) i zbieramy dane;
- Wykonujemy analizę statystyczną.

Banalnie proste co udowodnimy na przykładach.

## 5.2 Wiedza na temat szkodliwości palenia i jej uwarunkowania wśród studentów PSW

### 5.2.1 Cel

Celem jest ocena wielkości zjawiska palenia tytoniu oraz poziom wiedzy na temat szkodliwości palenia tytoniu wśród studentów PSW oraz zweryfikowanie wpływu wybranych czynników warunkujących ten nałóg.

**Postawiono następujące hipotezy badawcze:**

1. Jaka jest wielkość zjawiska palenie tytoniu wśród studentów PSW?
2. Jaka jest wiedza na temat szkodliwości palenia tytoniu wśród studentów PSW?
3. Czy palenie jest skorelowane z płcią, stażem pracy i miejscem pracy?
4. Czy wiedza na temat szkodliwości palenie jest skorelowana z płcią, stażem pracy i miejscem pracy?
5. Czy palenie jest skorelowane z wiedzą na temat szkodliwości palenia?

### 5.2.2 Metoda

Badanie ankietowe wśród studentów Ratownictwa Medycznego (RM) oraz Pielęgniarstwa (PO) przeprowadzono w styczniu 2023. Ankieta zawierała pytania dotyczące palenia tytoniu (pali/nie pali/palił, jak długo pali itd), test wiedzy na temat szkodliwości palenia oraz pytania o rodzaj miejsca pracy, staż pracy i płeć itd.

Pięć następujących pytań oceniało wiedzę ankietowanego na temat szkodliwości palenia:

- Czy bardziej szkodliwe dla zdrowia jest czynne czy bierne palenie? (JW);
- Jakie według Ciebie choroby układu oddechowego mogą być spowodowane bezpośrednio przez palenie papierosów? (WW);
- Czy palenie papierosów powoduje choroby układu pokarmowego? (JW);
- Jakie według Ciebie choroby kardiologiczne mogą być spowodowane bezpośrednio przez palenie papierosów? (WW);
- Jaki według Ciebie ma wpływ palenie papierosów na narządy zmysłów? (WW).

W przypadku pytań jednokrotnego wyboru (JW), za wskazanie poprawnej odpowiedzi respondent otrzymywał 1 punkt. W przypadku pytań wielokrotnego wyboru (WW) za wskazanie prawidłowej odpowiedzi respondent otrzymywał 1 punkt, ale za wskazanie nieprawidłowej otrzymywał (minus) -1 punkt (aby nie opłacała się strategia zaznaczenia wszystkich odpowiedzi). Maksymalna możliwa do uzyskania liczba punktów wynosiła 19.

### 5.2.3 Zastosowane metody statystyczne

- Hipotezę 1 weryfikowano na podstawie wielkości odsetka respondentów palących.
- Hipotezę 2 weryfikowano na podstawie wielkości odsetka respondentów wykazujących się dobrą i bardzo dobrą wiedzą na temat palenia.
- Hipotezy 3–5 zweryfikowano z wykorzystaniem tablic wielodzielnych/testu chi-kwadrat oraz porównania średniego poziomu depresji w grupach za pomocą testów Manna-Whitneya oraz Kruskala-Wallisa.

### 5.2.4 Metryczka (analiza respondentów)

Rozkład ankietowanych wg wyniku testu na znajomość szkodliwości palenia przedstawiono na rysunku 5.2.

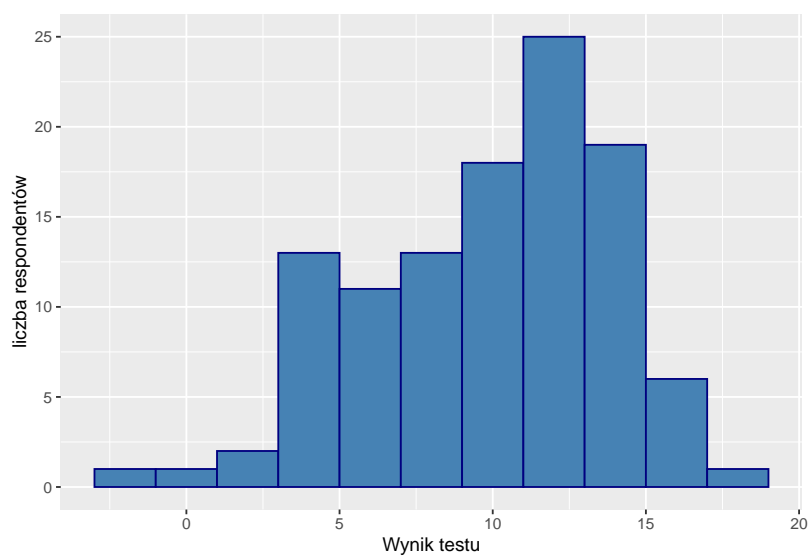
W badaniu wzięło udział 110 studentów. Otrzymano 110 poprawnie wypełnionych ankiet. Średnia wartość testu oceniającego wiedzę wyniosła 10,364 (odchylenie standardowe 3,997).

Rozkład ankietowanych ze względu na status względem palenia przedstawiono na rysunku 5.3.

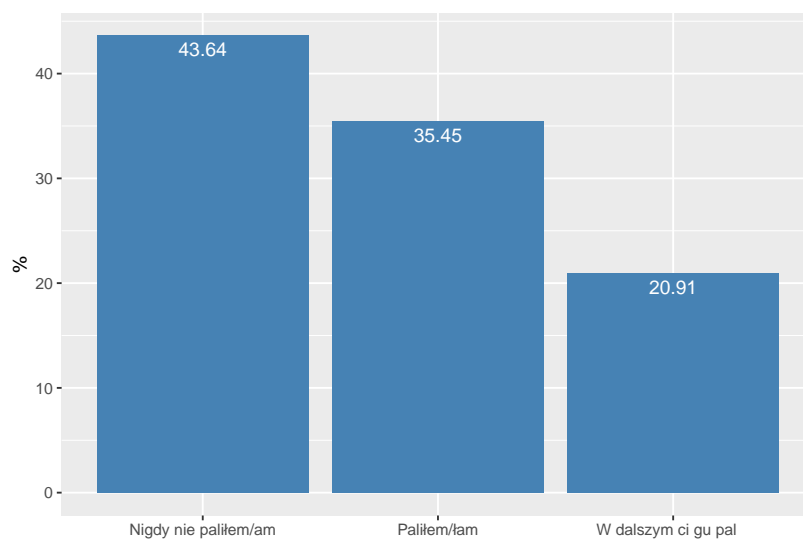
Rozkład ankietowanych ze względu na staż pracy przedstawiono na rysunku 5.4.

Zmienna płeć i miejsce pracy są dwuwartościowe, więc (oszczędzając lasy) proponujemy zamiast wykresu poprzestanie na podaniu odsetka kobiet/mężczyzn oraz pracowników szpitali/przychodni.

Zatem wśród respondentów było 70% kobiet i 30% mężczyzn. Pracę w szpitalu zadeklarowało 90% respondentów a w przychodni 10% respondentów.



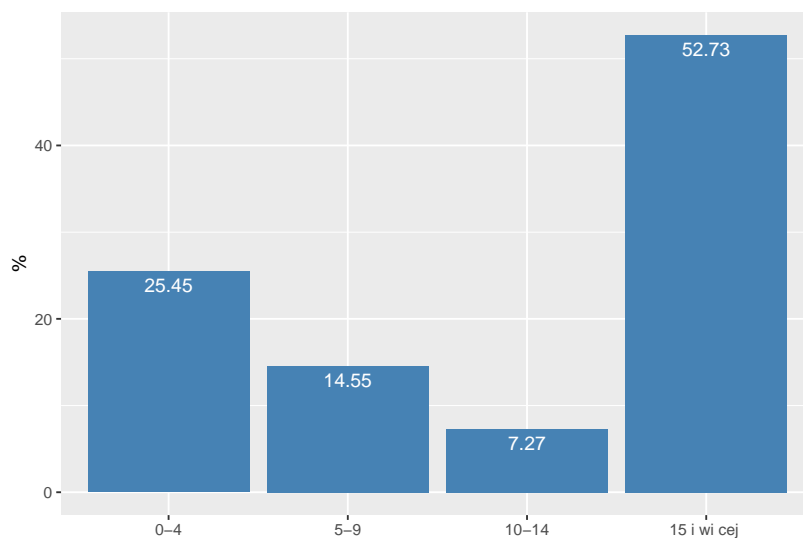
Rysunek 5.2: Studenci wg wyniku testu na znajomość szkodliwości palenia



Rysunek 5.3: Studenci wg statusu palenia (%)

### 5.2.5 Weryfikacja hipotezy 1

Palą lub paliło 62 respondentów (56.36 %). Żeby stwierdzić czy to jest dużo czy mało to na przykład można by porównać z jakąś średnią ogólnopolską.



Rysunek 5.4: Studenci wg stażu pracy (%)

### 5.2.6 Weryfikacja hipotezy 2

Średnia wartość uzyskana w teście wyniosła 10,36 (mediana 11); 3/4 respondentów nie uzyskało więcej niż 13 (czyli 68,4 %).

### 5.2.7 Weryfikacja hipotez 3-5

Czy palenie jest skorelowane z płcią?

	K	M
Nigdy nie paliłem/am	31	17
Paliłem/łam	27	12
W dalszym ciągu palę	19	4

```
##
```

```
## Pearson's Chi-squared test
```

```
##
```

```
## data: t.sex.f
```

```
## X-squared = 2,4228, df = 2, p-value = 0,2978
```

Nie jest o czym świadczy wysoka wartość p (0,2978).

Czy palenie jest skorelowane ze stażem pracy?

	0-4	5-9	10-14	15 i więcej
Nigdy nie paliłem/am	14	6	4	24
Paliłem/łam	8	6	2	23
W dalszym ciągu palę	6	4	2	11

```
##
## Pearson's Chi-squared test
##
## data:  t.staz.f
## X-squared = 1,7687, df = 6, p-value = 0,9397
```

Nie jest o czym świadczy wysoka wartość p (0,9397).  
Czy palenie jest skorelowane z miejscem pracy?

	Przychodnia	Szpital
Nigdy nie paliłem/am	5	43
Paliłem/łam	3	36
W dalszym ciągu palę	3	20

```
##
## Pearson's Chi-squared test
##
## data:  t.praca.f
## X-squared = 0,47674, df = 2, p-value = 0,7879
```

Nie jest o czym świadczy wysoka wartość p (0,7879).  
Czy wiedza na temat palenia jest skorelowana z płcią?

płeć	średnia	mediana	n
K	10,831169	12	77
M	9,272727	10	33

Sprawdzamy czy rozkłady są normalne:

płeć	S-W	p
K	0,9325443	0,0005147
M	0,9295526	0,0339828

Nie są. Należy zastosować test U Manna-Whitneya:

```
## $p.value
## [1] 0,04883299
```

Wartość p wynosi 0,048833, zatem hipotezę o braku zależności na poziomie 0,05 należy odrzucić. Istnieje zależność pomiędzy wiedzą na temat palenia a płcią. Kobiety wykazują się wyższą wiedzą na temat szkodliwości palenia od mężczyzn.

Czy wiedza na temat palenia jest skorelowana z miejscem pracy?

m.pracy	średnia	mediana	n
Przychodnia	10,36364	10	11
Szpital	10,36364	11	99

Sprawdzamy czy rozkłady są normalne:



m.pracy	S-W	p
Przychodnia	0,9272191	0,3834099
Szpital	0,9590504	0,0036563

Nie są. Należy zastosować test U Manna-Whitneya:

```
## [1] 0,8026325
```

Wartość p wynosi 0,8026325 – nie ma podstaw od odrzucenia hipotezy o identyczności rozkładów na poziomie  $p = 0,05$ . Nie ma różnicy w poziomie wiedzy wśród pracowników przychodni i szpitali.

Czy wiedza na temat palenia jest skorelowana ze stażem?

staż	średnia	mediana	n
0-4	9,928571	10,0	28
5-9	9,812500	11,0	16
10-14	9,875000	11,0	8
15 i więcej	10,793103	11,5	58

Sprawdzamy czy rozkłady są normalne:

staż	S-W	p
0-4	0,9644460	0,4418353
5-9	0,8847339	0,0460152
10-14	0,8985547	0,2804097
15 i więcej	0,9409402	0,0071265

Nie są. Należy zastosować test Kruskala-Wallisa:

```
## [1] 0,6771844
```

Wartość p wynosi 0,6771844 – nie ma podstaw od odrzucenia hipotezy o identyczności rozkładów na poziomie  $p = 0,05$ . Nie ma różnicy w poziomie wiedzy wśród pracowników o różnym stażu pracy.

Czy wiedza o szkodliwości palenia jest skorelowana ze statusem względem palenia? Chcemy zastosować tablicę wielodzielną/test chi kwadrat. Musimy zatem zamienić skalę liczbową zmiennej mierzącej wiedzę nt szkodliwości palenia na nominalną, np tak: 0–5 mała; 6–10 średnia; 11–15 duża, 16–19 ogromna:

	Nigdy nie paliłem/am	Paliłem/łam	W dalszym ciągu palę
duża	22	16	14
mała	7	6	4
ogromna	2	3	2
średnia	17	14	3

```
##
## Pearson's Chi-squared test
##
## data:  wiedza.status.t
## X-squared = 4,9954, df = 6, p-value = 0,5444
```

Wiedza i status względem palenia nie są skorelowane o czym świadczy wysoka wartość  $p$  (0,5444).

Można to samo zweryfikować stosując test Kruskala-Wallisa:

status	średnia	mediana	n
Nigdy nie paliłem/am	10,31250	10,5	48
Paliłem/łam	10,07692	10,0	39
W dalszym ciągu palę	10,95652	12,0	23

## [1] 0,7787663

Wynik jest identyczny (wysoka wartość  $p$  0,7788).

### 5.2.8 Wnioski

- Ponad połowa studentów pali lub paliła.
- Istnieje zależność pomiędzy wiedzą o szkodliwości palenia a płcią.
- Nie ma związku pomiędzy statusem względem palenia a płcią, miejscem pracy i stażem.
- Nie ma związku pomiędzy wiedzą o szkodliwości palenia a miejscem pracy i stażem.

## 5.3 Depresja i jej uwarunkowania wśród studentów PSW

### 5.3.1 Cel

Celem jest ustalenie czy depresja jest istotnym problemem wśród studentów PSW oraz zweryfikowanie wybranych czynników warunkujących depresję.

### 5.3.2 Metoda

Badanie ankietowe wśród studentów Ratownictwa Medycznego (RM) oraz Pielęgniarstwa (PO) przeprowadzono w styczniu 2023. Ankieta zawierała test samooceny depresji Becka oraz pytania o rodzaj miejsca pracy, staż pracy i płeć.

Test samooceny depresji Becka składa się z 21 pytań. W każdym pytaniu możliwe są 4 warianty odpowiedzi, odpowiadające zwiększonej intensywności objawów depresji, którym w związku z tym przypisuje się wartości od zera do 3 punktów. Maksymalna liczba punktów w teście wynosi 63 a minimalna 0.

Interpretacja wyników testu Becka: 0–19 brak/łagodna depresja; 20–25 umiarkowana; 26–63 ciężka depresja.

Postawiono następujące hipotezy badawcze:

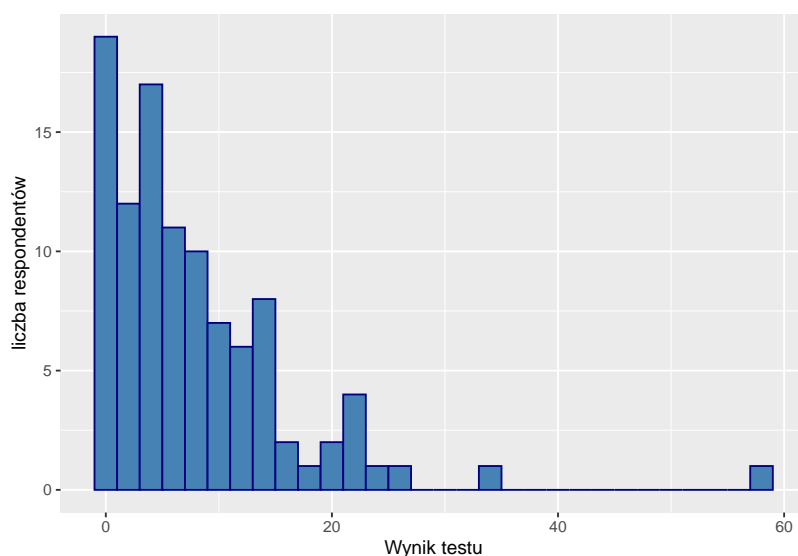
1. Depresja stanowi duży problem wśród studentów PSW.
2. Problem depresji zależy od miejsca pracy.
3. Problem depresji zależy od stażu pracy.
4. Problem depresji zależy od płci.

### 5.3.3 Zastosowane metody statystyczne

- Hipotezę 1 oceniono na podstawie odsetka respondentów wykazujących ciężką postać depresji.
- Hipotezy 2–4 zweryfikowano z wykorzystaniem tablic wielodzielnych/testu chi-kwadrat oraz porównania poziomu depresji w grupach za pomocą testów Manna-Whitneya oraz Kruskala-Wallisa.

### 5.3.4 Metryczka

W badaniu wzięło udział 103 studentów. Otrzymano 103 poprawnie wypełnionych ankiet. Średnia wartość testu Becka wyniosła 8,379 (odchylenie standardowe 8,577). Rozkład ankietowanych wg wyniku testu Becka przedstawiono na rysunku 5.5.



Rysunek 5.5: Studenci wg wyniku testu Becka

Rozkład ankietowanych ze względu na staż przedstawiono na rysunku 5.6.

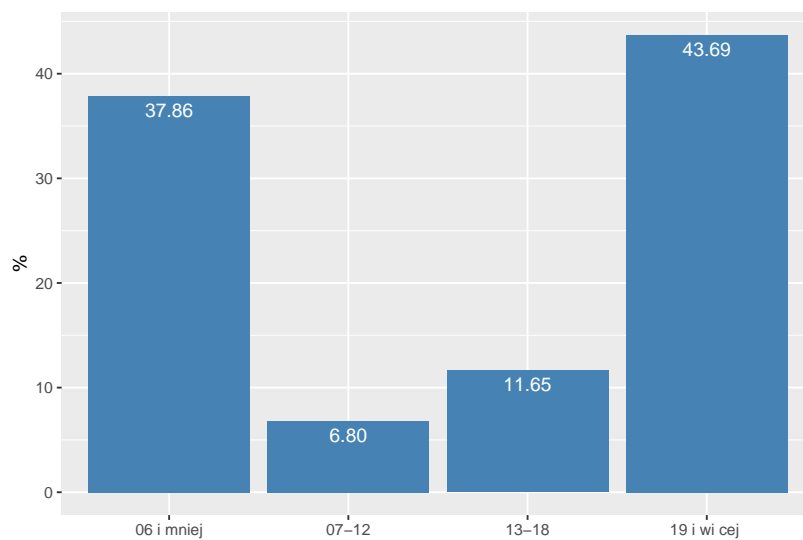
### 5.3.5 Weryfikacja hipotezy 1

Rozkład studentów wg stanu psychicznego przedstawiono na wykresie 5.7.

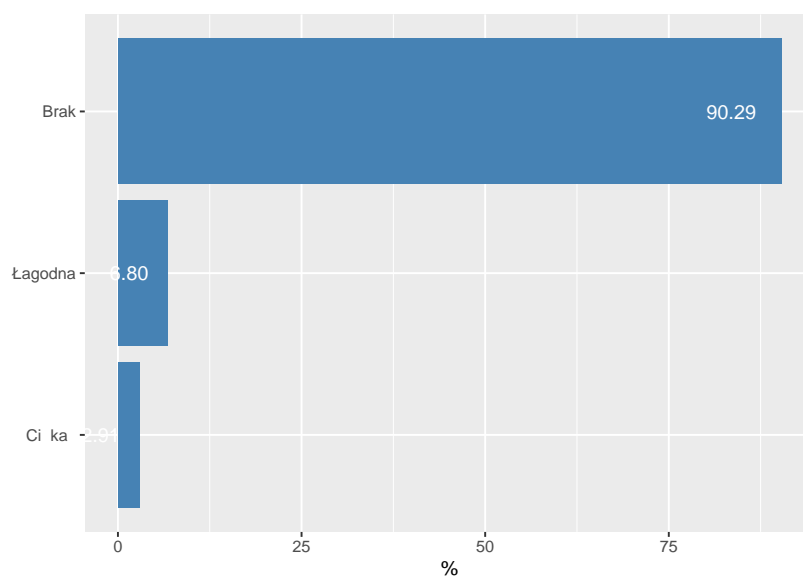
Ciężką postać depresji wykazuje zaledwie 3% studentów. Należy odrzucić hipotezę że depresja stanowi poważny problem wśród studentów RM/PO PSW.

### 5.3.6 Weryfikacja hipotez 2–4

Aby móc zastosować metody tablicy wielodzielnej i testu chi-kwadrat oryginalne wartości liczbowe depresji zamieniono na skalę porządkową: 0–19 brak/łagodna depresja (Brak); 20–25 umiarkowana (Łagodna); 26–63 ciężka (Ciężka).



Rysunek 5.6: Studenci wg stażu pracy (%)



Rysunek 5.7: Studenci wg stanu psychicznego (%)

### 5.3.7 Depresja a płeć

Tablica wielodzielna i test chi-kwadrat:

	K	M
Brak	64	29
Łagodna	5	2
Ciężka	2	1

```
##
## Pearson's Chi-squared test
##
## data:  dep.sex.f
## X-squared = 0,028134, df = 2, p-value = 0,986
```

Nie stwierdzono zależności pomiędzy depresją a płcią ( $p = 0,986$ ).

### 5.3.8 Depresja a staż

Tablica wielodzielna i test chi-kwadrat:

	06 i mniej	07-12	13-18	19 i więcej
Brak	36	6	9	42
Łagodna	2	1	2	2
Ciężka	1	0	1	1

```
##
## Pearson's Chi-squared test
##
## data:  dep.staz.f
## X-squared = 4,719, df = 6, p-value = 0,5803
```

Nie stwierdzono zależności pomiędzy depresją a stażem pracy ( $p = 0,5803$ ).

Jeżeli depresję mierzymy na (oryginalnej) skali liczbowej można zastosować test ANOVA lub Kruskala-Wallisa.

staż	średnia	mediana	n
06 i mniej	8,512821	7	39
07-12	7,857143	4	7
13-18	7,666667	3	12
19 i więcej	8,533333	6	45

Sprawdzamy czy rozkłady są normalne:

staż	S-W	p
06 i mniej	0,9008198	0,0023292
07-12	0,8565271	0,1408865
13-18	0,7596157	0,0033736
19 i więcej	0,6780397	0,0000000

Nie są. Należy zastosować test Kruskala-Wallisa:

```
## [1] 0,678923
```

Wynik jest ten sam (brak zależności). Nie ma różnicy w poziomie depresji wśród pracowników o różnym stażu pracy.

### 5.3.9 Depresja a rodzaj miejsca pracy

Tablica wielodzielna i test chi-kwadrat:

	Przychodnia	Szpital
Brak	12	81
Łagodna	0	7
Ciężka	0	3

```
##
## Pearson's Chi-squared test
##
## data:  dep.praca.f
## X-squared = 1,4605, df = 2, p-value = 0,4818
```

Nie stwierdzono zależności pomiędzy depresją a miejscem pracy ( $p = 0,4818$ ).  
Jeżeli depresję mierzymy na skali liczbowej można zastosować test Welcha lub test Manna-Whitneya.

m-pracy	średnia	mediana	n
Przychodnia	7,833333	7	12
Szpital	8,450549	6	91

Sprawdzamy czy rozkłady są normalne:

m-pracy.	S-W	p
Przychodnia	0,9256178	0,3359655
Szpital	0,7865090	0,0000000

Nie są. Należy zastosować test Manna-Whitneya:

```
## $p.value
## [1] 0,8528214
```

Wynik jest ten sam (brak zależności). Nie ma różnicy w poziomie depresji wśród pracowników szpitali i przychodni.

### 5.3.10 Wnioski

- Depresja nie jest istotnym problemem wśród studentów PSW.
- Nie ma związku pomiędzy depresją a stażem, płcią i miejscem pracy.

## 5.4 Satysfakcja, przywiązanie i zamiar odejścia

### 5.4.1 Cel

Czy satysfakcja z pracy, satysfakcja z wynagrodzenia, konflikt personalny ze współpracownikami oraz konflikt personalny z przełożonym warunkują przywiązanie do miejsca pracy oraz zamiar zmiany miejsca pracy (zamiar odejścia) w środowisku pielęgniarek/pielęgniarzy.

**Postawiono następujące hipotezy badawcze:**

1. Wysoka satysfakcja z pracy oraz satysfakcja z wynagrodzenia zmniejszają zamiar zmiany miejsca pracy.
2. Konflikt personalny zwiększa zamiar zmiany miejsca pracy.
3. Duże przywiązanie do miejsca pracy zmniejsza zamiar zmiany miejsca pracy.
4. Zamiar zmiany miejsca pracy zależy od płci i stażu pracy.
5. Praca na oddziale ratunkowym lub intensywnej terapii zwiększa zamiar zmiany miejsca pracy
6. Zmianę miejsca pracy zależy od satysfakcji z pracy, satysfakcji z wynagrodzenia, konfliktu personalnego, przywiązania do miejsca pracy, pracy na oddziale ratunkowym lub intensywnej terapii, stażu i poziom satysfakcji.

### 5.4.2 Metoda

Badanie ankietowe wśród studentów PSW przeprowadzono w 2023/2024 roku. Zamiar zmiany pracy (ZZP), przywiązanie do miejsca pracy (przywiązanie organizacyjne PO), satysfakcja z pracy (SP), satysfakcja z wynagrodzenia (SW), konflikt personalny ze współpracownikami (KPW) oraz konflikt personalny z przełożonym (KPP) były mierzone za pomocą stosownych skal pomiarowych. Ankietowani byli także pytani o płeć, staż oraz czy pracują na oddziale ratunkowym/intensywnej terapii (roddzial).

### 5.4.3 Zastosowane metody statystyczne

- Do weryfikacji hipotez 1–5 wykorzystano model regresji liniowej.
- Do weryfikacji hipotezy 6 wykorzystano model regresji logistycznej.

### 5.4.4 Wyniki

Pomniemy część opisową wyników żeby się nie powtarzać i przejdziemy od razu do weryfikacji hipotez 1–6.

Macierz korelacji dla zmiennych zzp, sp, sw, kpw, kpp oraz po:

##	zzp	sp	sw	kpw	kpp	po	staz
## zzp	1,0000	-0,5095	-0,1644	0,3245	0,3796	-0,3325	0,0666
## sp	-0,5095	1,0000	0,2815	-0,1041	-0,3446	0,3103	-0,0918
## sw	-0,1644	0,2815	1,0000	-0,4524	0,0163	0,1801	-0,1994
## kpw	0,3245	-0,1041	-0,4524	1,0000	0,4234	-0,1339	0,0030
## kpp	0,3796	-0,3446	0,0163	0,4234	1,0000	-0,3380	-0,1015
## po	-0,3325	0,3103	0,1801	-0,1339	-0,3380	1,0000	0,2808
## staz	0,0666	-0,0918	-0,1994	0,0030	-0,1015	0,2808	1,0000

Można zaobserwować wiele wysokich wartości współczynnika korelacji (zpp/sp, zpp/kpw, pp/kpp), ale też niektóre są beznadziejnie niskie (np. zpp/staż).

### Regresja liniowa

Zakładamy, że zpp jest objaśniana przez sp, sw, kpw, kpp, po, staz, plec oraz roddzial. (Odważnie dołączamy staz mimo że w świetle analizy macierzy korelacji raczej nie jest to dobry pomysł.) Oszacowanie tego modelu daje następujące wyniki:

Zmienna	B	SE	z	p	Beta	CI
(Intercept)	13,035	3,706	3,517	0,001	NA	5,57 20,50
sp	-1,099	0,331	-3,325	0,002	-0,48	-1,76 -0,43
sw	0,124	0,115	1,075	0,288	0,16	-0,11 0,36
kpw	0,778	0,376	2,068	0,044	0,31	0,02 1,54
kpp	0,032	0,235	0,137	0,891	0,02	-0,44 0,50
po	-0,040	0,026	-1,545	0,129	-0,21	-0,09 0,01
staz	0,045	0,047	0,953	0,346	0,13	-0,05 0,14
plecM	0,624	2,437	0,256	0,799	0,04	-4,28 5,53
roddzialtak	1,138	0,896	1,270	0,211	0,15	-0,67 2,94

Współczynnik zbieżności wynosi 40,85%. Tylko dwie zmienne sp oraz kpw okazały się istotne.

Stosując metodę regresji krokowej usuwamy iteracyjnie (jedną na raz) wszystkie zmienne nieistotne. Kolejno należało wyeliminować sw, kpp, po, staz, plec oraz roddzial. W modelu ostatecznym zpp jest objaśniana przez sp oraz kpw:

Zmienna	B	SE	z	p	Beta	CI
(Intercept)	13,987	2,969	4,710	0,000	NA	8,03 19,95
sp	-1,112	0,263	-4,226	0,000	-0,48	-1,64 -0,58
kpw	0,682	0,283	2,412	0,019	0,27	0,11 1,25

Współczynnik zbieżności wynosi 33,41%. Satysfakcja jest większym predyktorem zamiaru zmiany pracy (kolumna Beta).

### Regresja logistyczna

Żeby zademonstrować przykład wykorzystania regresji logistycznej przyjmijmy, że jak ktoś bardzo chce zmienić pracę to ją zmieni. Niech to **bardzo chce** będzie wtedy kiedy wartość zpp wynosi co najmniej 12. Zmienną zpp należy w tym celu przekodować na zmienną dwuwartościową (nazwijmy ją zp), która przyjmuje wartość 0 jeżeli zpp jest większe równe od 12 lub wartość 1 jeżeli zpp jest mniejsze od 12. Model nie objaśnia teraz zależności pomiędzy zamiarem, ale prognozuje zmianę miejsca pracy (zmieni=0, nie zmieni=1).

Zakładamy, że zp, podobnie jak w przypadku zamiaru, jest objaśniana przez sp, sw, kpw, kpp, po, staz, plec oraz roddzial.



Parametr	Ocena	SE	z	p	OR	CI
(Intercept)	-4,702	4,599	-1,022	0,307	0,01	0,00 82,27
sp	1,082	0,480	2,254	0,024	2,95	1,30 9,15
sw	-0,082	0,143	-0,577	0,564	0,92	0,68 1,22
kpw	-0,415	0,410	-1,013	0,311	0,66	0,28 1,51
kpp	0,154	0,386	0,399	0,690	1,17	0,62 2,53
po	0,037	0,036	1,028	0,304	1,04	0,97 1,13
staz	-0,107	0,071	-1,512	0,130	0,90	0,76 1,02
plecM	13,102	2596,962	0,005	0,996	489784,87	0,00 NA
roddzialtak	-1,817	1,074	-1,692	0,091	0,16	0,02 1,31

Wartości większości współczynników okazały się nieistotne. Mówiąc konkretniej na poziomie istotności 0,05 tylko sp jest istotna (roddzialtak jest także istotny ale na poziomie 0,1).

Postępując podobnie jak w przypadku „normalnej” regresji eliminujemy iteracyjnie wszystkie zmienne nieistotne ostatecznie dochodząc do modelu w którym odejście z pracy objaśnia tylko satysfakcja:

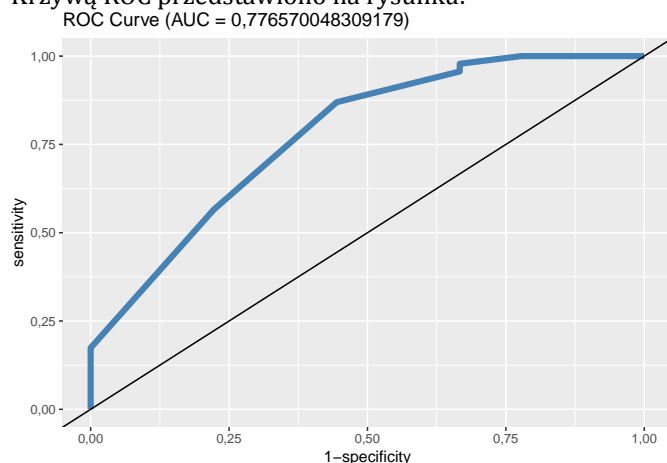
Parametr	Ocena	SE	z	p	OR	CI
(Intercept)	-5,448	2,460	-2,215	0,027	0,00	0,00 0,36
sp	0,803	0,283	2,838	0,005	2,23	1,35 4,26

Macierz pomyłek:

```
##          zp
## Prognoza  0  1
##           0  3  1
##           1  6 45
```

Stąd: czułość 0,98; swoistość 0,33.

Krzywą ROC przedstawiono na rysunku:



### 5.4.5 Wnioski

- Satysfakcja jest większym predyktorem zamiaru zmiany pracy niż konflikt personalny ze współpracownikami. Pozostałe czynniki okazały się nieistotne (hipotezy 1–5).
- Satysfakcja jest jedynym czynnikiem istotnie wpływającym na zmianę miejsca pracy. Pozostałe czynniki okazały się nieistotne (hipoteza 6).

## 5.5 Formularze ankiet

### 5.5.1 Skala Depresji Becka

#### Pytania

##### Pytanie 1

0. Nie jestem smutny ani przygnębiony.
1. Odczuwam często smutek, przygnębienie
2. Przeżywam stale smutek, przygnębienie i nie mogę uwolnić się od tych przeżyć.
3. Jestem stale tak smutny i nieszczęśliwy, że jest to nie do wytrzymania.

##### Pytanie 2

0. Nie przejmuję się zbyt przyszlnością.
1. Często martwię się o przyszlność.
2. Obawiam się, że w przyszlności nic dobrego mnie nie czeka.
3. Czuję, że przyszlność jest beznadziejna i nic tego nie zmieni.

##### Pytanie 3

0. Sądę, że nie popełniam większych zaniedbań.
1. Sądę, że czynię więcej zaniedbań niż inni.
2. Kiedy spoglądam na to, co robiłem, widzę mnóstwo błędów i zaniedbań.
3. Jestem zupełnie niewydolny i wszystko robię źle.

##### Pytanie 4

0. To, co robię, sprawia mi przyjemność.
1. Nie cieszy mnie to, co robię.
2. Nic mi teraz nie daje prawdziwego zadowolenia.
3. Nie potrafię przeżywać zadowolenia i przyjemności; wszystko mnie nuży.

##### Pytanie 5

0. Nie czuję się winnym ani wobec siebie, ani wobec innych.
1. Dość często miewam wyrzuty sumienia.
2. Często czuję, że zawiniłem.
3. Stale czuję się winny.

##### Pytanie 6

0. Sądę, że nie zasługuję na karę
1. Sądę, że zasługuję na karę
2. Spodziewam się ukarania
3. Wiem, że jestem karany (lub ukarany)

##### Pytanie 7

0. Jestem z siebie zadowolony
1. Nie jestem z siebie zadowolony
2. Czuję do siebie niechęć

3. Nienawidzę siebie

Pytanie 8

0. Nie czuję się gorszy od innych ludzi

1. Zarzucam sobie, że jestem nieudolny i popełniam błędy

2. Stale potępiam siebie za popełnione błędy

3. Winię siebie za wszelkie zło, które istnieje

Pytanie 9

0. Nie myślę o odebraniu sobie życia

1. Myślę o samobójstwie – ale nie mógłbym tego dokonać

2. Pragnę odebrać sobie życie

3. Popełnię samobójstwo, jak będzie odpowiednia sposobność

Pytanie 10

0. Nie płaczę częściej niż zwykle

1. Płaczę częściej niż dawniej

2. Ciągłe chce mi się płakać

3. Chciałbym płakać, lecz nie jestem w stanie

Pytanie 11

0. Nie jestem bardziej podenerwowany niż dawniej

1. Jestem bardziej nerwowy i przykry niż dawniej

2. Jestem stale zdenerwowany lub rozdrażniony

3. Wszystko, co dawniej mnie drażniło, stało się obojętne

Pytanie 12

0. Ludzie interesują mnie jak dawniej

1. Interesuję się ludźmi mniej niż dawniej

2. Utraciłem większość zainteresowań innymi ludźmi

3. Utraciłem wszelkie zainteresowanie innymi ludźmi

Pytanie 13

0. Decyzje podejmuję łatwo, tak jak dawniej

1. Częściej niż kiedyś odwlekam podjęcie decyzji

2. Mam dużo trudności z podjęciem decyzji

3. Nie jestem w stanie podjąć żadnej decyzji

Pytanie 14

0. Sądzę, że wyglądam nie gorzej niż dawniej

1. Martwię się tym, że wyglądam staro i nieatrakcyjnie

2. Czuję, że wyglądam coraz gorzej

3. Jestem przekonany, że wyglądam okropnie i odpychająco

Pytanie 15

0. Mogę pracować jak dawniej

1. Z trudem rozpoczynam każdą czynność

2. Z wielkim wysiłkiem zmuszam się do zrobienia czegokolwiek

3. Nie jestem w stanie nic zrobić

Pytanie 16

0. Sypiam dobrze, jak zwykle

1. Sypiam gorzej niż dawniej

2. Rano budzę się 1–2 godziny za wcześnie i trudno jest mi ponownie usnąć

3. Budzę się kilka godzin za wcześnie i nie mogę usnąć

Pytanie 17

0. Nie męczę się bardziej niż dawniej
1. Męczę się znacznie łatwiej niż poprzednio.
2. Męczę się wszystkim, co robię.
3. Jestem zbyt zmęczony, aby cokolwiek robić.

## Pytanie 18

0. Mam apetyt nie gorszy niż dawniej
1. Mam trochę gorszy apetyt
2. Apetyt mam wyraźnie gorszy
3. Nie mam w ogóle apetytu

## Pytanie 19

0. Nie tracę na wadze (w okresie ostatniego miesiąca)
1. Straciłem na wadze więcej niż 2 kg
2. Straciłem na wadze więcej niż 4 kg
3. Straciłem na wadze więcej niż 6 kg

## Pytanie 20

0. Nie martwię się o swoje zdrowie bardziej niż zawsze
1. Martwię się swoimi dolegliwościami, mam rozstrój żołądka, zaparcie, bóle
2. Stan mojego zdrowia bardzo mnie martwi, często o tym myślę
3. Tak bardzo martwię się o swoje zdrowie, że nie mogę o niczym innym myśleć

## Pytanie 21

0. Moje zainteresowania seksualne nie uległy zmianom
1. Jestem mniej zainteresowany sprawami płci (seksu)
2. Problemy płciowe wyraźnie mnie mniej interesują
3. Utraciłem wszelkie zainteresowanie sprawami seksu

**Treść pytań** (nie prezentowana ankietowanym): Odczuwanie smutku i przygnębienia (1), Martwienie się o przyszłość (2), Uważasz, że zaniedbujesz swoje obowiązki? (3), Jesteś zadowolony z siebie? (4), Czy często masz poczucie winy? (5), Czy zasługujesz na karę? (6), Zadowolenie z siebie (7), Czy czujesz się gorszy od innych? (8), Czy masz myśli samobójcze? (9), Często chce Ci się płakać? (10), Jesteś ostatnio bardziej nerwowy i rozdrażniony? (11), Czy zmieniło się coś w Twoim zainteresowaniu innymi ludźmi? (12), Czy ostatnio miewasz większe problemy z podejmowaniem różnych decyzji? (13), Czy uważasz, że wyglądasz gorzej i mniej atrakcyjnie niż kiedyś? (14), Czy masz większe trudności z wykonywaniem różnych prac i zadań? (15), Masz kłopoty ze snem? (16), Czy męczysz się bardziej niż zwykle? (17), Czy masz kłopoty z apetytem? (18), W ciągu ostatniego miesiąca nie stosowałem diety, aby schudnąć, lecz straciłem na wadze (19), Czy ostatnio bardziej martwisz się swoim stanem zdrowia? (20), Czy masz kłopoty z potencją? (21).

**Interpretacja wyników**

Punkty	Depresja
0–11	Brak
12–19	Łagodna
20–25	Umiarkowana
26–63	Ciężka

**Źródło:** <https://psychiatra.bydgoszcz.eu/publikacje-dla-pacjenta/depresja/skala->

depresji-becka/ oraz [http://centrum-psychologiczne.com/files/files/Skala\\_Depresji\\_Beck\\_i\\_Beck\\_a\\_word.pdf](http://centrum-psychologiczne.com/files/files/Skala_Depresji_Beck_i_Beck_a_word.pdf)

### 5.5.2 Ankieta na temat szkodliwości palenia

Poziom wiedzy personelu pielęgniarskiego na temat szkodliwości palenia tytoniu

#### Pytania

1. Płeć ☐ Kobieta ☐ Mężczyzna
2. Wiek
3. Pochodzenie ☐ wieś ☐ Miasto do 20 tys. mieszkańców ☐ Miasto powyżej 20 tys. mieszkańców
4. Wykształcenie ☐ Średnie medyczne ☐ Licencjat pielęgniarstwa ☐ Magister pielęgniarstwa ☐ Inne wyższe
5. Staż pracy ☐ Mniej niż rok ☐ 1-10 lat ☐ 10-15 lat ☐ Więcej niż 15 lat
6. Miejsce pracy ☐ Oddział zabiegowy ☐ Oddział zachowawczy ☐ Przychodnia/poradnia
7. Czy kiedykolwiek paliłeś/aś papierosy? ☐ Paliłem/łam ☐ W dalszym ciągu palę ☐ Nigdy nie paliłem/am
8. Od ilu lat palisz? ☐ Nie palę/Nie dotyczy ☐ Mniej niż rok ☐ 1-10 lat ☐ 11-15 lat ☐ Więcej niż 15 l
9. Czy zdarza Ci się palić w miejscu pracy? ☐ Tak ☐ Nie ☐ Nie dotyczy
10. Czy próbowałeś kiedykolwiek rzucić palenie? ☐ Tak, udało się ☐ Tak, ale wróciłem/am do nałogu ☐ Nie ☐ Nie dotyczy
11. Czy paląc przyznajesz się do uzależnienia ☐ Tak ☐ Nie ☐ Nie dotyczy
12. Uważasz, że bardziej szkodliwe dla zdrowia jest: ☐ Palenie czynne - dym tytoniowy wdychany bezpośrednio przez palacza ☐ Palenie bierne-boczny strumień dymu ☒ Każda forma kontaktu z dymem jest równie szkodliwa ☐ Nie wiem/Nie mam zdania
13. Czy uważasz, że przepisy prawa powinny zabraniać palenia w obecności dzieci poniżej 15 roku życia? ☐ Tak ☐ Nie ☐ Nie wiem/nie mam zdania
14. Czy przerwa na papierosa pomaga Ci w sytuacjach stresowych? ☐ Tak ☐ Nie ☐ Nie dotyczy
15. Czy palenie nikotyny sprawia Ci przyjemność? ☐ Tak ☐ Nie ☐ Nie dotyczy
16. Jakże według Ciebie choroby układu oddechowego mogą być spowodowane bezpośrednio przez palenie papierosów?(wielokrotnego) ☒ Przewlekła obturacyjna choroba płuc ☒ Astma oskrzelowa ☒ Alergie wziewne ☐ Gruźlica ☐ Zapalenie płuc ☒ Przewlekłe zapalenie oskrzeli ☒ Infekcje dróg oddechowych ☐ Palenie nie powoduje chorób układu oddechowego ☐ Nie wiem
17. Jakże według Ciebie choroby kardiologiczne mogą być spowodowane bezpośrednio przez palenie papierosów? (wielokrotnego) ☒ Nadciśnienie tętnicze krwi ☒ Zawał mięśnia sercowego ☒ Udar mózgu ☒ Choroba niedokrwienna serca ☒ Miażdżyca tętnic obwodowych ☒ Zaburzenie rytmu serca ☒ Choroba Buergera ☐ Hipercholesterolemia ☐ Tętniak aorty ☐ Palenie nie powoduje chorób kardiologicznych
18. Czy palenie papierosów powoduje choroby układu pokarmowego? ☒ Tak ☐ Nie ☐ Nie wiem
19. Jaki według Ciebie ma wpływ palenie papierosów na narządy zmysłów? (wie-

lokrotnego) ☒ Upośledza węch i smak ☒ Powoduje podrażnienie spojówek  
☒ Obniża apetyt ☒ Niszczy struny głosowe ☒ Zmniejsza ostrość wzroku ☐  
 Palenie nie ma negatywnego wpływu na narządy zmysłu.

20. Jak oceniasz swoją wiedzę na temat palenia papierosów i jego wpływu na zdrowie człowieka? ☐ Bardzo dobrze ☐ Dobrze ☐ Przeciętnie ☐ Źle

21. Czy wzorce sięgania po tytoń wyniosłaś/wyniosłeś z domu rodzinnego ☐ Tak  
☐ Nie ☐ Nie dotyczy

(Prawidłowe odpowiedzi zaznaczono ☒)

**Źródło:** Na podstawie ankiety z pracy Izabeli Marendowskiej *Ocena poziomu wiedzy pielęgniarek na temat szkodliwości palenia papierosów* (promotor: dr Marzena Barton); Kwidzyn 2022.

### 5.5.3 Ankieta satysfakcja, przywiązanie i zamiar odejścia

#### Zamiar zmiany pracy/ZZP

(Zdecydowanie się nie zgadzam/Nie zgadzam się/Trudno powiedzieć/Zgadzam się/Zdecydowanie się zgadzam.)

1. Często poważnie rozważam odejście z obecnej pracy
2. Zamierzam rzucić obecną pracę
3. Zacząłem szukać innej pracy

#### Satysfakcja z pracy/SP

(Zdecydowanie się nie zgadzam/Nie zgadzam się/Trudno powiedzieć/Zgadzam się/Zdecydowanie się zgadzam.)

4. Ogólnie rzecz biorąc nie lubię swojej pracy
5. Ogólnie rzecz biorąc jestem zadowolony ze swojej pracy
6. Ogólnie rzecz biorąc, lubię tu pracować

#### Satysfakcja z pracy/SSP/alternatywna skala

(Zdecydowanie się nie zgadzam/Nie zgadzam się/Trudno powiedzieć/Zgadzam się/Zdecydowanie się zgadzam.)

7. Pod bardzo wieloma względami moja praca bliska jest ideału
8. Mam świetne warunki pracy
9. Jestem zadowolony z pracy
10. Jak dotąd w pracy udawało mi się osiągać to, czego chciałem
11. Gdybym miał decydować raz jeszcze, wybrałbym tę samą pracę

#### Satysfakcja z wynagrodzenia/SW

(Zdecydowanie niesatysfakcjonująca-Niesatysfakcjonująca-TP-Satysfakcjonująca-Zdecydowanie niesatysfakcjonująca.)

12. Moja pensja na rękę jest
13. Wielkość mojej obecnej pensji jest
14. Moje obecne wynagrodzenie jest
15. Poziom mojego łącznego wynagrodzenia jest

#### Konflikt ze współpracownikami/KPW

(Raz na miesiąc lub mniej/1–2 razy w miesiącu/1–2 razy w tygodniu/1–2 razy dziennie/wiele razy dziennie.)

16. Jak często w pracy kłócisz się z innymi osobami?
17. Jak często inni ludzie krzyczą na Ciebie w pracy?
18. Jak często ludzie są wobec Ciebie niemili w pracy?

19. Jak często inni ludzie robią Ci nieprzyjemne rzeczy w pracy?

**Konflikt z przełożonym/KPP**

(Raz na miesiąc lub mniej/1–2 razy w miesiącu/1–2 razy w tygodniu/1–2 razy dziennie/wiele razy dziennie.)

20. Jak często w pracy kłócisz się z przełożonym?

21. Jak często twój przełożony krzyczy na ciebie?

22. Jak często twój przełożony jest niemiły w pracy?

23. Jak często Twój przełożony robi ci nieprzyjemne rzeczy w pracy?

**Przywiązanie do organizacji/PO** (PO-1–PO-18), składające się z następujących aspektów: przywiązanie afektywne (PO-1–PO-6), przywiązanie normatywne (PO-7–PO-12), przywiązanie trwania (PO-13–PO-18).

(Zdecydowanie się nie zgadzam/nie zgadzam się/raczej się nie zgadzam/ trudno powiedzieć/raczej się zgadzam/zgadzam się/zdecydowanie się zgadzam.)

24. Zakład pracy, w której pracuję ma dla mnie duże znaczenie osobiste. (PO-2)

25. Czułbym się winny, gdybym teraz odszedł z mojego zakładu pracy (PO-18) 26  
.Zbyt wiele straciłbym w moim życiu, decydując się teraz na odejście z mojego zakładu pracy (PO-9)

26. Mogę powiedzieć, że czuję się w moim zakładzie pracy jak w rodzinie (PO-5)

27. Gdybym dostał ofertę lepszej pracy, czułbym się nie w porządku odchodząc z mojego zakładu pracy (PO-14)

28. Byłbym bardzo zadowolony gdybym do emerytury mógł pracować w moim zakładzie pracy (PO-11)

29. Wiele zawdzięczam mojemu zakładowi pracy (PO-4)

30. Jednym z głównych powodów, dla których wciąż pracuję w tym zakładzie pracy jest wiara w znaczenie lojalności

dająca mi poczucie moralnego obowiązku pozostania (PO-16)

31. Czuje, że problemy mojego zakładu pracy są rzeczywiście moimi własnymi problemami (PO-6)

32. Nie odszedłbym teraz z mojego zakładu pracy, ponieważ mam zobowiązania w stosunku do ludzi, którzy w niej pracują (PO-17)

33. Jedną z kilku negatywnych konsekwencji odejścia z mojego zakładu pracy mógłby być brak dostępnych możliwości zatrudnienia (PO-10)

34. Byłoby mi bardzo ciężko odejść teraz z mojego zakładu pracy, nawet gdybym chciał (PO-3)

35. Ten zakład pracy zasługuje na to, żebym był wobec niego w porządku (PO-15)

36. Mam poczucie, że pozostanie w zakładzie pracy jest dla mnie koniecznością (PO-12)

37. Sprawia mi przyjemność, kiedy mogę porozmawiać o moim zakładzie pracy z ludźmi z zewnątrz (PO-1)

38. Sądzę, że odchodząc z mojego zakładu pracy, mam zbyt mało innych możliwości do wyboru (PO-8)

39. Nawet, gdyby to było dla mnie korzystne, nie czułbym się w porządku odchodząc teraz z mojego zakładu pracy (PO-13)

40. Lepiej było, kiedy ludzie większość swojego życia zawodowego wiązali z jedną firmą (PO-7)

**Satysfakcja z pracy/SP0/skala jednopozycyjna**

42. W jakim stopniu jesteś ogólnie zadowolony ze swojej pracy?

(Bardzo niezadowolony/Niezadowolony/Trudno powiedzieć/Zadowolony/Bardzo zadowolony)

**Zamiar odejścia z pracy/ZZP0/skala jednopozycyjna**

(Nigdy/Rzadko/Czasami/Często/Bardzo często)

43. Jak często poważnie rozważałeś odejście z obecnej pracy?

**„Metryczka”**

44. Płeć (K/M)

45. Rodzaj studiów

46. Staż łączny (lata pracy)

47. Ile razy zmieniałeś miejsce pracy?

48. Praca na oddziale ratunkowym lub intensywnej terapii (T/N)

**Źródło:** Liu Y i inni, *Turnover intention and its associated factors among nurses: a multi-center cross-sectional study* (Zamiar zmiany pracy oraz czynniki z nim związane wśród pielęgniarek: badanie przekrojowe.) *Front. Public Health* 11:1141441 (2023) <https://www.frontiersin.org/articles/10.3389/fpubh.2023.1141441/full> oraz Bańka A., Bazińska R., Wołowska A., *Polska wersja Meyera i Allen Skali Przywiązania do Organizacji*, *Czasopismo Psychologiczne* 2002, nr 8 (1), s. 65–74.



## Rozdział 6

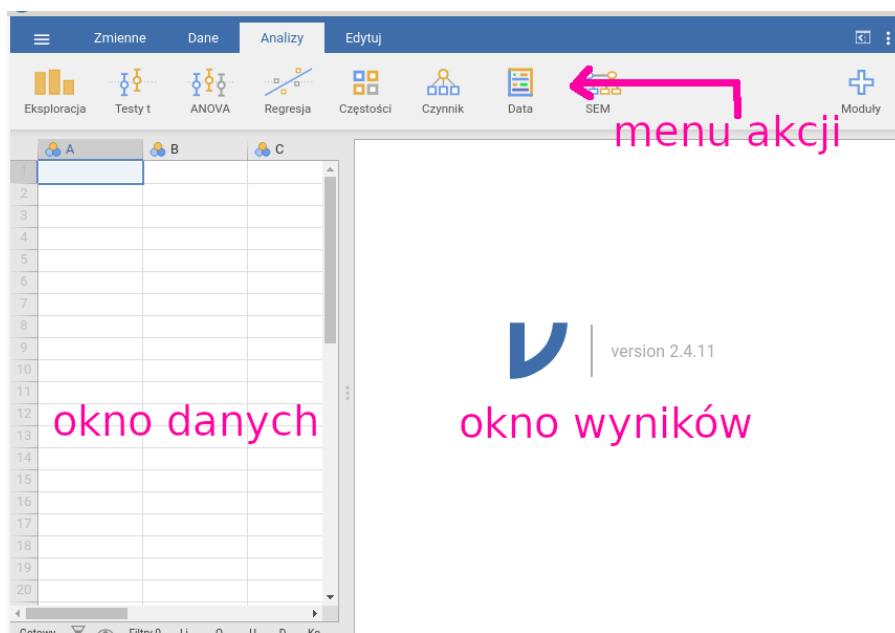
# Praca z programem Jamovi

Jak wspomniano w rozdziale 1, statystykę można uprawiać (tj. liczyć statystyki w drugim tego słowa znaczeniu :-)) wykorzystując różne programy. My zdecydowaliśmy się promować **Jamovi**, program który naszym zdaniem jest najlepszym – z punktu widzenia większości studentów Nauk o Zdrowiu – połączeniem ceny, możliwości, prostoty i łatwości nauki.

## 6.1 Podstawy pracy z Jamovi

Jamovi jest oprogramowaniem rozpowszechnianym na licencji typu *Open Source*, a więc można go używać za darmo. Program jest dostępny ze strony <https://www.jamovi.org/download.html>. Klikamy, ściągamy, uruchamiamy instalator. Program jest dość duży, ale to nie jest aż tak wielki problem w czasach kiedy pojemności dysków w domowym komputerze zaczynają się od 250 gigabajtów. Jest też wersja Jamovi365, którą można się posługiwać w przeglądarce i wtedy nic nie trzeba instalować na swoim komputerze.

W naszym podręczniku pokazujemy jak uprawiać statystykę z tradycyjną wersją Jamovi. Po zainstalowaniu uruchamiamy program, którego ekran startowy wygląda jak na rysunku 6.1.



Rysunek 6.1: Ekran startowy Jamovi

Menu akcji umożliwia wykonanie podstawowych akcji:

- wczytanie danych i zapisanie danych (pierwsza pozycja menu oznaczona jako trzy poziome kreski);
- podgląd (w sensie skontrolowania wartości zmiennych) i modyfikację danych (pozycje **Zmienne** oraz **Dane**);
- wykonanie obliczeń (pozycja **Analizy**);
- modyfikowanie raportu (pozycja **Edit**).

Typowa sesja w **Jamovi**:

1. Wczytanie danych z pliku o praktycznie dowolnym formacie. Jeżeli przykładowo dane są wynikiem wykonania badania ankietowego z wykorzystaniem Formularzy Google to zalecamy posługiwanie się formatem CSV.
2. Transformacja danych. Przekodowanie wartości nominalnych na rangi. Przekodowanie wartości liczbowych na nominalne. Odwrócenie pytań odwróconych. Obliczenie sum/średnich rang dla wielu zmiennych.
3. Wykonanie obliczeń:
  1. Analiza struktury (**Eksploracja**).
  2. Analiza zależności między zmienną liczbową a nominalną (*testy t/ANOVA*).
  3. Analiza zależności między zmiennymi liczbowymi: współczynnik korelacji liniowej/macierz korelacji (*Regresja*).
  4. Analiza zależności między zmienną liczbową a zmiennymi liczbowymi/nominalnymi: regresja liniowa i logistyczna (**Regresja**).
  5. Analiza zależności między zmiennymi nominalnymi: tablica wielodzielna, test chi-kwadrat zgodności (*Częstości*).

Wykonanie obliczeń jest banalnie proste i sprowadza się do wybrania myszką odpowiednich zmiennych oraz procedury, która ma być wykonana. Wynik obliczeń pojawia się natychmiast w **oknie wyników**. Jeżeli coś nam nie wyszło można procedurę poprawić a poprzedni wynik usunąć z okna wynikowego.
4. Zapisania danych (pozycja trzy poziome kreski). Po skończeniu pracy wynik można zapisać, żeby np. wysłać wykładowcy lub nie zaczynać od zera jeżeli będziemy musieli pracę kontynuować, bo wykładowca chciał żebyśmy coś poprawili.

## 6.2 Analiza ankiety: satysfakcja – wiedza o paleniu – zamiar odejścia

Przykład nieco absurdalny, ale za to w zwartej postaci ilustrujący praktyczne sposoby transformacji danych oraz wykorzystania wszystkich procedur omawianych w podręczniku.

### 6.2.1 Wczytanie danych

W wyniku przeprowadzenia badania ankietowego zebrano za pomocą Formularza Google dane dotyczące satysfakcji/zamiaru odejścia oraz wiedzy nt. szkodliwości palenia tytoniu. Wyniki wyeksportowano do arkusza kalkulacyjnego, którego początek wygląda jak na rysunku 6.2:

Ankieta składa się z 10 następujących pytań:

Ogólnie rzecz biorąc nie lubię swojej pracy (kolumna B), Ogólnie rzecz biorąc jestem

	A	B	C	D	E	
1	Sygnatura c	Ogólnie rzecz biorąc	Ogólnie rzecz biorąc j	Ogólnie rzecz biorąc, lu	Jakie według Ciebie choroby u	bepośrednio przez palenie pę Cze
2	2023-12-01	Zdecydowanie się nie	Nie zgadzam się	Zgadzam się	Alergie wziewne, Zapalenie płuc, In	
3	2023-12-02	Zdecydowanie się nie	Zdecydowanie się zg	Zdecydowanie się zgac	Przewlekła obturacyjna choroba	Zdei
4	2023-12-02	Nie zgadzam się	Zgadzam się	Zgadzam się	Przewlekła obturacyjna choroba	Zdei
5	2023-12-02	Zdecydowanie się nie	Nie zgadzam się	Nie zgadzam się	Astma oskrzelowa	Nie :
6	2023-12-02	Zdecydowanie się nie	Zdecydowanie się nie	Zdecydowanie się nie z	Przewlekła obturacyjna choroba	Zdei
7	2023-12-02	Nie zgadzam się	Zgadzam się	Trudno powiedzieć	Astma oskrzelowa	Trud
8	2023-12-02	Zdecydowanie się nie	Zdecydowanie się zg	Zdecydowanie się zgac	Przewlekła obturacyjna choroba	Zdei
9	2023-12-02	Nie zgadzam się	Zgadzam się	Zgadzam się	Zapalenie płuc	Nie :
10	2023-12-02	Trudno powiedzieć	Trudno powiedzieć	Trudno powiedzieć	Przewlekła obturacyjna choroba	Trud
11	2023-12-02	Nie zgadzam się	Zgadzam się	Zgadzam się	Przewlekła obturacyjna choroba	Nie :

Rysunek 6.2: Fragment przykładowej ankiety

zadowolony ze swojej pracy (C), Ogólnie rzecz biorąc, lubię tu pracować (D), Jakie według Ciebie choroby układu oddechowego mogą być spowodowane bezpośrednio przez palenie papierosów? (E), Często poważnie rozważam odejście z obecnej pracy (F), Zamierzam rzucić obecną pracę (G), Zacząłem szukać innej pracy (H), Płeć (I), Wiek (w latach) (J), oraz Staż pracy (K).

Ponadto Formularz Google dodał automatycznie sygnaturę czasową jako zawartość pierwszej kolumny (A).

Zmieniamy wartości w pierwszym wierszu, który powinien zawierać nazwy zmiennych. Nazwy zmiennych powinny być jednowyrazowe i w miarę krótkie żeby się później można nimi wygodnie posługiwać. Jednocześnie nie powinny być za krótkie żeby od razu było widać jakie dane zawiera zmienna.

Jak widać pytania z kolumn B–D mierzą to samo (satysfakcję) więc zmieniamy im nazwę na bardziej zwartą s1, s2 oraz s3 (s od satysfakcja). Podobnie ponieważ pytania z kolumn F–H też mierzą to samo (zamiar odejścia), to też zmieniamy nazwy na coś krótszego: zo1, zo2, zo3. Kolumnę E nazywamy wiedza\_nt\_palenia a kolumny I, J oraz K odpowiednio: plec, wiek oraz staz.

Teraz arkusz wygląda jak na rysunku 6.3.

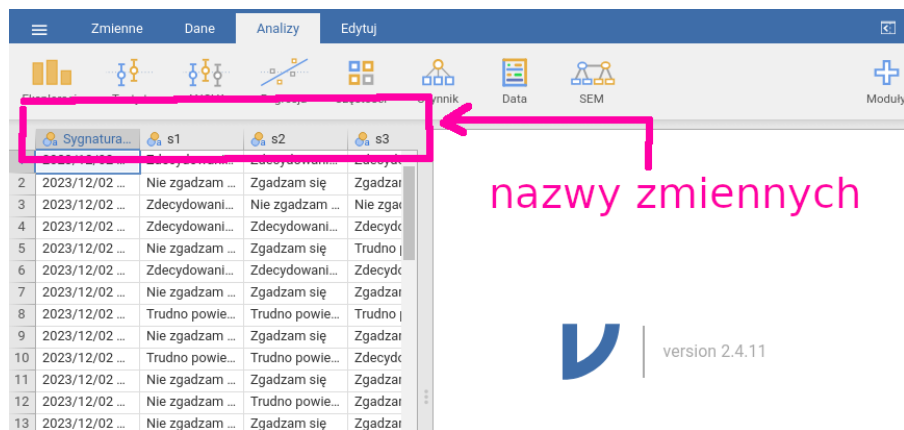
Arkusz eksportujemy wybierając format CSV. Bez problemu powinniśmy go wczytać do Jamovi (trzy poziome kreski →Otwórz). Jeżeli import się powiódł, to powinniśmy zobaczyć coś podobnego do tego co widać na rysunku 6.4.

Reasumując:

- Pytania oznaczone jako s1/s2/s3 mierzą **satysfakcję z pracy**; pytania zo1/zo2/zo3 mierzą **zamiar odejścia z pracy**. Pytania s1–s3 oraz zo1–zo3 są pytaniami jednokrotnego wyboru.
- Pytanie oznaczone jako wiedza\_nt\_palenia mierzy wiedzę na temat palenia tytoniu. Jest to przykład wykorzystania pytania z wielokrotnym wyborem.
- Pytania plec, wiek, staz mierzą płeć (kobieta/mężczyzna), wiek (lata ukończone) oraz staż pracy (lata przepracowane).
- Pierwsza kolumna nie jest potrzebna, ale jest dodawana przez aplikację Formularze Google.

	A	B	C	D	E	
1	Sygnatura czasowa	s1	s2	s3	wiedza_nt_palenia	zo1
2	2023-12-01 15:45:5	Zdecydowanie się n	Nie zgadzam się	Zgadzam się	Alergie wziewne, Zapalenie	
3	2023-12-02 11:55:1	Zdecydowanie się n	Zdecydowanie się z	Zdecydowanie się z	Przewlekła obturacy	Zdecyc
4	2023-12-02 11:55:2	Nie zgadzam się	Zgadzam się	Zgadzam się	Przewlekła obturacy	Zdecyc
5	2023-12-02 11:56:0	Zdecydowanie się n	Nie zgadzam się	Nie zgadzam się	Astma oskrzelowa	Nie zg
6	2023-12-02 11:56:1	Zdecydowanie się n	Zdecydowanie się n	Zdecydowanie się n	Przewlekła obturacy	Zdecyc
7	2023-12-02 11:56:3	Nie zgadzam się	Zgadzam się	Trudno powiedzieć	Astma oskrzelowa	Trudnc
8	2023-12-02 11:56:3	Zdecydowanie się n	Zdecydowanie się z	Zdecydowanie się z	Przewlekła obturacy	Zdecyc
9	2023-12-02 11:56:3	Nie zgadzam się	Zgadzam się	Zgadzam się	Zapalenie płuc	Nie zg
10	2023-12-02 11:56:4	Trudno powiedzieć	Trudno powiedzieć	Trudno powiedzieć	Przewlekła obturacy	Trudnc
11	2023-12-02 11:56:4	Nie zgadzam się	Zgadzam się	Zgadzam się	Przewlekła obturacy	Nie zg

Rysunek 6.3: Fragment przykładowej ankiety



Rysunek 6.4: Import danych

### 6.2.2 Przekodowanie danych

Zwykle zawartość arkusza zawierającego wyniki ankiety wymaga przekodowania. W naszym przykładzie należy wykonać:

- Zmienne s1–s3 oraz zo1–zo3 są mierzone w skali porządkowej. Wartości tych zmiennych chcemy zmienić (przekodować) na rangi wg schematu: Zdecydowanie się nie zgadzam = 1; Nie zgadzam się = 2; Trudno powiedzieć = 3 itd. Dodatkowo zauważmy że s1 jest pytaniem odwróconym. W takich pytaniach należy przeliczyć rangi wg prostej formuły  $s1r = 6 - s1$ .
  - Miarą satysfakcji będzie suma rang  $s1r + s2 + s3$ .
  - Miarą zamiaru odejścia będzie suma rang  $zo1 + zo2 + zo3$ .
- Zmienna plec jest mierzona w skali nominalnej. Nie musimy jej przekodowywać
- Wartość zmiennej wiedza\_nt\_palenia należy przekodować na liczbę wg schematu: za wybranie poprawnej odpowiedzi plus jeden punkt; za wybranie błędnej odpowiedzi minus jeden punkt.

- Miarą wiedzy nt. palenia będzie suma punktów uzyskanych za odpowiedzi prawidłowe minus suma punktów uzyskanych za odpowiedzi nieprawidłowe.

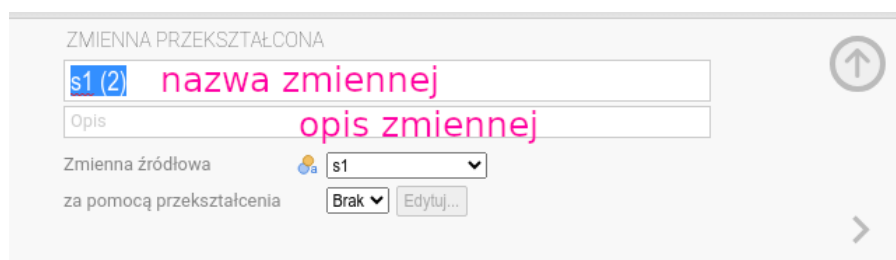
Uwaga: Sposób mierzenia wiedzy nt. palenia jest niepotrzebnie pokręcony; zamiast pytania z wielokrotnym wyborem spośród 8 możliwości/wariantów prościej jest zastosować 8 pytań Tak/Nie po czym pytania poprawne zsumować a pytania niepoprawne też dodać a wartość odjąć od sumy uzyskanej dla pytań poprawnych. My o tym wiemy, że tak jest bez sensu ale pokazujemy jako przykład przekodowania pytania z wielokrotnym wyborem.

- Wartości zmiennych wiek oraz staz są liczbami. Mogą być analizowane tak-jak-są (regresja/korelacja), ale można też je przekodować na wartości nominalne (mały-średni-duży staż) i zastosować metody z grupy zmienna-liczbowa/zmienna nominalna (takie jak test ANOVA czy Kruskala-Wallisa).

Przekodowanie wykonujemy wybierając **Dane** w menu głównym.

1. Klikamy w nazwę zmienną, którą zamierzamy przekodować. Niech to będzie s1. Kolumna po kliknięciu zmieni kolor.
2. Wybieramy ikonę Przekształć. Wypełniamy jak na rysunku 6.5.

Uwaga: Jamovi nie zmienia wartości zmiennej s1 tylko utworzy nową zmienną z przekodowanymi wartościami. Zmienna na podstawie której jest tworzona nowa zmienna nazywa się źródłową (s1 w naszym przykładzie jest źródłową).



Rysunek 6.5: Przekształcenie

Wpisujemy sensowną nazwę (na przykład s1p od przekodowana). Jak będziemy używać sensownych nazwa łatwiej będzie nam się pracowało. Dobrze jest też podać w opisie co zawiera zmienna.

Klikamy w pole wyboru na dole (obok napisu za pomocą przekształcenia). Powinniśmy zobaczyć coś podobnego do tego co widać na rysunku 6.6.

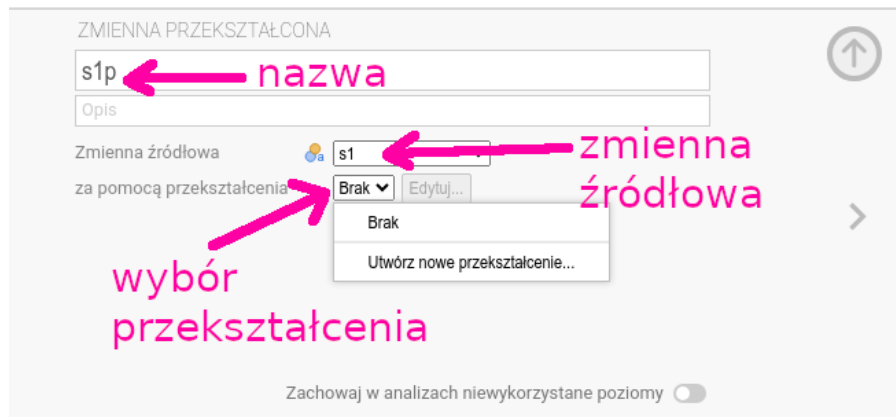
Wybieramy Utwórz nowe przekształcenie. Wpisujemy sensowną nazwę przekształcenia (na przykład Likert2R5) oraz formułę przekształcenia:

```
IF ($source=="Zdecydowanie się nie zgadzam", 1,
  IF ($source=="Nie zgadzam się", 2,
    IF ($source=="Trudno powiedzieć", 3,
      IF ($source=="Zgadzam się", 4, 5))))
```

Formuła może wydawać się przerażająca, ale jest koncepcyjnie bardzo prosta:

```
IF (warunek, jeżeli-prawda, jeżeli-fałsz)
```

Warunek to fragment `$source=="Zdecydowanie się nie zgadzam"`:



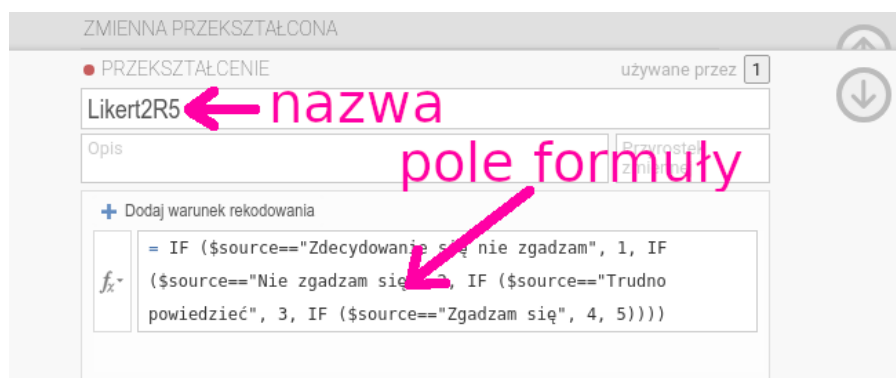
Rysunek 6.6: Przekształcenie

- `$source` oznacza bieżącą wartość zmiennej źródłowej;
- `==` to **operator** równości; jest więcej operatorów, które można wybrać z menu;
- `$source=="Zdecydowanie się nie zgadzam"` oznacza, że jeżeli bieżącą wartością w kolumnie źródłowej jest Zdecydowanie się nie zgadzam to wykonaj jeżeli-prawda; w wypadku przeciwnym wykonaj jeżeli-fałsz.

jeżeli-prawda to zwykle wstawienie nowej wartości; jeżeli-fałsz to często następna formuła IF albo wstawienie innej nowej wartości. Przykładowo jeżeli bieżącą wartością w kolumnie źródłowej jest Zdecydowanie się nie zgadzam to wstaw 1, jeżeli nie jest to wstaw 0:

```
IF ($source=="Zdecydowanie się nie zgadzam", 1,0)
```

Ponieważ w naszym przykładzie mamy do przekodowania nie dwie a 5 wartości musimy użyć 4 warunków, które są zagnieżdżone jeden w drugim. Można powyższe przepisać, można też skopiować z podręcznika i wkleić do Jamovi (por. rys. 6.7).

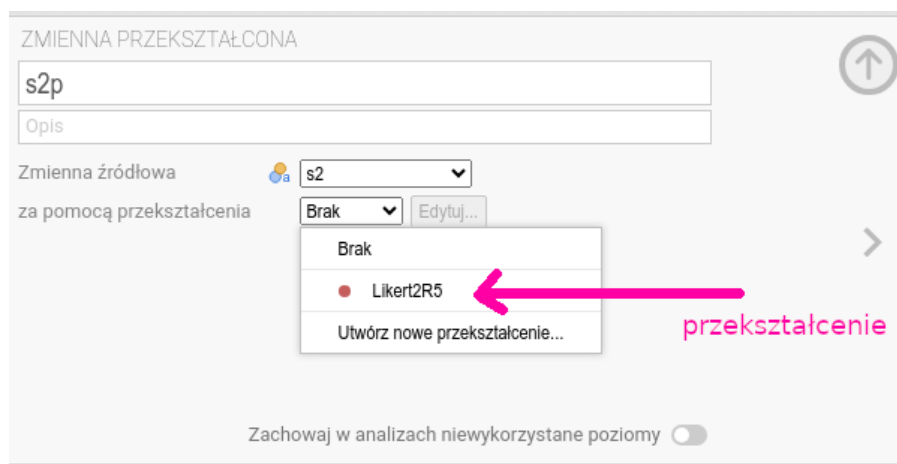


Rysunek 6.7: Przekształcenie

Naciskamy Enter i gotowe. Zostaje utworzona zmienna s1p zawierająca zamiast

napisów rangi.

Jeżeli uporaliśmy się z przekodowaniem s1 ustawiamy kursor na s2 w okno danych. Naciskamy ikonę Przekształć. Upewniamy się, że zmienną źródłową jest s2. Zmieniamy nazwę nowej zmiennej na s2p. Klikamy w pole wyboru przekształcenia. Poprzednio były tam tylko dwie pozycje Brak oraz Utwórz nowe przekształcenie teraz jest trzecia pozycja Likert2R5 czyli przekształcenie, które zdefiniowaliśmy dla zmiennej s1p. Wybieramy Likert2R5 bo zmienną s2 chcemy przekodować dokładnie w ten sam sposób jak s1. Po wybraniu przekształcenia w oknie danych pojawia się nowa zmienna s2p (por. rys. 6.8).



Rysunek 6.8: Przekształcenie

W podobny łatwy sposób przekodowujemy s3 oraz zo1, zo2, zo3.

**Uwaga:** polecenie IF wpisujemy używając dużych liter. Słowo \$source wpisujemy tak jak jest to zademonstrowane (\$Source jest błędem).

Przekodowanie pytania z możliwością wielokrotnego wyboru jest równie proste, tyle że pisanie jest więcej. Zmienna wiedza\_na\_temat\_palenia może zawierać do ośmiu następujących napisów oddzielonych średnikami: Przewlekła obturacyjna choroba płuc, Astma oskrzelowa, Alergie wziewne, Gruźlica (B), Zapalenie płuc (B), Przewlekłe zapalenie oskrzeli, Infekcje dróg oddechowych, Palenie nie powoduje chorób układu oddechowego (B).

Odpowiedzi błędne oznaczono jako (B).

W arkuszu lub oknie danych Jamovi ta zmienna wygląda jakoś tak:

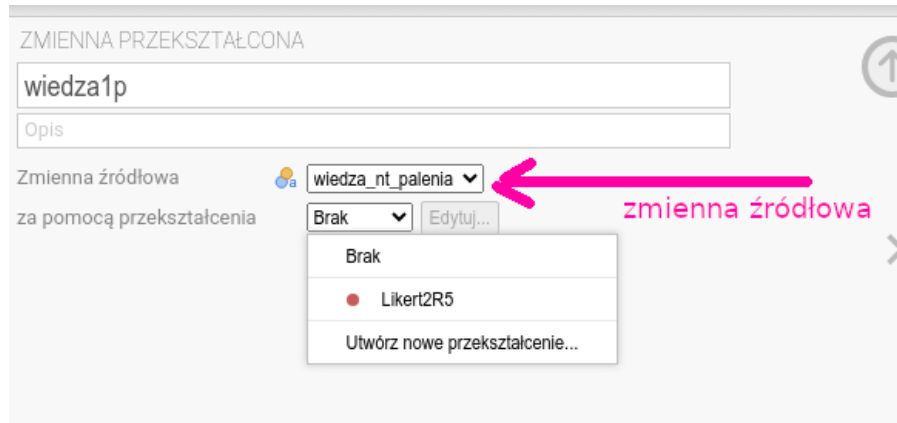
```
...,Przewlekła obturacyjna choroba płuc,...
...,Przewlekła obturacyjna choroba płuc;Astma oskrzelowa;...
...,Astma oskrzelowa,
...,Astma oskrzelowa;Gruźlica;Przewlekłe zapalenie oskrzeli,...
...,Przewlekła obturacyjna choroba płuc;Astma oskrzelowa;...
```

Należy zsumować wystąpienia poprawne i wystąpienia błędne. W tym celu trzeba utworzyć tyle nowych zmiennych ile jest wariantów odpowiedzi, czyli w naszym przykładzie osiem. Każda nowa zmienna jest przekodowywana za pomocą prostej formuły

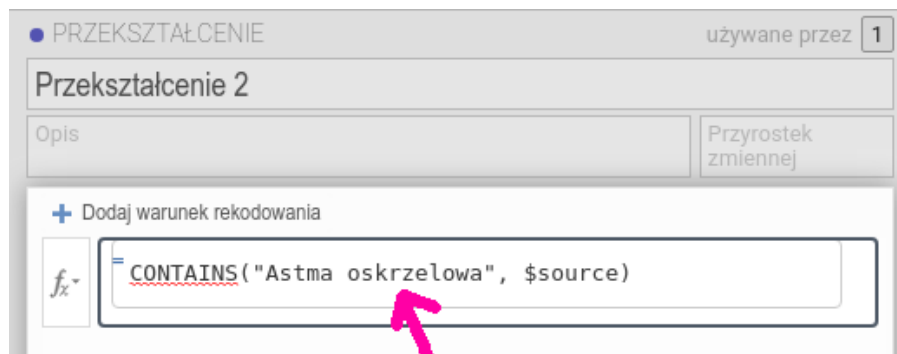
wykorzystującą funkcję CONTAINS (zawiera). Przykładowo pierwsza (nazwijmy ją wiedza1p) powinna być utworzona w oparciu o następujące przekształcenie:

```
CONTAINS("Przewlekła obturacyjna choroba płuc", $source)
```

Jak to wygląda w oknie programu Jamovi przedstawiono na rysunkach 6.9 oraz 6.10.



Rysunek 6.9: Przekształcenie



Rysunek 6.10: Przekształcenie

Funkcja CONTAINS wstawi 1 jeżeli \$source zawiera Przewlekła obturacyjna choroba płuc. Oczywiście następna zmienna powinna zawierać Astma oskrzelowa:

```
CONTAINS("Asthma oskrzelowa", $source)
```

I tak dalej aż do ostatniego wariantu odpowiedzi:

```
CONTAINS('Alergie wziewne', $source)
CONTAINS('Gruźlica', $source)
CONTAINS('Zapalenie płuc', $source)
CONTAINS('Przewlekłe zapalenie oskrzeli', $source)
```



```
CONTAINS('Infekcje dróg oddechowych', $source)
CONTAINS('Palenie nie powoduje chorób układu oddechowego', $source)
```

Każda zmienna wiedza1...wiedza8 zawiera 1 jeżeli ankietowany wskazał dany wariant lub zero jeżeli nie wskazał.

Ostatnia sprawa to przekodowanie liczb na wartości nominalne. Przykładowo chcemy podzielić ankietowanych na grupy stażowe: mały (do pięciu lat), średni (5–15 lat), duży (16 i więcej) staż pracy.

Wartości liczbowe stażu pracy zawiera zmienna staz. Aby ją przekodować należy użyć następującego przekształcenia:

```
IF ($source < 5, "M",
    IF ($source < 16, "S", "D"))
```

Polecenie IF musi być o jedno mniej niż mamy klas. W naszym przykładzie zatem dwa. Jeżeli staż jest mniejszy od 5 wstawiony zostanie napis M, jeżeli staż jest mniejszy od 16 wstawiony zostanie napis S a w przeciwnym wypadku zostanie wstawiony napis D.

Gdyby ktoś się niepokoił że 3 spełnia jednocześnie  $\$source < 5$  oraz  $\$source < 16$  to dodamy, że pierwszy się liczy. Przekształcenie kończy działanie po spełnieniu pierwszego warunku i nie wykonuje dalszych porównań. Dlatego liczba 3 zostanie zamieniona na M a nie na S.

Podobnie przekodujemy zmienną wiek.

### 6.2.3 Wyliczenie nowych zmiennych

**Przekodowanie** to była w zasadzie zamiana sposobu mierzenia. **Wyliczenie** to utworzenie nowej zmiennej, zwykle w oparciu o jakąś formułę matematyczną. Na przykład odwrócenie pytanie s1p realizuje  $s1pr = 6 - s1p$ . Satysfakcja to suma rang z trzech pytań:  $satysfakcja = s1pr + s2p + s3p$ .

W celu wyliczenie nowych zmiennych należy wybrać Dane Oblicz. Pojawia się okno zmiennej wyliczonej zatytułowane ZMIENNA WYLICZONA.

Pierwszy pasek zawiera nazwę zmienną (domyślnie nazwę kolumny w konwencji arkusza kalkulacyjnego, w przykładzie jest to litera H) W polu definiowania zmiennej należy wpisać stosowną formułę matematyczną. W przypadku odwracania pytania s1p będzie to:

$6 - s1p$

W przypadku liczenia łącznej satysfakcji (por. rys. 6.11):

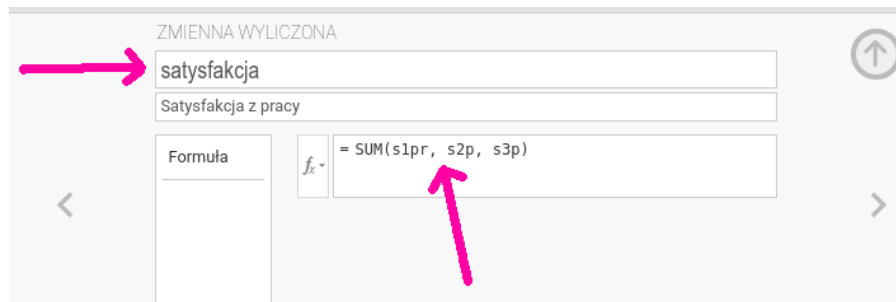
$SUM(s1pr, s2p, s3p)$

Oczywiście wcześniej musimy utworzyć zmienną s1pr, inaczej Jamovi zgłosi błąd. Jeżeli nie chcemy sumy, ale np. średnią powinniśmy użyć:

$MEAN(s1pr, s2p, s3p)$

Inne funkcje matematyczne są dostępne po kliknięciu w pole wyboru znajdujące się po lewej stronie pola definiowania zmiennej.

Powiedzieliśmy, że miarą wiedzy nt. palenia będzie suma punktów uzyskanych za odpowiedzi prawidłowe minus suma punktów uzyskanych za odpowiedzi nieprawidłowe. Odpowiedzi prawidłowe to w1p, w2p, w3p, w6p oraz w7p. Odpowiedzi błędne to w4p, w5p, w8p. Zatem w polu definiowania zmiennej wpisujemy:



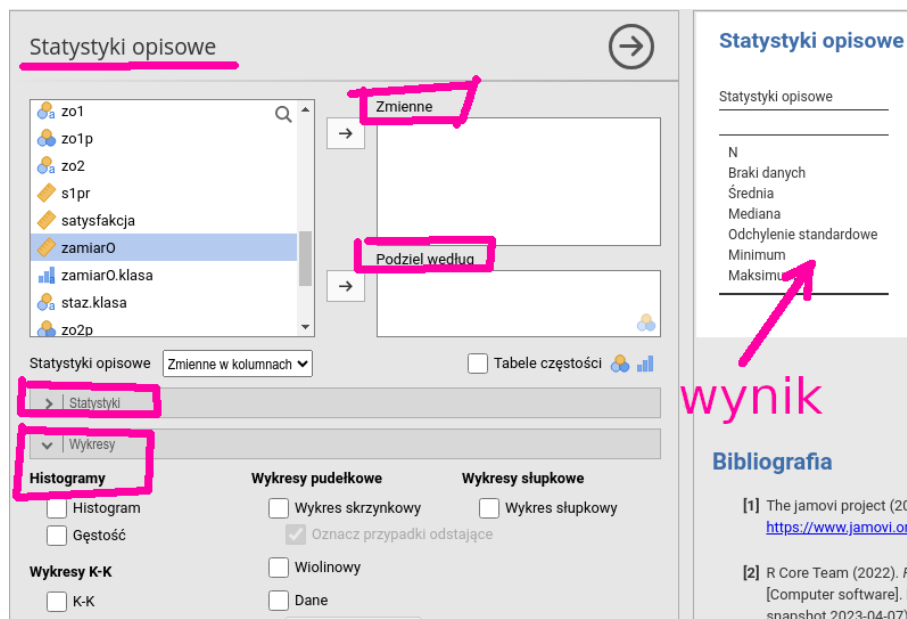
Rysunek 6.11: Obliczanie nowej zmiennej

$SUM(w1p, w2p, w3p, w6p, w7p) - SUM(w4p, w5p, w8p)$

### 6.2.4 Analiza struktury

Wybieramy Analizy → Eksploracja → Statystyki opisowe.

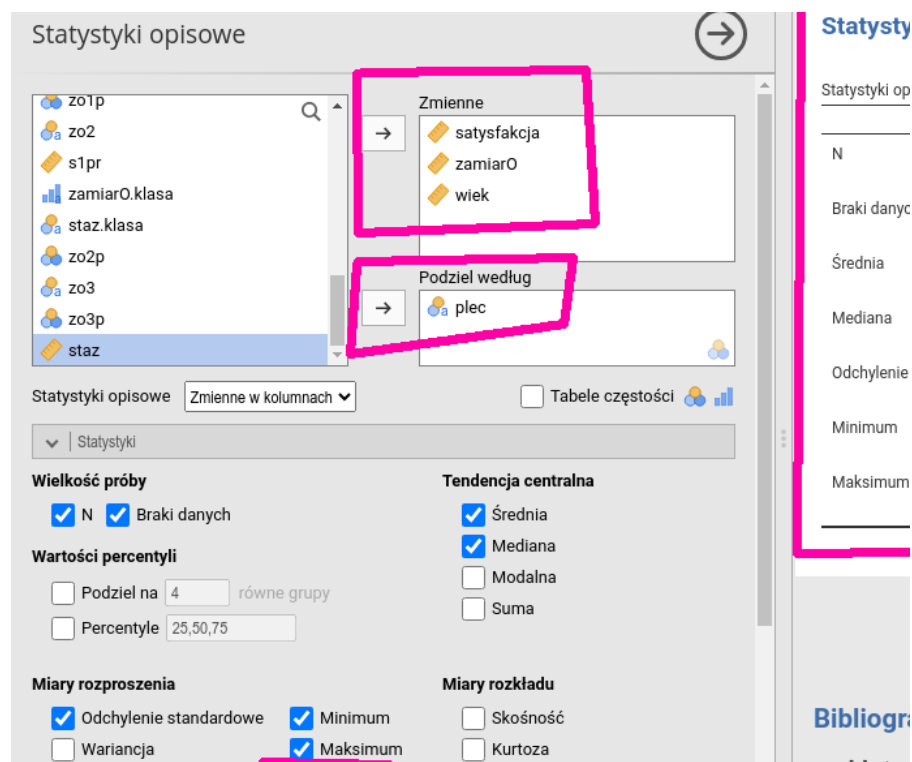
W wyświetlonym oknie po lewej deklarujemy co ma być liczone. Wynik pojawi się po prawej (por. rys. 6.12).



Rysunek 6.12: Statystyki opisowe

Ustawiamy kursor na zmiennej, która nas interesuje i klikamy w strzałkę górną. Jeżeli chcemy podzielić wartości zmiennej na grupy według jakiejś zmiennej nominalnej, to ustawiamy kursor na tej zmiennej nominalnej (na przykład plec) i klikamy strzałkę dolną.

Można analizować wiele zmiennych na raz (por. rys. 6.13). Wystarczy w tym celu ustawić kursor na zmiennej i kliknąć w odpowiednią strzałkę. Zawartość okna wynikowego zostanie automatycznie uaktualniona.



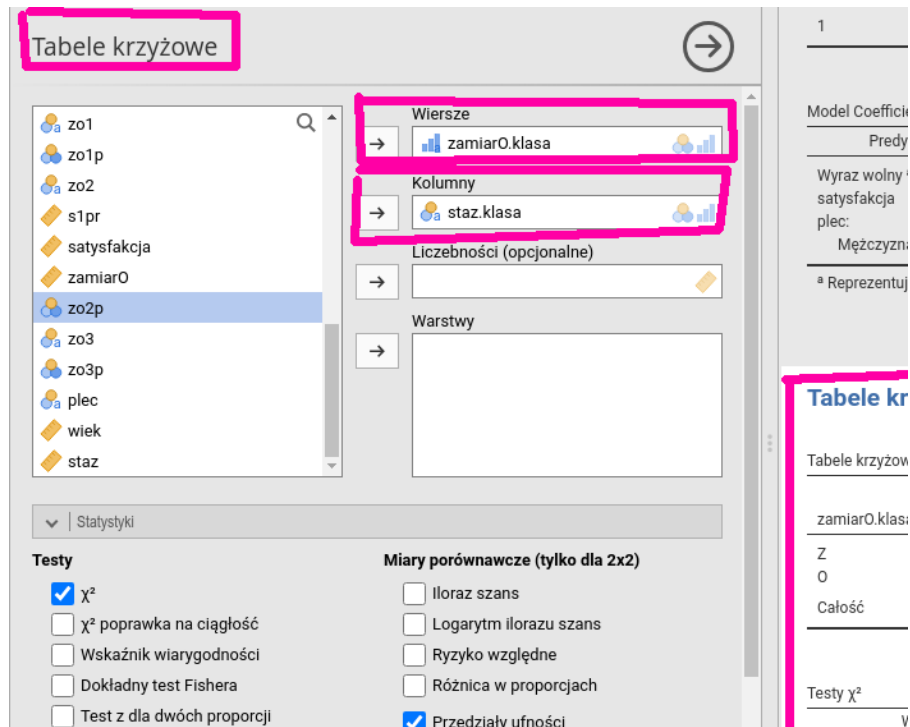
Rysunek 6.13: Statystyki opisowe

Poniżej okien wyboru zmiennych są zakładki określające precyzyjnie, co ma być obliczone oraz jakie wykresy mają zostać wyrysowane. Przykładowo domyślny wydruk nie zawiera rozstępu ćwiartkowego. Żeby go dodać do wyniku należy w zakładce Statystyki zaznaczyć przycisk IQR.

### 6.2.5 Analiza zależności: zmienne nominalne

Wybieramy Analizy→Częstości→Próby niezależne.

Podobnie jak w przypadku analizy struktury jest wyświetlana lista zmiennych oraz okna i strzałki pozwalające wygodnie wybrać to co ma być analizowane. Jest to tak proste że wystarczy przyjrzeć się przykładowemu rysunkowi żeby wiedzieć jak postępować. Przykładową analizę zależności pomiędzy zmiennymi nominalnymi zamiar0klasa oraz staz.klasa przedstawia rysunek 6.14.



Rysunek 6.14: Tabele krzyżowe

### 6.2.6 Analiza zależności: zmienna liczbowa/zmienna nominalna

Jeżeli zmienna nominalna (zwana grupującą) przyjmuje dwie wartości wybieramy Analizy→Testy t→Test t dla prób niezależnych. Zmienne zależne to zmienne, których wartości zostaną podzielone na grupy. Ustawiamy kursor kolejno na zmiennej, która nas interesuje i klikamy w strzałkę górną. Ustawiamy kursor na zmiennej grupującej (na przykład *plec*) i klikamy strzałkę dolną (por. rys. 6.15).

Zaznaczamy przyciski Test Welcha, U Manna-Whitneya (w sekcji **Testy**) oraz Test normalności (w sekcji **Weryfikacja założeń**)

Jeżeli zmienna nominalna więcej niż dwie wartości wybieramy Analizy→ANOVA→Jednoczynnikowa ANOVA. Zmienne zależne i grupujące wybieramy w identyczny sposób jak w przypadku testu Welcha (por. rys. 6.16). Zaznaczamy przyciski Nie zakładaj równości (Test Welcha) (w sekcji **Wariancje**) oraz Test normalności (w sekcji **Weryfikacja założeń**) i Tabela statystyk opisowych (w sekcji **Dodatkowe statystyki**).

Jeżeli wynik testu Shapiro-Wilka wskaże, że rozkład zmiennej zależnej nie jest normalny należy wykonać test Kruskala-Wallisa wybierając Jednoczynnikowa ANOVA Test Kruskala-Wallisa.

Rysunek 6.15: Test dla prób niezależnych

### 6.2.7 Analiza zależności: zmienna liczbowa/zmienna liczbowa lub nominalna

Wybieramy Analizy→Regresja→Regresja liniowa.

Interfejs jest podobny do poprzednio opisywanych. Wybieramy zmienną zależną (musi oczywiście być liczbowa) klikając w górną strzałkę. Zmienne niezależne mierzone w skali liczbowej klikając w środkową strzałkę. Zmienne niezależne mierzone w skali nominalnej klikając w dolną strzałkę. Wynik automatycznie pojawia się w lewym oknie (por. rys. 6.17).

### 6.2.8 Regresja logistyczna

Wybieramy Analizy→Regresja→Regresja logistyczna→Dwie wartości.

Interfejs jest ładząco podobny do analizy regresji. Wybieramy zmienną zależną klikając w górną strzałkę. Zmienna ta **musi** być zmienną dwuwartościową.

Zmienne niezależne mierzone w skali liczbowej wybieramy klikając w środkową strzałkę, a zmienne niezależne mierzone w skali nominalnej klikając w dolną strzałkę.

**Jednoczynnikowa ANOVA**

Zmienne zależne  
 satysfakcja

Zmienna grupująca  
 staz.klasa

**Wariancje**  
☒ Nie zakładaj równości (test Welch)  
☐ Załóż równość (test Fishera)

**Dodatkowe statystyki**  
☒ Tabela statystyk opisowych  
☐ Wykresy opisowe

**Braki danych**  
☒ Sprawdzaj dla każdej analizy oddzielnie  
☐ Wyklucz obserwacje ze wszystkich analiz

**Weryfikacja założeń**  
☐ Test homogeniczności  
☒ Test normalności  
☐ Wykres K-K

**ANOVA - satysfakcja**

	Suma k
staz.klasa	
Reszty	

**Jednoczynnik**

Jednoczynnikowa ANCOVA

	F
staz.klasa	0.27

Statystyki opisowe dla

	staz.kl
satysfakcja	M S

Rysunek 6.16: Jednoczynnikowa ANOVA

**Regresja liniowa**

Zmienna zależna  
 zmiarO

Współzmiennie  
 satysfakcja

Czynniki  
 plec

Wagi (opcjonalnie)

**Założenia**

Test normalności

satysfakcja

Uwaga. Niska wartość naruszenie założeń normalności rozkładu

**Regresja liniowa**

Miary dopasowania modelu

Model	R
1	0.601

Rysunek 6.17: Regresja liniowa

### 6.2.9 Redagowanie raportu

Zwykle dobrze jest dodać jakieś dodatkowe objaśnienia do wyników obliczeń wygenerowanych przez program i Jamovi nam to umożliwia. Wybierając Edytuj przechodzimy do prostego edytora umożliwiającego redagowanie raportu w oknie wyników a obsługa tego menu jest tak banalnie prosta, że nie wymaga jakiś specjalnych objaśnień.

# Literatura

Beaglehole, Robert, Ruth Bonita, and Tord Kjellstrom. 1996. *Podstawy Epidemiologii*. Instytut Medycyny Pracy.

Bland, Martin. 2015. *An Introduction to Medical Statistics*. Oxford University Press.

Machin, David, Michael J Campbell, and Stephen J Walters. 2007. *Medical Statistics*. John Wiley.