

Uniwersytet Wrocławski
Wydział Matematyki i Informatyki
Instytut Matematyczny
specjalność: Analiza Danych

Kamil Zaborniak

Optymalny próg wykrycia w modelach typu:
„igła w stogu siana”

Praca licencjacka
napisana pod opieką
dr. Liudmyly Zaitsevy

Wrocław, 2023

Spis treści

1	Wstęp	3
2	Podstawy do teorii testowania hipotez	4
2.1	Wprowadzenie	4
2.2	Procedura weryfikacji hipotez statystycznych	6
2.3	Lemat Neymana-Pearsona	9
3	Rozkład p-wartości przy założeniu hipotezy zerowej	11
3.1	Rozkład p -wartości opartej na ciągłej statystyce testowej	11
3.2	Rozkład p -wartości opartej na dyskretnej statystyce testowej	12
4	Globalna hipoteza zerowa oraz sposób jej weryfikacji	14
4.1	Globalna hipoteza zerowa	14
4.2	Sposób testowania globalnej hipotezy zerowej	15
5	Modele typu „igła w stogu siana” oraz optymalne progi ich wykrycia	16
5.1	Model oparty o dane pochodzące z rozkładu normalnego	17
5.1.1	Optymalny próg wykrycia	17
5.1.2	Asymptotyczna moc korekty Bonferroniego	20
5.1.3	Asymptotyczna moc testu opartego o lemat Neymana-Pearsona	20
5.2	Model oparty o dane pochodzące z rozkładu Poissona	23
5.2.1	Asymptotyczna wartość kwantyla	24
5.2.2	Asymptotyczna moc korekty Bonferroniego	26
6	Wnioski	29

1 Wstęp

Głównym celem niniejszej pracy jest znalezienie przez nas optymalnych progów wykrycia dla modeli typu „igła w stogu siana”. Uzyskane wyniki przeanalizujemy i sprawdzimy jak wpływają na optymalizację procesów testowania wcześniej wspomnianych modeli statystycznych.

Rozważania rozpoczniemy od poznania i przypomnienia podstaw teorii testowania hipotez statystycznych. Rozważymy czym są hipotezy statystyczne oraz w jaki sposób prawidłowo je konstruować. W kolejnym kroku określimy typy rozpatrywanych sądów, ze względu na to do czego się odnoszą, jaki kierunek obierają oraz jak złożone są. Przemyślenia te pozwolą nam poprawnie łączyć je w pary, czyli tworzyć modele statystyczne. Następnym krokiem będzie poznanie przez nas idei ich testowania. Uzyskamy wiedzę na temat sposobów konstrukcji oraz klasyfikacji testów statystycznych ze względu na postać modelu statystycznego. Zapoznamy się z rodzajami błędów, które możemy uzyskać w wyniku weryfikacji hipotez statystycznych, oraz ich prawdopodobieństwami. Po przeanalizowaniu podstaw, poznamy narzędzia testów oraz postaramy się przeprowadzić kilka prostych weryfikacji przykładowych modeli statystycznych, a następnie wyznaczmy ich moce. Po tym, poznamy metodę konstrukcji testu najmocniejszego, opartą o lemat Neymana-Pearsona, oraz rozpatrzmy przykład jej wykorzystania.

W następnym kroku postaramy się odpowiedzieć na pytanie: „*Jaki rozkład przyjmuje p -wartość?*”. Rozpatrzmy przypadki w zależności od rodzaju rozkładu prawdopodobieństwa statystyki testowej. Problem rozwiążemy przy pomocy własności dystrybuanty poszczególnych rozkładów. W rezultacie dowiemy się, że jeśli statystyka testowa pochodzi z rozkładu ciągłego, to rozkład jej przekształcenia jest rozkładem jednostajnym na przedziale $[0, 1]$. Ustalimy również rozkład p -wartości w sytuacji, kiedy statystyka testowa będzie pochodziła z rozkładu Poissona. Wtedy, dowiemy się, że p -wartość przyjmuje asymptotycznie rozkład $\mathcal{U}[0, 1]$.

Następnym zagadnieniem, które poznamy, będzie globalna hipoteza zerowa. Zapoznamy się ze sposobem jej konstrukcji. Dowiemy się jak korekta Bonferroniego pozwala nam sprawdzić prawdziwość globalnej hipotezy zerowej. W konsekwencji przeprowadzonych rozważań na temat rozkładu p -wartości, udowodnimy, że korekta Bonferroniego kontroluje prawdopodobieństwo odrzucenia globalnej hipotezy zerowej na ustalonym poziomie istotności $\alpha \in (0, 1)$.

Konsekwencją przeprowadzonych przez nas rozważań będzie definicja modelu typu „igła w stogu siana”. Pierwszym takim modelem, który poddamy analizie będzie model oparty o dane pochodzące z rozkładu normalnego o parametrze średniej — μ . Model będzie się składał z globalnej hipotezy zerowej: $H_0 : \forall i \in \{1, \dots, n\} : \mu_i = 0$, gdzie n określa liczbę prób, na podstawie, których dokonamy weryfikacji, oraz alternatywy: $H_1 : \exists! i \sim \mathcal{U}\{1, \dots, n\} : \mu_i > 0$. Wyznamy asymptotyczną wartość kwantyla rozkładu normalnego, czyli asymptotycznego progu wykrycia. W oparciu o to wyznaczmy asymptotyczną moc korekty Bonferroniego. Następnie przekształcimy model i wyznaczmy asymptotyczną moc testu opartego o lemat Neymana-Pearsona. Wynikiem przeprowadzonych testów będzie fakt, że będziemy mogli ustalić czy pośród obserwacji z rozkładu $\mathcal{N}(0, 1)$, znajduje się obserwacja z rozkładu $\mathcal{N}(\mu, 1)$, wyłącznie gdy $\mu(1 + \varepsilon)\sqrt{2 \ln(n)}$, gdzie $\varepsilon \in (0, 1)$. Sprawdzimy również przypadek, kiedy model typu „igła w stogu siana” będzie oparty o dane pochodzące z rozkładu Poissona o parametrze λ . Analizowany model będzie się składał z hipotez: $H_0 : \forall i \in \{1, \dots, n\} : \lambda_i = \lambda'$, $H_1 : \exists! i \sim \mathcal{U}\{1, \dots, n\} : \lambda_i > \lambda'$. Wyaprosymujemy wartość asymptotycznego progu wykrycia, oraz postaramy się sprawdzić jak korekta Bonferroniego pozwala nam wykryć obserwację odstającą oraz jaką moc osiąga, przy liczbie prób: $n \rightarrow \infty$. W rezultacie dowiemy się, że korekta Bonferroniego nie wykryje obserwacji odstającej, pochodzącej z rozkładu Poissona, jeśli jej parametr rozkładu będzie równy $(1 - \varepsilon) \frac{\ln(n)}{\ln(\ln(n))}$.

2 Podstawy do teorii testowania hipotez

Jednym z głównych działów statystyki matematycznej jest wnioskowanie statystyczne. Obejmuje ono zasady i metody uogólniania wyników otrzymanych poprzez badanie próby losowej na całą populację generalną, z której pochodzi dana próba. Głównymi domenami wnioskowania statystycznego są:

Estymacja — szacowanie postaci lub parametrów rozkładu zmiennej losowej w populacji generalnej na podstawie próby badawczej,

Weryfikacja hipotez statystycznych — sprawdzenie przypuszczeń dotyczących rozkładu populacji generalnej poprzez badanie jej części.

W przypadku estymacji, wnioski na temat populacji generalnej stwierdzamy w oparciu o wyniki próby. W odróżnieniu, przeprowadzając testowanie hipotez statystycznych, najpierw wysuwamy określone sądy na temat zbiorowości generalnej, a w następnych krokach sprawdzamy ich prawdziwość. W niniejszym rozdziale, opierając się o wiedzę zawartą w książce pt. „Introduction to Mathematical Statistics (6th Edition)”^[1], przeprowadzimy głównie rozważania na temat teorii weryfikacji hipotez.

2.1 Wprowadzenie

Hipotezami statystycznymi nazywamy założenia dotyczące rozkładu prawdopodobieństwa badanej cechy populacji, która jest zmienną losową. Na nich oparty jest proces ich weryfikacji. Hipotezy statystyczne dzielimy na parametryczne i nieparametryczne. **Hipotezy nieparametryczne** dotyczą postaci funkcji gęstości rozkładu parametru populacji, którego badamy. Ich przykładem jest hipoteza: „Rozkład zmiennych losowych pochodzących z populacji jest rozkładem wykładniczym.”. **Hipotezy parametryczne** dotyczą wartości parametrów rozkładu cechy populacji generalnej. Przykładem takiej hipotezy jest założenie: „Średnie dwóch populacji o rozkładzie normalnym są równe.”. Przypuszczenia te dzielimy na hipotezy niekierunkowe oraz hipotezy kierunkowe. **Hipoteza niekierunkowa** zakłada tylko różnice wielkości badanego parametru względem konkretnej wartości. Przykładem jest hipoteza: „Odsetek osób palących w populacji generalnej jest różny od 0.2”. Natomiast **hipoteza kierunkowa** określa kierunek spodziewanych wyników. Taka hipoteza może mieć formę lewostronną lub prawostronną. **Hipoteza lewostronna** zakłada, że badany parametr jest mniejszy od konkretnej wartości. Jej przykładem jest hipoteza: „Odsetek osób palących w populacji generalnej jest mniejszy niż 0.2”. **Hipoteza prawostronna** stwierdza natomiast, że wartość badanego parametru należy do zbioru wartości większych od konkretnej wartości. Przykładem takiej hipotezy jest założenie: „Odsetek osób palących jest większy od 0.2”. Rodzinę rozkładów lub zbiór wszystkich wartości badanego parametru, określanych przez wszystkie możliwe hipotezy statystyczne, dotyczące tego parametru, będziemy oznaczać symbolem Θ .

Każda wiedza dotycząca populacji generalnej pozwala nam ograniczyć zbiór możliwych hipotez do **zbioru hipotez dopuszczalnych**. Przykładowo, jeśli wiemy, że badana cecha populacji — zmienna losowa, pochodzi z rozkładu Poissona, to naszym zbiorem hipotez dopuszczalnych będzie zbiór wszystkich rozkładów Poissona różniących się wyłącznie parametrem λ . Bezcelowe jest w takim przypadku testowanie hipotezy, czy dana cecha ma rozkład inny niż rozkład Poissona.

Jeśli hipoteza ze zbioru hipotez dopuszczalnych jednoznacznie określa rozkład populacji, taką hipotezę nazywamy **hipotezą prostą**. Przykładem jest hipoteza: „Populacja pochodzi z rozkładu Poissona o parametrze λ równym 2.”. Jakakolwiek hipoteza, która nie jest hipotezą prostą jest **hipotezą złożoną**, na przykład: „Populacja pochodzi z rozkładu Poissona o parametrze λ większym od 2.”.

Założenie dotyczące rozkładu zmiennej losowej X w populacji generalnej, którego prawdziwość chcemy sprawdzić nazywamy **hipotezą zerową** i oznaczamy symbolem H_0 . Zazwyczaj hipoteza zerowa będzie określała założenie, które wynika z podstawowej wiedzy jaką posiadamy na temat rozkładu badanej populacji. W przypadku badania związku między dwiema populacjami, za hipotezę zerową będziemy stawiać założenie, które nie zakłada żadnych różnic między nimi. Zaprzeczeniem hipotezy zerowej jest **hipoteza alternatywna**, którą oznaczamy symbolem H_1 . Określa ona wszystkie lub część hipotez ze zbioru hipotez dopuszczalnych nie będących hipotezą zerową. W wyniku odrzucenia hipotezy zerowej, możemy uznać hipotezę alternatywną za prawdziwą. W przypadku hipotezy zerowej, zbiór wartości parametru lub rodzinę rozkładów, które ona zakłada, będziemy oznaczać symbolem Θ_0 . W przypadku hipotezy alternatywnej — symbolem Θ_1 .

W podjęciu decyzji o odrzuceniu hipotezy zerowej pomaga nam test statystyczny. **Testem statystycznym** jest ściśle określone postępowanie, w wyniku którego, dla każdej wylosowanej próby z populacji, podejmujemy decyzję o odrzuceniu hipotezy zerowej dotyczącej całej populacji. Test ten oznaczamy symbolem φ . Wyróżniamy dwa rodzaje testów statystycznych:

Testy parametryczne służą do testowania hipotez parametrycznych. W celu ich przeprowadzenia potrzebna jest znajomość rozkładu badanego parametru w populacji generalnej. Przykładami takich testów są m.in. test Studenta.

Testy nieparametryczne służą do weryfikacji hipotez nieparametrycznych. Przykładami takich testów są: test Kołmogorowa-Smirnowa oraz test Shapiro-Wilka.

W celu prawidłowego przeprowadzenia testu statystycznego próba musi być losowo wybrana z populacji generalnej, a jej obserwacje nie mogą być od siebie zależne.

Biorąc pod uwagę fakt, że próba jest wybierana w sposób losowy, ten sam test statystyczny może prowadzić do podjęcia dwóch różnych decyzji dla dwóch różnych prób, przez co możemy popełnić jeden z dwóch rodzajów błędów. **Błędem I rodzaju** nazywamy decyzję o odrzuceniu hipotezy zerowej, wtedy kiedy jest ona prawdziwa. **Błąd II rodzaju** popełniamy natomiast, kiedy przyjmujemy za prawdziwą hipotezę zerową, kiedy w rzeczywistości jest ona fałszywa.

Ważnymi elementami przeprowadzenia weryfikacji hipotez są prawdopodobieństwa popełnienia powyższych błędów. Ryzyko popełnienia błędu I rodzaju, nazywamy **poziomem istotności** oraz oznaczamy symbolem α , a prawdopodobieństwo popełnienia błędu II rodzaju oznaczamy symbolem β . Prawdopodobieństwo wystąpienia zdarzenia przeciwnego do sytuacji, w której popełniamy błąd II rodzaju, nazywamy **mocą testu**. Innymi słowy moc testu określa prawdopodobieństwo przyjęcia hipotezy alternatywnej, kiedy jest ona prawdziwa. Współczynnik ten oznaczamy symbolem γ .

Niestety, nie możemy skonstruować testu, który jednocześnie minimalizowałby współczynniki α i β . W celu udowodnienia tego, założmy, że przyjmujemy $\alpha = 0$ i chcemy kontrolować współczynnik β na poziomie 0. Z czego wynika, że zawsze przyjmujemy hipotezę zerową, nawet kiedy jest ona fałszywa. przez co ryzyko popełnienia błędu II rodzaju jest większe od 0 i przeczy wcześniejszemu założeniu, że $\beta = 0$.

W celu prawidłowego przeprowadzenia weryfikacji hipotez statystycznych, powinniśmy najpierw ustalić poziom istotności. Zazwyczaj, będziemy ustalać jego wartość na poziomie 0.10, 0.05. Bezpieczniejsze jest przeprowadzenie przez nas testu hipotez statystycznych w oparciu o ryzyko popełnienia błędu I rodzaju, niż w oparciu o ryzyko popełnienia błędu II rodzaju. Rozpatrzmy sytuację, w której chcemy sprawdzić skuteczność nowo opracowanego leku zmniejszającego ciśnienie krwi, przeznaczonego dla osób chorych na nadciśnienie. Za hipotezę zerową przyjmujemy założenie, że lek nie zmniejsza ciśnienia krwi, zaś za alternatywę — założenie, że lek zmniejsza ciśnienie krwi. Skutkiem podjęcia decyzji o odrzuceniu hipotezy zerowej na rzecz hipotezy alternatywnej może być dopuszczenie leku do sprzedaży. Korzystniejsze jest w tym przypadku popełnienie błędu II rodzaju, czyli

stwierdzenia, że lek nie działa, kiedy w rzeczywistości jest on skuteczny, czego konsekwencją jest niedopuszczenie leku do dystrybucji. Z kolei przyjęcie skuteczności leku, kiedy w rzeczywistości on nie działa, czyli popełnienie błędu I rodzaju, wiąże się z ryzykiem podania tego leku osobie chorej w sytuacji zagrażającej jej życiu. Sytuacja ta pokazuje, że bezpieczniejsze jest, podczas przeprowadzania testowania hipotez statystycznych, kontrolowanie ryzyka popełnienia błędu I rodzaju.

Jako że dla jednej pary hipotez oraz ustalonego poziomu istotności możemy skonstruować wiele różnych testów statystycznych, decyzję o odrzuceniu hipotezy zerowej podejmujemy w oparciu o test, który ma największą moc. Test, który spośród innych testów na tym samym poziomie istotności α , przyjmuje największą moc, nazywamy **testem jednostajnie najmocniejszym** na poziomie istotności α . Niestety, nie zawsze taki test istnieje.

2.2 Procedura weryfikacji hipotez statystycznych

Proces testowania hipotez statystycznych rozpoczynamy od określenia hipotezy zerowej, dotyczącej rozkładu badanego parametru populacji generalnej, oraz jej alternatywy. Układ tych hipotez nazywamy modelem statystycznym. W kolejnym kroku przyjmujemy poziom istotności na podstawie, którego testowanie ma przebiegać. Następnie losujemy n -elementową próbę obserwacji z populacji generalnej. Zazwyczaj taką próbę oznaczamy symbolem \underline{X} , a i -tą obserwację symbolem x_i . W celu poznania kolejnych kroków procedury weryfikacji hipotez, spróbujmy przeprowadzić ją na poniższym przykładzie.

Przyjmijmy, że dysponujemy próbą \underline{X} złożoną z 16 niezależnych zmiennych losowych, pochodzących z rozkładu normalnego, o nieznanym średniej μ oraz wariancji $\sigma^2 = 1$. Wiemy, że średnia arytmetyczna tej próby — $\hat{\mu}$ wynosi 0.4375. Chcemy sprawdzić, na poziomie istotności $\alpha = 0.05$, czy parametr μ jest równy 0. W tym celu tworzymy układ hipotez statystycznych, które chcemy przetestować:

$$\begin{aligned} H_0 : \quad & \mu = 0, \\ H_1 : \quad & \mu \neq 0. \end{aligned} \tag{A}$$

Kolejnym krokiem, w celu przeprowadzenia testu, jest wyznaczenie wartości statystyki testowej. Taką statystyką nazywamy statystykę pozwalającą orzekać w powyższym i jakimkolwiek innym problemie testowania hipotez, oraz oznaczamy:

$$T(\underline{X}) = T((x_1, x_2, \dots, x_n)).$$

W zależności od przeprowadzanego eksperymentu statystyka testowa ma ściśle określony wzór, w tym przypadku ma postać:

$$T(\underline{X}) = \frac{\hat{\mu} - \mu_{H_0}}{\sigma} \sqrt{n} = \frac{0.4375 - 0}{1} \cdot \sqrt{16} = 1.75.$$

Widzimy, że statystyka testowa zależy od hipotezy zerowej, a także od wyników próby, z czego wynika, że każda statystyka testowa jest zmienną losową, której rozkład określamy na podstawie założenia prawdziwości hipotezy zerowej. W tym przypadku statystyka T ma rozkład standardowy normalny.

Znając rozkład statystyki testowej możemy przejść do kolejnej czynności, czyli wyznaczenia obszaru krytycznego. Zbiór ten najczęściej oznaczamy symbolem C . Tworzą go wartości, które, przy założeniu prawdziwości hipotezy zerowej, są małoprawdopodobne do osiągnięcia przez statystykę testową. Wartości te jednocześnie zaprzeczają hipotezie zerowej. W przypadku zaobserwowania, że wartość naszej statystyki testowej należy do obszaru krytycznego, możemy odrzucić hipotezę zerową. Na położenie obszaru krytycznego ma wpływ hipoteza alternatywna. W zależności od jej formy wyróżniamy trzy rodzaje obszaru krytycznego:

- Dwustronny obszar krytyczny:

$$C = (-\infty, t_{\frac{\alpha}{2}}) \cup (t_{1-\frac{\alpha}{2}}, \infty),$$

gdzie $t_{\frac{\alpha}{2}}$, $t_{1-\frac{\alpha}{2}}$ są kwantylami rozkładu statystyki testowej, rzędów odpowiednio: $\frac{\alpha}{2}$, $1 - \frac{\alpha}{2}$. Obszar ten wyznaczamy na podstawie hipotezy alternatywnej niekierunkowej.

- Lewostronny obszar krytyczny:

$$C = (-\infty, t_{\alpha}),$$

gdzie t_{α} jest kwantylem rzędu α rozkładu statystyki T . Przedział ten wyznaczamy, kiedy alternatywa jest lewostronna.

- Prawostronny obszar krytyczny:

$$C = (t_{1-\alpha}, \infty),$$

gdzie $t_{1-\alpha}$ jest kwantylem rzędu $1 - \alpha$ rozkładu statystyki testowej. Powyższy obszar wyznacza hipoteza alternatywna prawostronna.

Wartości brzegowe obszaru krytycznego nazywamy wartościami krytycznymi. Dopełnienie obszaru krytycznego stanowią wartości statystyki testowej, które potwierdzają prawdziwość hipotezy zerowej. Zbiór ten nazywamy obszarem przyjęcia hipotezy zerowej.

Kolejnym sposobem, który pozwala nam zweryfikować, czy hipoteza zerowa jest prawdziwa, jest wyznaczenie ***p*-wartości**. Nazywamy w ten sposób przekształcenie statystyki testowej. Tak samo jak w przypadku obszaru krytycznego, wzór *p*-wartości zależy od kierunku, wyznaczanego przez hipotezę alternatywną. Jeśli alternatywa jest niekierunkowa, *p*-wartość definiujemy jako:

$$p := 2 \min \{ \mathbb{F}_{H_0}(T(\underline{X})), 1 - \mathbb{F}_{H_0}(T(\underline{X})) \},$$

gdzie \mathbb{F}_{H_0} jest dystrybucją rozkładu statystyki testowej, przy założeniu prawdziwości hipotezy zerowej. W sytuacji gdzie hipoteza alternatywna jest lewostronna, *p*-wartość ma wzór:

$$p := \mathbb{F}_{H_0}(T(\underline{X})),$$

a w sytuacji kiedy alternatywa jest prawostronna:

$$p := 1 - \mathbb{F}_{H_0}(T(\underline{X})).$$

Hipotezę zerową odrzucamy, gdy wartość *p* jest mniejsza od poziomu istotności testu.

W naszym modelu statystycznym (A) hipoteza alternatywna jest niekierunkowa, przez co wyznacza dwustronny obszar krytyczny. Wiedząc, że statystyka testowa pochodzi z rozkładu $N(0, 1)$, a poziom istotności równy jest 0.05, obliczmy wartości brzegowe:

$$t_{\frac{\alpha}{2}} = t_{0.025} = -1.96,$$

$$t_{1-\frac{\alpha}{2}} = t_{0.975} = 1.96.$$

Z powyższego wynika, że obszar krytyczny ma postać:

$$C_A = (-\infty, -1.96) \cup (1.96, \infty).$$

Obliczmy jeszcze *p*-wartość:

$$p_A = 2 \min \{ \Phi(1.75), 1 - \Phi(1.75) \} = 2 \cdot (1 - \Phi(1.75)) = 2 \cdot 0.04 = 0.08$$

Widzimy, że wartość statystyki testowej — $T(\underline{X}) = 1.75$, nie należy do obszaru krytycznego oraz p -wartość jest większa od poziomu istotności. Możemy przez to odrzucić prawdziwość hipotezy zerowej — $H_0 : \mu = 0$.

W kolejnym przypadku, bazując ciągle na danych z modelu (A) oraz zakładając ten sam poziom istotności, rozpatrzmy modele:

$$\begin{aligned} H_0 : \quad & \mu = 0, \\ H_1 : \quad & \mu < 0, \end{aligned} \tag{B}$$

oraz:

$$\begin{aligned} H_0 : \quad & \mu = 0, \\ H_1 : \quad & \mu > 0. \end{aligned} \tag{C}$$

Dostrzegamy, że powyższe modele różnią się od modelu (A) tylko formami hipotez alternatywnych, przez co statystyka testowa przyjmuje tą samą wartość dla wszystkich tych modeli. W przypadku modelu (B) wyznaczamy lewostronny obszar krytyczny, którego wartością brzegową jest kwantyl rzędu 0.05 rozkładu $N(0, 1)$.

$$t_\alpha = t_{0.05} = -1.64.$$

$$C_B = (-\infty, t_\alpha) = (-\infty, -1.64).$$

Wyznamy jeszcze wartość p :

$$p = \Phi(1.75) = 0.96.$$

Ponieważ $T(\underline{X}) \notin C_B$ i $p > 0.05$, uznajemy hipotezę zerową modelu (B) za prawdziwą. Dla modelu (C) wyznaczamy prawostronny obszar krytyczny, którego wartość krytyczna równa jest kwantylowi rzędu 0.95 standardowego rozkładu normalnego.

$$t_{1-\alpha} = t_{0.95} = 1.64.$$

$$C_C = (t_\alpha, \infty) = (1.64, \infty).$$

P -wartość w tym przypadku liczymy jako:

$$p = 1 - \Phi(1.75) = 0.04.$$

W tym modelu wartość statystyki testowej należy do obszaru krytycznego, a p -wartość jest mniejsza od poziomu istotności. Odrzucamy więc hipotezę zerową modelu (C) na rzecz jej alternatywy.

Widzimy, że test i -tego modelu, gdzie $i \in \{(A), (B), (C)\}$, możemy utożsamić z funkcją:

$$\varphi_i : \quad \underline{X} \rightarrow \{0, 1\},$$

$$\varphi_i(\underline{X}) = \mathbb{1}_{C_i}(T(\underline{X})).$$

Uzyskaną przez funkcję γ_i wartość „1” interpretujemy jako decyzję o odrzuceniu hipotezy zerowej na rzecz hipotezy alternatywnej, a wartość „0” jako decyzję o przyjęciu hipotezy zerowej. Moc testu i -tego modelu wyznaczamy funkcją:

$$\gamma_i : \quad \mathbb{R} \rightarrow [0, 1],$$

$$\gamma_i(\mu) = \mathbb{E}_\mu[\varphi_i(\underline{X})],$$

gdzie $\mu \in \Theta_{1i}$.

2.3 Lemat Neymana-Pearsona

Jeden ze skutecznych sposobów przeprowadzenia weryfikacji hipotez statystycznych podaje nam Lemat Neymana-Pearsona. Wykorzystanie tego twierdzenia pozwala nam skonstruować test jednostajnie najmocniejszy na poziomie istotności α . W przypadku, kiedy nasz model statystyczny opiera się na parze hipotez prostych, a rozmiar próby, w oparciu, o którą test ma zostać przeprowadzony, jest znany. Test ten opiera się na ilorazie wiarygodności.

Dla próby $\underline{X} = (x_1, \dots, x_n)$, z rozkładu o funkcji prawdopodobieństwa $f(x, \theta)$, gdzie parametr $\theta \in \Theta$, ilorazem wiarygodności nazywamy wyrażenie:

$$r(H_0, H_1) = \frac{L(\theta_1, \underline{X})}{L(\theta_0, \underline{X})},$$

gdzie $H_0 : \theta = \theta_0$, $H_1 : \theta = \theta_1$. Natomiast funkcję:

$$L : \Theta \rightarrow \mathbb{R},$$

$$L(\theta, \underline{X}) = \prod_{i=1}^n f(x_i),$$

nazywamy funkcją wiarygodności.

Lemat 2.1 (Neymana-Pearsona). *Niech $\underline{X} = (x_1, \dots, x_n)$, przy $n \geq 1$, będzie próbą niezależnych zmiennych z rozkładu prawdopodobieństwa F o funkcji gęstości lub funkcji masy prawdopodobieństwa $f(x, \theta)$. Model statystyczny niech będzie oparty na prostych hipotezach:*

$$H_0 : \theta = \theta_0,$$

$$H_1 : \theta = \theta_1.$$

Wtedy, test φ jednostajnie najmocniejszy na poziomie istotności α określa funkcja:

$$\varphi(\underline{X}) = \begin{cases} 1 & \text{gdy} \quad \prod_{i=1}^n f_0(x_i) < k \prod_{i=1}^n f_1(x_i) \\ \delta & \text{gdy} \quad \prod_{i=1}^n f_0(x_i) = k \prod_{i=1}^n f_1(x_i) \\ 0 & \text{gdy} \quad \prod_{i=1}^n f_0(x_i) > k \prod_{i=1}^n f_1(x_i) \end{cases},$$

gdzie $k \geq 0$ oraz $\delta \in [0, 1]$, są statymi wyznaczonymi z warunku:

$$\mathbb{E}_{H_0}[\varphi(\underline{X})] = \alpha.$$

Obszar krytyczny testu φ ma postać:

$$C = \left\{ (x_1, \dots, x_n) : \frac{L(\theta_1, \underline{X})}{L(\theta_0, \underline{X})} > k \right\}.$$

W celu lepszego zrozumienia Lematu Neymana-Pearsona rozpatrzmy poniższy przykład.

Przykład 2.1. Załóżmy, próbę niezależnych zmiennych losowych: $\underline{X} = \{-0.50, 0.13, -0.08, 0.89, 0.12, 0.32, -0.58, 0.71\}$, pochodzi z rozkładu normalnego, o nieznanym parametrze średniej μ oraz wariancji równej 1. Przyjmijmy poziom istotności $\alpha = 0.05$. Przetestujmy, przy pomocy lematu Neymana-Pearsona, poniższy model statystyczny:

$$H_0 : \mu = 0,$$

$$H_1 : \mu = 4.$$

Na początku wyznaczmy funkcje wiarygodności dla powyższych hipotez:

$$\begin{aligned} L_0(\theta_0, \underline{X}) &= \prod_{i=1}^8 \frac{1}{\sqrt{2\pi}} \exp\left(\frac{-x_i^2}{2}\right) \\ &= \left(\frac{1}{\sqrt{2\pi}}\right)^8 \exp\left(\frac{-\sum_{i=1}^8 x_i^2}{2}\right), \end{aligned}$$

$$\begin{aligned} L_1(\theta_1, \underline{X}) &= \prod_{i=1}^8 \frac{1}{\sqrt{2\pi}} \exp\left(\frac{-(x_i - 4)^2}{2}\right) \\ &= \left(\frac{1}{\sqrt{2\pi}}\right)^8 \exp\left(\frac{-\sum_{i=1}^8 (x_i - 4)^2}{2}\right). \end{aligned}$$

W kolejnym kroku wyznaczmy iloraz wiarygodności:

$$\begin{aligned} r(H_0, H_1) &= \frac{L_1(\theta_1, \underline{X})}{L_0(\theta_0, \underline{X})} = \exp\left(-\frac{1}{2}\left(\sum_{i=1}^8 (x_i - 4)^2 - \sum_{i=1}^8 x_i^2\right)\right) \\ &= \exp\left(-\frac{1}{2}\left(-8\sum_{i=1}^8 x_i + 128\right)\right) \\ &= \exp\left(4\sum_{i=1}^8 x_i - 64\right), \end{aligned}$$

gdzie:

$$\begin{aligned} \exp\left(4\sum_{i=1}^8 x_i - 64\right) &> k_1 \quad / \ln \\ 4\sum_{i=1}^8 x_i - 64 &> \ln k_1 \\ \sum_{i=1}^8 x_i &> \frac{\ln k_1 - 64}{4} = k. \end{aligned}$$

Zatem obszar krytyczny jednostajnie najmocniejszego testu φ ma postać:

$$C = \left\{ (x_1, \dots, x_8) : \exp\left(4\sum_{i=1}^8 x_i - 64\right) > k_1 \right\} = \left\{ (x_1, \dots, x_8) : \sum_{i=1}^8 x_i > k \right\}$$

Obliczmy wartość k z warunku:

$$P(\underline{X} \in C \mid H_0) = \alpha.$$

Wiemy, że dla niezależnych zmiennych losowych y_1, \dots, y_n ze standardowego rozkładu nor-

malnego, $\sum_{i=1}^n y_i \sim N(0, n)$. Z tego wynika:

$$\begin{aligned} P(\underline{X} \in C \mid H_0) &= P\left(\sum_{i=1}^8 x_i > k \mid H_0\right) \\ &= P\left(\frac{1}{\sqrt{8}} \sum_{i=1}^8 x_i > \frac{k}{\sqrt{8}} \mid H_0\right) \\ &= 1 - \Phi\left(\frac{k}{\sqrt{8}}\right). \end{aligned}$$

Wobec tego wartość k wynosi:

$$\begin{aligned} \frac{k}{\sqrt{8}} &= \Phi^{-1}(1 - 0.05) \\ k &= \sqrt{8}\Phi^{-1}(0.95) \approx 4.6523. \end{aligned}$$

Widzimy, że:

$$\sum_{i=1}^8 x_i = 1.01 < k.$$

W skutek tego możemy przyjąć hipotezę zerową.

3 Rozkład p -wartości przy założeniu hipotezy zerowej

Z poprzedniego rozdziału wiemy, że p -wartość jest przekształceniem statystyki testowej. Inaczej, p -wartość określa prawdopodobieństwo, przy założeniu prawdziwości hipotezy zerowej, przyjęcia przez statystykę testową wartości „ekstremalnych”, sygnalizujących, aby w wyniku testu odrzucić hipotezę zerową. Z tego wynika, że czynnikiem, który wpływa na rozkład p -wartości, jest rozkład statystyki testowej. Spróbujmy określić rozkład p -wartości w dwóch przypadkach: kiedy statystyka testowa przyjmuje rozkład ciągły oraz, kiedy statystyka testowa przyjmuje rozkład dyskretny.

3.1 Rozkład p -wartości opartej na ciągłej statystyce testowej

Rozpatrzmy sytuację, kiedy statystyka testowa pochodzi z ciągłego rozkładu prawdopodobieństwa. Rozwiązanie problemu określenia rozkładu p -wartości w tej sytuacji podaje nam poniższe twierdzenie.

Twierdzenie 3.1. *Kiedy statystyka testowa pochodzi z rozkładu bezwzględnie ciągłego, wyznaczona na jej podstawie p -wartość przyjmuje rozkład jednostajny na przedziale $[0, 1]$.*

Dowód. Załóżmy, że w wyniku lewostronnego testu prostej hipotezy zerowej, opartego na próbie niezależnych zmiennych losowych: $\underline{X} = (x_1, \dots, x_n)$, gdzie $n \geq 1$, pochodzącej z ciągłego rozkładu prawdopodobieństwa \mathcal{G} , otrzymaliśmy statystykę testową: $T(\underline{X})$. Niech przy założeniu hipotezy zerowej:

$$T(\underline{X}) \sim \mathcal{F}.$$

Jako że próba \underline{X} pochodzi z rozkładu ciągłego, bez straty ogólności rozważań, załóżmy, że \mathcal{F} jest ciągłym rozkładem prawdopodobieństwa, którego dystrybuanta \mathbb{F} jest ściśle rosnąca dla nośnika rozkładu \mathcal{F} . Następstwem tego jest fakt, że istnieje funkcja kwantylowa \mathbb{F}^{-1} na zbiorze wartości, dla których gęstość rozkładu \mathcal{F} przyjmuje wartości niezerowe, czyli nośnika rozkładu \mathcal{F} . Przyjmijmy zmienną losową $z \sim \mathcal{F}$. Wyznaczmy p -wartość:

$$p = P(z < T(\underline{X}) \mid H_0) = \mathbb{F}(T(\underline{X})).$$

Z tego wynika, że $p \in (0, 1)$. Przyjmijmy zmienną $x \in (0, 1)$, oraz obliczmy prawdopodobieństwo, przy założeniu prostej hipotezy zerowej, tego że $p \leq x$:

$$\begin{aligned} P(p \leq x \mid H_0) &= P(\mathbb{F}(T(\underline{X})) \leq x \mid H_0) \\ &= P(T(\underline{X}) \leq \mathbb{F}^{-1}(x) \mid H_0) \\ &= \mathbb{F}(\mathbb{F}^{-1}(x)) \\ &= x. \end{aligned}$$

Wiedząc, że dla każdej zmiennej losowej $y \sim \mathcal{U}[0, 1]$, spełniony jest warunek:

$$P(y \leq x) = x,$$

równość:

$$P(p \leq x \mid H_0) = x$$

dowodzi, że p -wartość pochodzi z rozkładu $\mathcal{U}[0, 1]$. ■

Wykorzystanie testu prawostronnego lub testu niekierunkowego nie zmienia faktu, że wartość p , wyznaczona na podstawie statystyki o rozkładzie ciągłym, ma rozkład $\mathcal{U}[0, 1]$. Dowód oparty na teście lewostronnym jest najmniej skomplikowany i najłatwiejszy do przeprowadzenia.

3.2 Rozkład p -wartości opartej na dyskretnej statystyce testowej

Rozpatrzmy przypadek, kiedy statystyka testowa pochodzi z rozkładu dyskretnego. W tej sytuacji, rozkładu p -wartości nie można określić jednoznacznie, jednak w niektórych sytuacjach, rozważane przekształcenie statystyki testowej przyjmuje rozkład asymptotycznie jednostajny na przedziale $[0, 1]$. Oznacza to, że kiedy rozmiar próby, na podstawie której określona jest wartość statystyki testowej, dąży do nieskończoności, to rozkład p -wartości dąży do rozkładu $\mathcal{U}[0, 1]$. Poniższa uwaga wynika z Prawa Wielkich Liczb oraz Twierdzenia Gliwienki-Cantellego.

Uwaga 3.1. *Kiedy statystyka testowa pochodzi z rozkładu dyskretnego, p -wartość wyznaczona na jej podstawie może przyjąć rozkład asymptotycznie jednostajny na przedziale $[0, 1]$.*

Aby sprawdzić prawdziwość powyższej uwagi, rozważmy poniższy przykład.

Przykład 3.1. Niech $\underline{X} = (x_1, \dots, x_n)$, gdzie $n \geq 1$, będzie próbą niezależnych zmiennych losowych z rozkładu Poissona o parametrze λ . Rozważmy, na poziomie istotności $\alpha \in [0, 1]$, prawostronny test hipotez:

$$\begin{aligned} H_0 : \quad & \lambda = \lambda_0, \\ H_1 : \quad & \lambda > \lambda_0, \end{aligned}$$

gdzie $\lambda_0 \in (0, \infty)$. Wyznaczmy statystykę testową:

$$T(\underline{X}) = \sum_{i=1}^n x_i.$$

Widzimy, że przy założeniu hipotezy zerowej, statystyka $T(\underline{X})$ przyjmuje rozkład Poissona o parametrze $n\lambda_0$. Załóżmy, że zmienna losowa z pochodzi z rozkładu Poissona o parametrze $n\lambda_0$. Wyznaczmy p -wartość:

$$p = P(z > T(\underline{X}) \mid H_0) = 1 - \mathbb{F}(T(\underline{X})),$$

gdzie $\mathbb{F}(\bullet)$ jest dystrybuantą rozkładu Poissona o parametrze $n\lambda_0$. Niech $x \in [0, 1]$. Obliczmy prawdopodobieństwo, przy założeniu prawdziwości hipotezy zerowej, że $p \leq x$:

$$\begin{aligned} P(p \leq x \mid H_0) &= P(1 - \mathbb{F}(T(\underline{X})) \leq x \mid H_0) \\ &= P(\mathbb{F}(T(\underline{X})) \geq 1 - x \mid H_0). \end{aligned} \quad (1)$$

Niech

$$k = \min \left\{ m : \sum_{i=0}^m \frac{(n\lambda_0)^i e^{-n\lambda_0}}{i!} \geq 1 - x \right\}.$$

Z tego wynika, że:

$$\begin{aligned} P(p \leq x \mid H_0) &= P(T(\underline{X}) \geq k \mid H_0) \\ &= 1 - \mathbb{F}(k), \end{aligned} \quad (2)$$

gdzie:

$$1 - x \leq \mathbb{F}(k) \leq 1 - x + \frac{(n\lambda_0)^k e^{-n\lambda_0}}{k!}. \quad (3)$$

Oznaczmy $L = n\lambda_0$, $a_k = \frac{L^k e^{-L}}{k!}$. Wykażmy, że przy $n \rightarrow \infty$, $\mathbb{F}(k) \rightarrow 1 - x$. Ponieważ $L > 0$ oraz $k \geq 0$, to $\forall k : a_k > 0$. Sprawdźmy monotoniczność ciągu a_k . Ciąg jest rosnący, kiedy:

$$a_k < a_{k+1} \iff \frac{a_{k+1}}{a_k} > 1.$$

Wyznaczmy:

$$\frac{a_{k+1}}{a_k} = \frac{\frac{L^{k+1} e^{-L}}{(k+1)!}}{\frac{L^k e^{-L}}{k!}} = \frac{L}{k+1}.$$

Z tego wynika, że:

$$\frac{a_{k+1}}{a_k} > 1 \iff k < \lfloor L \rfloor,$$

gdzie $\lfloor r \rfloor$ oznacza część całkowitą liczby rzeczywistej r . Natomiast ciąg a_k jest malejący, kiedy:

$$a_k > a_{k+1} \iff \frac{a_{k+1}}{a_k} < 1 \iff k > \lfloor L \rfloor.$$

Z tego wynika, że:

$$\sup_{k \geq 0} \left(\frac{L^k e^{-L}}{k!} \right) = \frac{L^{\lfloor L \rfloor} e^{-L}}{\lfloor L \rfloor!}.$$

Wobec tego:

$$\begin{aligned} \lim_{n \rightarrow \infty} \left(\sup_{k \geq 0} \left(\frac{(n\lambda_0)^k e^{-n\lambda_0}}{k!} \right) \right) &= \lim_{L \rightarrow \infty} \left(\sup_{k \geq 0} \left(\frac{L^k e^{-L}}{k!} \right) \right) \\ &= \lim_{L \rightarrow \infty} \left(\frac{L^{\lfloor L \rfloor} e^{-L}}{\lfloor L \rfloor!} \right). \end{aligned}$$

Oszacujmy:

$$0 \leq \frac{L^{\lfloor L \rfloor} e^{-L}}{\lfloor L \rfloor!} \leq \frac{(\lfloor L \rfloor + 1)^{\lfloor L \rfloor} e^{-\lfloor L \rfloor}}{\lfloor L \rfloor!}. \quad (4)$$

Korzystając ze wzoru Stirlinga, dla $i \in \mathbb{N}_+$:

$$i! \approx \left(\frac{i}{e} \right)^i \sqrt{2\pi i}, \quad \text{gdzie } \lim_{i \rightarrow \infty} \frac{i!}{\left(\frac{i}{e} \right)^i \sqrt{2\pi i}} = 1, \quad (5)$$

obliczmy:

$$\begin{aligned}
\lim_{L \rightarrow \infty} \left(\frac{([\![L]\!] + 1)^{[\![L]\!]} e^{-[\![L]\!]}}{[\![L]\!]!} \right) &= \lim_{L \rightarrow \infty} \left(\frac{([\![L]\!] + 1)^{[\![L]\!]} e^{-[\![L]\!]}}{[\![L]\!]!} \cdot \frac{[\![L]\!] + 1}{[\![L]\!] + 1} \right) \\
&= \lim_{L \rightarrow \infty} \left(\frac{([\![L]\!] + 1)^{[\![L]\!] + 1} e}{([\![L]\!] + 1)! \cdot e^{[\![L]\!] + 1}} \right) \\
&\stackrel{(5)}{=} \lim_{L \rightarrow \infty} \left(\frac{([\![L]\!] + 1)^{[\![L]\!] + 1} e}{\left(\frac{[\![L]\!] + 1}{e} \right)^{[\![L]\!] + 1} \sqrt{2\pi([\![L]\!] + 1)} \cdot e^{[\![L]\!] + 1}} \right) \\
&= \lim_{L \rightarrow \infty} \left(\frac{e}{\sqrt{2\pi([\![L]\!] + 1)}} \right) \\
&= 0.
\end{aligned} \tag{6}$$

Z Twierdzenia o trzech ciągach, (4) oraz (6) wynika, że:

$$\lim_{n \rightarrow \infty} \left(\sup_{k \geq 0} \left(\frac{(n\lambda_0)^k e^{-n\lambda_0}}{k!} \right) \right) = 0. \tag{7}$$

Poprzez (3), (2) pokazaliśmy, że:

$$\mathbb{F}(k) \xrightarrow{n \rightarrow \infty} 1 - x. \tag{8}$$

Z (1), (2), (8) wynika, że:

$$P(p \leq x \mid H_0) \xrightarrow{n \rightarrow \infty} 1 - (1 - x) = x.$$

W ten sposób zareprezentowaliśmy, że rozkład p -wartości, dla tego przykładu, jest asymptotycznie jednostajny na przedziale $[0, 1]$.

4 Globalna hipoteza zerowa oraz sposób jej weryfikacji

W celu znalezienia w zbiorze danych próby, która w stosunku do innych prób, różni się rozkładem prawdopodobieństwa, możemy posłużyć się wielokrotnym testowaniem hipotez. Sposób ten polega na przeprowadzeniu jednakowego testu hipotez każdej próby należącej do badanego zbioru danych. W sytuacji, kiedy liczba prób jest duża, sposób ten jest nieoptymalny. Wynika to z tego, że przy pomocy wielokrotnej weryfikacji hipotez możemy znaleźć nie tylko jedną próbę odstającą, ale również czas przeprowadzenia testu, może być zbyt długi. Bardziej optymalnym rozwiązaniem powyższego problemu jest weryfikacja modelu typu „igła w stogu siana”, który jest oparty o globalną hipotezę zerową.

4.1 Globalna hipoteza zerowa

Globalna hipoteza zerowa jest pojęciem związanym z wielokrotnym testowaniem hipotez statystycznych. Załóżmy, że dysponujemy zbiorem danych złożonym z prób: X_1, \dots, X_n , gdzie $n > 1$. Niech każda próba X_i , gdzie $i \in \{1, \dots, n\}$, pochodzi z rozkładu \mathcal{F} o parametrze $\theta \in \mathbb{R}$. Dla każdej próby X_i chcemy przeprowadzić test oparty na modelu:

$$\begin{aligned}
H_{0_i} &: \theta_i = 0, \\
H_{1_i} &: \theta_i \neq 0.
\end{aligned}$$

Globalną hipotezą zerową nazywamy wyrażenie:

$$H_0 = \bigcap_{i=1}^n H_{0_i},$$

z tego wynika, że globalna hipoteza zerowa, w tym przypadku, ma postać:

$$H_0 : \theta_1 = \dots = \theta_n = 0.$$

Globalną hipotezę zerową przyjmujemy za prawdziwą wyłącznie wtedy, gdy każda z hipotez: H_{0_1}, \dots, H_{0_n} jest prawdziwa. Alternatywa powyższej, globalnej hipotezy zerowej ma postać:

$$H_1 : \exists i \in \{1, \dots, n\} : \theta_i \neq 0.$$

Zakłada ona, że przynajmniej jedna próba pochodzi z rozkładu \mathcal{F} o parametrze $\theta \neq 0$, ale nie określa ile i które to są konkretnie próby.

4.2 Sposób testowania globalnej hipotezy zerowej

Jeden ze sposobów przeprowadzenia wielokrotnego testowania hipotez statystycznych i weryfikacji prawdziwości globalnej hipotezy zerowej, podaje nam **korekta Bonferroniego**. Procedura ta odrzuca globalną hipotezę zerową na poziomie istotności α , gdy:

$$\min_{i \in \{1, \dots, n\}} p_i \leq \frac{\alpha}{n},$$

gdzie p_i jest p -wartością wyznaczoną na podstawie statystyki testowej: $T(\underline{X}_i)$, uzyskanej w wyniku testu i -tego modelu statystycznego: (H_{0_i}, H_{1_i}) . Uzyskanie prawidłowych wyników przez przeprowadzenie przez nas testu opartego na powyższej procedurze, potwierdza następujące twierdzenie.

Twierdzenie 4.1. *Korekta Bonferroniego kontroluje na poziomie α prawdopodobieństwo odrzucenia globalnej hipotezy zerowej, kiedy jest ona prawdziwa.*

Dowód. Niech $p_1, \dots, p_n \sim \mathcal{U}[0, 1]$, gdzie $n \geq 2$, będą p -wartościami uzyskanymi odpowiednio w wyniku testów par hipotez: $(H_{0_1}, H_{1_1}), \dots, (H_{0_n}, H_{1_n})$. Globalna hipoteza zerowa ma postać:

$$H_0 = \bigcap_{i=1}^n H_{0_i}.$$

Niech poziom istotności testu globalnej hipotezy zerowej wynosi α . Obliczmy prawdopodobieństwo popełnienia błędu I rodzaju, w wyniku testu globalnej hipotezy zerowej, opartego na korekcie Bonferroniego:

$$\begin{aligned} P\left(\min_{i \in \{1, \dots, n\}} p_i \leq \frac{\alpha}{n} \mid H_0\right) &= P\left(\bigcup_{i=1}^n \left\{p_i \leq \frac{\alpha}{n}\right\} \mid H_0\right) \\ &\leq \sum_{i=1}^n P\left(p_i \leq \frac{\alpha}{n} \mid H_0\right) \\ &\leq \sum_{i=1}^n \frac{\alpha}{n} \\ &= \alpha. \end{aligned}$$

W sytuacji, kiedy wartości: p_1, \dots, p_n są niezależnymi zmiennymi losowymi, pochodzącymi z rozkładu $\mathcal{U}[0, 1]$, powyższe prawdopodobieństwo obliczamy jako:

$$\begin{aligned} P\left(\min_{i \in \{1, \dots, n\}} p_i \leq \frac{\alpha}{n} \mid H_0\right) &= 1 - P\left(\bigcap_{i=1}^n \left\{p_i \geq \frac{\alpha}{n}\right\} \mid H_0\right) \\ &= 1 - \prod_{i=1}^n P\left(p_i \geq \frac{\alpha}{n} \mid H_0\right) \\ &= 1 - \left(1 - \frac{\alpha}{n}\right)^n \\ &\approx 1 - e^{-\alpha} \\ &\leq \alpha. \end{aligned}$$

Z powyższego wynika, że korekta Bonferroniego kontroluje poziom istotności testu globalnej hipotezy zerowej. ■

Korekta Bonferroniego pozwala prawidłowo zweryfikować globalną hipotezę zerową. Oznacza to również, że test oparty na tej procedurze jest jednym z bardziej optymalnych sposobów testowania modeli typu „igła w stogu siana”. W celu uzyskania większej wiedzy na temat sposobów weryfikacji globalnej hipotezy zerowej, możemy sięgnąć po artykuł pt. „An omnibus test for the global null hypothesis”^[3].

5 Modele typu „igła w stogu siana” oraz optymalne progi ich wykrycia

Definiując globalną hipotezę zerową oraz objaśniając sposób jej weryfikacji, możemy przejść do rozważań na temat modeli typu „igła w stogu siana”. Załóżmy, że próby: $\underline{X}_1, \dots, \underline{X}_n$, gdzie $n \geq 2$, pochodzą z rozkładu \mathcal{F} o parametrze $\theta \in \mathbb{R}$. Niech hipoteza zerowa: $H_{0_i} : \theta_i = 0$, gdzie $i \in \{1, \dots, n\}$, odpowiada próbie \underline{X}_i . Model statystyczny:

$$\begin{aligned} H_0 : \quad &\forall i \in \{1, \dots, n\} : \theta_i = k, \\ H_1 : \quad &\exists! i \sim \mathcal{U}\{1, \dots, n\} : \theta_i \neq k, \end{aligned}$$

gdzie $k \in \mathbb{R}$, nazywany **modelem typu „igła w stogu siana”**. Aby skonstruować model tego typu, musimy sformułować hipotezę alternatywną do globalnej hipotezy zerowej, tak aby zakładała, że tylko jedna spośród badanych prób różni się rozkładem w stosunku do reszty. Powodem takiej konstrukcji jest fakt, że hipoteza alternatywna:

$$H'_1 : \quad \exists i \in \{1, \dots, n\} : \theta_i \neq k,$$

jest zbyt złożona, aby korekta Bonferroniego mogła odrzucić hipotezę H_0 na rzecz hipotezy H'_1 . Procedura ta odrzuca hipotezę H_0 w oparciu o najmniejszą wartość p_i , więc nie możemy stwierdzić przy jej pomocy, czy pośród innych prób znajdują się kolejne, których parametr rozkładu jest różny od 0. Przykładem użycia modelu typu „igła w stogu siana” jest próba wykrycia, czy w zbiorze operacji finansowych znajduje się transakcja, która jest nieuczciwa. Jednym z rozwiązań tego problemu jest weryfikacja odpowiednio skonstruowanego modelu typu „igła w stogu siana”. Aby znaleźć odpowiedni sposób weryfikacji rozważanych modeli, spróbujmy w następnych podrozdziałach przeanalizować modele oparte na danych pochodzących z rozkładów normalnego oraz Poissona.

5.1 Model oparty o dane pochodzące z rozkładu normalnego

Założmy, że zbiór danych, który chcemy poddać analizie, składa się z jednoelementowych prób: $\underline{X}_1, \dots, \underline{X}_n$, gdzie $n \geq 2$. Niech, dla $i \in \{1, \dots, n\}$, próba \underline{X}_i pochodzi z rozkładu $\mathcal{N}(\mu_i, 1)$. Niech każdej próbie \underline{X}_i odpowiada model statystyczny:

$$\begin{aligned} H_{0_i} &: \mu_i = 0 \\ H_{1_i} &: \mu_i > 0. \end{aligned}$$

Z tego wynika, że statystykę testową i -tej hipotezy zerowej obliczamy w następujący sposób:

$$T(\underline{X}_i) = \frac{\underline{X}_i - 0}{1} = \underline{X}_i.$$

Widzimy, że przy założeniu prawdziwości i -tej hipotezy zerowej, statystyka testowa $T(\underline{X}_i)$ pochodzi ze standardowego rozkładu normalnego. Wyznamy p -wartość opartą na i -tej statystyce testowej:

$$p_i = 1 - \Phi(T(\underline{X}_i)) = |\Phi(T(\underline{X}_i))|$$

Głównym problemem, o który oprzemy dalsze rozważania, jest sprawdzenie, czy spośród prób: $\underline{X}_1, \dots, \underline{X}_n$, znajduje się wyłącznie jedna, której parametr rozkładu μ jest większy od 0. Opierając się na tym problemie, skonstruujemy model typu „igła w stogu siana”:

$$\begin{aligned} H_0 &: \forall i \in \{1, \dots, n\} : \mu_i = 0, \\ H_1 &: \exists! i \sim \mathcal{U}\{1, \dots, n\} : \mu_i > 0. \end{aligned} \tag{D1}$$

Wiemy, że Korekta Bonferroniego odrzuca globalną hipotezę zerową H_0 na poziomie istotności α , kiedy:

$$\min_{i \in \{1, \dots, n\}} p_i \leq \frac{\alpha}{n},$$

co jest równoznaczne, że H_0 zostanie odrzucona, gdy:

$$\max_{i \in \{1, \dots, n\}} T(\underline{X}_i) \geq \left| \Phi^{-1}\left(\frac{\alpha}{n}\right) \right|,$$

gdzie $\Phi^{-1}(\cdot)$ jest funkcją kwantylową standardowego rozkładu normalnego.

5.1.1 Optymalny próg wykrycia

Znalezienie dokładnej wartości kwantyla: $|\Phi^{-1}(\frac{\alpha}{n})|$, stanowi utrudnienie, w przypadku kiedy liczba prób n jest zbyt duża, a wartość poziomu istotności α jest bliska zeru. W rozwiązaniu tego problemu pomaga nam przybliżenie wartości powyższego kwantyla. Jednak w sytuacji, kiedy różnica między rzeczywistą wartością, a przybliżeniem liczby $|\Phi^{-1}(\frac{\alpha}{n})|$ jest dostatecznie duża, skonstruowany przez nas test statystyczny, oparty na przybliżeniu, może doprowadzić do używania błędnych rezultatów. Jedną ze skutecznych metod uzyskania odpowiedniej aproksymacji wyrażenia: $|\Phi^{-1}(\frac{\alpha}{n})|$ podaje nam Nierówność Markowa.

Twierdzenie 5.1. *Niech, w przypadku testu jednostronnego, $t = |\Phi^{-1}(\frac{\alpha}{n})|$, lub $t = |\Phi^{-1}(\frac{\alpha}{2n})|$, w przypadku testu niekierunkowego. Dla zmiennej losowej Z , pochodzącej z rozkładu $\mathcal{N}(0, 1)$, prawdziwa jest nierówność:*

$$\frac{\phi(t)}{t} \left(1 - \frac{1}{t^2}\right) \leq P(Z > t) \leq \frac{\phi(t)}{t},$$

gdzie $\phi(\cdot)$ jest funkcją gęstości standardowego rozkładu normalnego.

Dowód. Niech Z będzie zmienną losową, pochodzącą z rozkładu $\mathcal{N}(0, 1)$, o funkcji gęstości: $\phi(\cdot)$, oraz dystrybuancie: $\Phi(\cdot)$. Bez straty ogólności założmy, że $t = |\Phi^{-1}(\frac{\alpha}{n})|$, gdzie $\alpha \in (0, 1)$ oraz $n \in \mathbb{N}_+$.

- (i) Udowodnijmy nierówność: $P(Z > t) \geq \frac{\phi(t)}{t} \left(1 - \frac{1}{t^2}\right)$, która jest jednoznaczna z nierównością: $P(Z > t) - \frac{\phi(t)}{t} \left(1 - \frac{1}{t^2}\right) \geq 0$. Wiemy, że $t \in (0, \infty)$. Przeprowadźmy proste badanie przebiegu zmienności funkcji:

$$f : (0, \infty) \rightarrow \mathbb{R},$$

$$f(t) = P(Z > t) - \frac{\phi(t)}{t} \left(1 - \frac{1}{t^2}\right) = 1 - \Phi(t) - \frac{\phi(t)}{t} + \frac{\phi(t)}{t^3}.$$

Wyznaczmy granice powyższej funkcji:

$$\begin{aligned} \lim_{t \rightarrow 0^+} f(t) &= \lim_{t \rightarrow 0^+} (1 - \Phi(t)) - \lim_{t \rightarrow 0^+} (\phi(t)) \cdot \lim_{t \rightarrow 0^+} \left(\frac{t^2 - 1}{t^3}\right) \\ &= \left[\frac{1}{2} - \phi(0) \cdot (-\infty)\right] = \infty, \end{aligned} \tag{1}$$

$$\begin{aligned} \lim_{t \rightarrow \infty} f(t) &= \lim_{t \rightarrow \infty} (1 - \Phi(t)) - \lim_{t \rightarrow \infty} (\phi(t)) \cdot \lim_{t \rightarrow \infty} \left(\frac{t^2 - 1}{t^3}\right) \\ &= [0 - 0 \cdot 0] = 0. \end{aligned} \tag{2}$$

Obliczmy pochodne pierwszego rzędu funkcji $\phi(t)$ oraz $f(t)$:

$$\begin{aligned} \frac{d}{dt} \phi(t) &= \frac{d}{dt} \left(\frac{1}{\sqrt{2\pi}} e^{-\frac{t^2}{2}}\right) \\ &= \frac{1}{\sqrt{2\pi}} e^{-\frac{t^2}{2}} \cdot \left(-\frac{1}{2} \cdot 2t\right) \\ &= -t \cdot \phi(t), \end{aligned} \tag{3}$$

$$\begin{aligned} \frac{d}{dt} f(t) &= 0 - \frac{d}{dt} \Phi(t) - \frac{d}{dt} \left(\phi(t) \left(\frac{t^2 - 1}{t^3}\right)\right) \\ &= -\phi(t) - \left(\left(\frac{t^2 - 1}{t^3}\right) \cdot \frac{d}{dt} \phi(t) + \phi(t) \cdot \frac{d}{dt} \left(\frac{t^2 - 1}{t^3}\right)\right) \\ &\stackrel{(3)}{=} -\phi(t) - \left(\frac{t^2 - 1}{t^3}\right) \cdot (-t \cdot \phi(t)) - \phi(t) \left(\frac{-t^2 + 3}{t^4}\right) \\ &= -\phi(t) \left(1 + \frac{t^2 - t^4}{t^4} + \frac{3 - t^2}{t^4}\right) \\ &= -\frac{3}{t^4} \phi(t). \end{aligned}$$

Widzimy, że:

$$\forall t \in (0, \infty) : \frac{d}{dt} f(t) < 0, \tag{4}$$

więc $f(t)$ jest funkcją malejącą, dla $t \in (0, \infty)$. Z (1), (2), (4) wynika, że prawdziwa jest nierówność:

$$P(Z > t) \geq \frac{\phi(t)}{t} \left(1 - \frac{1}{t^2}\right).$$

(ii) Wykażmy, że $P(Z > t) \leq \frac{\phi(t)}{t}$. Niech $k > t$, wtedy:

$$\begin{aligned} P(Z > t) &= \int_t^\infty \phi(u) du \stackrel{(u \geq t)}{\leq} \int_t^\infty \frac{u}{t} \phi(u) du = \frac{1}{t} \int_t^\infty u \cdot \phi(u) du \\ &\stackrel{(3)}{=} \frac{1}{t} \int_t^\infty \frac{d}{du} (-\phi(u)) du = -\frac{1}{t} \phi(u) \Big|_t^\infty \\ &= \frac{\phi(t)}{t}. \end{aligned}$$

Z rozważań zawartych w podpunktach (i) oraz (ii) wynika, że prawdziwa jest nierówność:

$$\frac{\phi(t)}{t} \left(1 - \frac{1}{t^2}\right) \leq P(Z > t) \leq \frac{\phi(t)}{t}.$$

■

Za pomocą powyższego twierdzenia, możemy wyaprosymować wartość: $t = \left| \Phi^{-1} \left(\frac{\alpha}{n} \right) \right|$. Wykazaliśmy w ten sposób, że dla $n \in \mathbb{N}_+$ oraz ustalonej wartości $\alpha \in (0, 1)$:

$$1 - \Phi(t) = \frac{\alpha}{n} \iff \frac{\phi(t)}{t} \approx \frac{\alpha}{n},$$

kiedy $t \xrightarrow{n \rightarrow \infty} \infty$. Rozpatrzmy następstwo tego faktu:

$$\begin{aligned} \frac{\phi(t)}{t} &\approx \frac{\alpha}{n} \\ \frac{1}{t\sqrt{2\pi}} e^{-\frac{t^2}{2}} &\approx \frac{\alpha}{n} \\ \frac{1}{\sqrt{2\pi}} e^{-\frac{t^2}{2} - \ln(t)} &\approx \frac{\alpha}{n} && \backslash \ln(\cdot) \\ -\ln(\sqrt{2\pi}) - \frac{t^2}{2} - \ln(t) &\approx \ln(\alpha) - \ln(n) \\ -\frac{t^2}{2} - \ln(t) &\approx \ln(\alpha) - \ln(n) + \ln(\sqrt{2\pi}) \end{aligned}$$

co, dla $n \rightarrow \infty$, jest równoważne z:

$$\begin{aligned} \frac{t^2}{2} &\approx \ln(n) \\ t &\approx \sqrt{2 \ln(n)}. \end{aligned}$$

Więc:

$$\frac{\phi(t)}{t} \approx \frac{\alpha}{n} \iff t = \left| \Phi^{-1} \left(\frac{\alpha}{n} \right) \right| \approx \sqrt{2 \ln(n)}.$$

Konsekwentnie, w przypadku modelu statystycznego (D1), stosując korektę Bonferronio, odrzucamy globalną hipotezę zerową H_0 , dla $n \rightarrow \infty$, gdy:

$$\max_{i \in \{1, \dots, n\}} |X_i| > \sqrt{2 \ln(n)}.$$

Zauważmy, że wynik tak skonstruowanego testu nie zależy od poziomu istotności α dla dużych n . Niezależnie od ustalonego poziomu istotności α testu modelu (D1), próg odrzucenia

dla wartości: $\max_{i \in \{1, \dots, n\}} |X_i|$, będzie asymptotycznie równy $\sqrt{2 \ln(n)}$. Potwierdza to fakt, że przy założeniu prawdziwości globalnej hipotezy zerowej:

$$\frac{\max_{i \in \{1, \dots, n\}} |X_i|}{\sqrt{2 \ln(n)}} \xrightarrow{D} 1.$$

5.1.2 Asymptotyczna moc korekty Bonferroniego

Sprawdźmy, jaką moc osiąga korekta Bonferroniego, w zależności od wielkości parametru średniej rozkładu normalnego, z którego pochodzą próby. Kolejne rozważania oprzyjmy o test modelu (D1) oraz założenia związane z tym modelem.

W celu wyznaczenia asymptotycznej mocy korekty Bonferroniego, bez straty ogólności, załóżmy, że $\mu_1 > 0$. Rozpatrzmy dwa przypadki:

1. Niech $\mu_1 = (1 - \varepsilon)\sqrt{2 \ln(n)}$, gdzie $\varepsilon \in (0, 1)$. Wiedząc, że $X_1 = z + \mu_1$, gdzie $z \sim \mathcal{N}(0, 1)$, obliczmy asymptotyczną moc korekty Bonferroniego:

$$\begin{aligned} P\left(\min_{i \in \{1, \dots, n\}} p_i \leq \frac{\alpha}{n} \mid H_1\right) &\leq P\left(p_1 \leq \frac{\alpha}{n} \mid H_1\right) + P\left(\min_{i \in \{2, \dots, n\}} p_i \leq \frac{\alpha}{n} \mid H_1\right) \\ &\leq P\left(X_1 > |\Phi^{-1}(\alpha/n)| \mid H_1\right) + 1 - \left(1 - \frac{\alpha}{n}\right)^{n-1} \\ &\leq P\left(z > |\Phi^{-1}(\alpha/n)| - \mu_1 \mid H_1\right) + 1 - \left(1 - \frac{\alpha}{n}\right)^{n-1} \\ &= P\left(z > \varepsilon\sqrt{2 \ln(n)} \mid H_1\right) + 1 - \left(1 - \frac{\alpha}{n}\right)^{n-1} \\ &\xrightarrow{n \rightarrow \infty} 0 + 1 - e^{-\alpha} \approx \alpha. \end{aligned}$$

2. Niech $\mu_1 = (1 + \varepsilon)\sqrt{2 \ln(n)}$. Moc korekty Bonferroniego w tym przypadku wynosi:

$$\begin{aligned} P\left(\min_{i \in \{1, \dots, n\}} p_i \leq \frac{\alpha}{n} \mid H_1\right) &\geq P\left(p_1 \leq \frac{\alpha}{n} \mid H_1\right) \\ &= P\left(X_1 > |\Phi^{-1}(\alpha/n)| \mid H_1\right) \\ &= P\left(z > |\Phi^{-1}(\alpha/n)| - \mu_1 \mid H_1\right) \\ &= P\left(z > -\varepsilon\sqrt{2 \ln(n)} \mid H_1\right) \\ &\xrightarrow{n \rightarrow \infty} 1 \end{aligned}$$

Z powyższych rozważań wynika, że korekta Bonferroniego odrzuci globalną hipotezę zerową H_0 , gdy parametr μ odstającej obserwacji będzie równy $(1 + \varepsilon)\sqrt{2 \ln(n)}$. W przypadku, kiedy $\mu = (1 - \varepsilon)\sqrt{2 \ln(n)}$, korekta Bonferroniego nie wykryje obserwacji. Wtedy test oparty na tej procedurze osiąga asymptotycznie moc równą α . Oznacza to, że korekta Bonferroniego nie jest dobrą metodą by przeprowadzać weryfikacje modeli takich jak model (D1), w przypadku, kiedy parametr średniej rozkładu obserwacji odstającej jest równy $(1 - \varepsilon)\sqrt{2 \ln(n)}$, gdzie $\varepsilon \in (0, 1)$.

5.1.3 Asymptotyczna moc testu opartego o lemat Neymana-Pearsona

Określmy, jaką moc osiąga test oparty o lemat Neymana-Pearsona. W tym celu, musimy dokonać zmian w modelu (D1), tak by hipoteza alternatywna była hipotezą prostą. Test oprzyjmy o dane związane z modelem (D1). Skonstruujmy model:

$$\begin{aligned} H_0 : \quad &\forall i \in \{1, \dots, n\} : \mu_i = 0, \\ H_1 : \quad &\exists! i \sim \mathcal{U}\{1, \dots, n\} : \mu_i = \hat{\mu}, \end{aligned} \tag{D2}$$

gdzie $\hat{\mu} = (1 - \varepsilon)\sqrt{2\ln(n)}$. Ustalmy poziom istotności $\alpha \in (0, 1)$. Oznaczmy $T_n := \sqrt{2\ln(n)}$. Wyznaczmy funkcje wiarygodności w oparciu o powyższe hipotezy:

$$\begin{aligned} L_0(0, (\underline{X}_1, \dots, \underline{X}_n)) &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}X_i^2\right) \\ &= \left(\frac{1}{\sqrt{2\pi}}\right)^n \prod_{i=1}^n \exp\left(-\frac{1}{2}X_i^2\right), \\ L_1(\hat{\mu}, (\underline{X}_1, \dots, \underline{X}_n)) &= \frac{1}{n} \sum_{i=1}^n \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(X_i - \hat{\mu})^2\right) \prod_{j:j \neq i} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}X_j^2\right) \\ &= \frac{1}{n} \left(\frac{1}{\sqrt{2\pi}}\right)^n \sum_{i=1}^n \exp\left(-\frac{1}{2}(X_i - \hat{\mu})^2\right) \prod_{j:j \neq i} \exp\left(-\frac{1}{2}X_j^2\right). \end{aligned}$$

Wyznaczmy iloraz wiarygodności:

$$r(H_0, H_1) = \frac{L_1(\hat{\mu}, (\underline{X}_1, \dots, \underline{X}_n))}{L_0(\hat{\mu}, (\underline{X}_1, \dots, \underline{X}_n))} = \frac{1}{n} \sum_{i=1}^n \exp\left(X_i \hat{\mu} - \frac{1}{2}\hat{\mu}^2\right).$$

Oznaczmy powyższy iloraz jako:

$$R = \frac{1}{n} \sum_{i=1}^n Y_i,$$

gdzie $Y_i = \exp\left(X_i \hat{\mu} - \frac{1}{2}\hat{\mu}^2\right)$. Prawdziwe jest stwierdzenie, że jeśli zmienna losowa: z , pochodząca z rozkładu $\mathcal{N}(\mu_z, \sigma_z^2)$, to e^z jest zmienną losową, pochodzącą z rozkładu logarytmicznie normalnego: $\mathcal{LN}(\mu_z, \sigma_z)$. Z tego wynika, że $Y_i \sim \mathcal{LN}\left(-\frac{1}{2}\hat{\mu}^2, \hat{\mu}\right)$.

W celu wyznaczenia asymptotycznej mocy rozważanego testu, nie możemy skorzystać z Centralnego Twierdzenia Granicznego, ponieważ $\hat{\mu}$ zależy od n . Musimy wykazać, że $R \xrightarrow{D} 1$, gdy hipoteza H_0 jest prawdziwa. Przyjmijmy zmienną:

$$R' = \frac{1}{n} \sum_{i=1}^n Y_i \mathbb{1}_{\{X_i \leq T_n\}}.$$

Prawdopodobieństwo, że $R \neq R'$ wynosi:

$$\begin{aligned} P(R \neq R') &\leq P\left(\max_{i \in \{1, \dots, n\}} X_i \geq T_n\right) \\ &\leq \sum_{i=1}^n P(X_i \geq T_n) \\ &\leq \sum_{i=1}^n \frac{1}{T_n} \cdot \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}T_n^2\right) \\ &\leq \sum_{i=1}^n \frac{1}{T_n} \cdot \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}\sqrt{2\ln(n)}^2\right) \\ &\leq n \cdot \frac{1}{T_n} \cdot \frac{1}{\sqrt{2\pi}} \\ &\xrightarrow{n \rightarrow \infty} 0. \end{aligned}$$

Obliczmy wartość oczekiwaną oraz variancję zmiennej R' :

$$\begin{aligned}
\mathbb{E}_{H_0} [R'] &= \frac{1}{n} \mathbb{E}_{H_0} \left[\sum_{i=1}^n Y_i \mathbb{1}_{\{X_i \leq T_n\}} \right] = \mathbb{E}_{H_0} [Y_1 \mathbb{1}_{\{X_1 \leq T_n\}}] \\
&= \int_{-\infty}^{T_n} \exp\left(z\hat{\mu} - \frac{1}{2}\hat{\mu}^2\right) \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}z^2\right) dz \\
&= \int_{-\infty}^{T_n} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(z - \hat{\mu})^2\right) dz \\
&= \int_{-\infty}^{T_n - \hat{\mu}} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}z^2\right) dz \\
&= \int_{-\infty}^{\varepsilon\sqrt{2\ln(n)}} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}z^2\right) dz \\
&= \Phi\left(\varepsilon\sqrt{2\ln(n)}\right)
\end{aligned}$$

$$\begin{aligned}
\text{Var}_{H_0} [R'] &= \frac{1}{n^2} \text{Var}_{H_0} \left[\sum_{i=1}^n Y_i \mathbb{1}_{\{X_i \leq T_n\}} \right] = \frac{1}{n} \text{Var}_{H_0} [Y_1 \mathbb{1}_{\{X_1 \leq T_n\}}] \\
&\leq \frac{1}{n} \mathbb{E}_{H_0} [Y_1^2 \mathbb{1}_{\{X_1 \leq T_n\}}] \\
&= \frac{1}{n} \int_{-\infty}^{T_n} \exp(2z\hat{\mu} - \hat{\mu}^2) \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}z^2\right) dz \\
&= \frac{1}{n} e^{\hat{\mu}^2} \int_{-\infty}^{T_n} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(z - 2\hat{\mu})^2\right) dz \\
&= \frac{1}{n} e^{\hat{\mu}^2} \int_{-\infty}^{T_n - 2\hat{\mu}} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(z)^2\right) dz \\
&= \frac{1}{n} e^{\hat{\mu}^2} \Phi(T_n - 2\hat{\mu}).
\end{aligned}$$

Przyjmijmy, że $T_n - 2\hat{\mu} < 0$, czyli $\varepsilon \in (0, \frac{1}{2})$. Wtedy:

$$\begin{aligned}
\frac{1}{n} e^{\hat{\mu}^2} \Phi(T_n - 2\hat{\mu}) &= \frac{1}{n} e^{\hat{\mu}^2} (1 - \Phi(-T_n + 2\hat{\mu})) \leq \frac{1}{n} e^{\hat{\mu}^2} \cdot \frac{\phi(2\hat{\mu} - T_n)}{T_n - 2\hat{\mu}} \\
&= \frac{1}{n} \cdot \frac{\exp\left(\hat{\mu}^2 - \frac{1}{2}(2\hat{\mu} - T_n)^2\right)}{(2\hat{\mu} - T_n) \sqrt{2\pi}} \\
&= \frac{1}{n} \cdot \frac{\exp(-2\varepsilon^2 \ln(n) + \ln(n))}{(2\hat{\mu} - T_n) \sqrt{2\pi}} \\
&= \frac{n^{-2\varepsilon^2}}{(2\hat{\mu} - T_n) \sqrt{2\pi}} \\
&\xrightarrow{n \rightarrow \infty} 0
\end{aligned}$$

Kiedy $T_n - 2\hat{\mu} > 0$, czyli $\varepsilon \in (\frac{1}{2}, 1)$, wtedy:

$$\frac{1}{n}e^{\hat{\mu}^2} = \frac{1}{n} \exp((1 - \varepsilon)^2 \cdot 2 \ln(n)) = \frac{1}{n} n^{2(1-\varepsilon)^2} \xrightarrow{n \rightarrow \infty} 0.$$

Z powyższego wynika, że:

$$\begin{aligned} \mathbb{E}_{H_0} [R'] &= \Phi(\varepsilon \sqrt{2 \ln(n)}), \\ \text{Var}_{H_0} [R'] &= o(1). \end{aligned}$$

Z nierówności Czebyszewa następuje:

$$R' = \Phi(\varepsilon \sqrt{2 \ln(n)}) + o(1).$$

Dowodzi to, że:

$$R \xrightarrow{D} 1.$$

Z tego wynika, że $L_0(0, (\underline{X}_1, \dots, \underline{X}_n)) \approx L_1(\hat{\mu}, (\underline{X}_1, \dots, \underline{X}_n))$, przy $n \rightarrow \infty$. Oznacza to, że H_0 i H_1 są takie same.

Oznaczmy zmienną G , taką, że $\lim_{n \rightarrow \infty} P(R \geq G | H_0) = \alpha$. Obliczmy asymptotyczną moc testu opartego o lemat Neymana-Pearsona:

$$\begin{aligned} P(R > G | H_1) &= 1 - P(R \leq G | H_1) \\ &= 1 - \int \mathbb{1}_{\{R \leq G\}} dP_{H_1} \\ &= \left| R = \frac{dP_{H_1}}{dP_{H_0}} \right| = 1 - \int R \mathbb{1}_{\{R \leq G\}} dP_{H_0} \\ &= 1 - \int \mathbb{1}_{\{R \leq G\}} dP_{H_0} - \int (R - 1) \mathbb{1}_{\{R \leq G\}} dP_{H_0} \\ &\xrightarrow{n \rightarrow \infty} = 1 - (1 - \alpha) - 0 = \alpha. \end{aligned}$$

Test modelu typu „igła w stogu siana”, oparty o lemat Neymana-Pearsona osiąga asymptotycznie moc równą α . Z faktu, że lemat Neymana-Pearsona pozwala nam skonstruować test jednostajnie najmocniejszy na poziomie istotności α , wynika, że nie możemy wyznaczyć testu o mocy większej niż α w przypadku kiedy, $\hat{\mu} = (1 - \varepsilon)\sqrt{2 \ln(n)}$. Konsekwencją tego jest fakt, że jeżeli dysponujemy zbiorem danych pochodzących z rozkładu $\mathcal{N}(0, 1)$ i chcemy sprawdzić, czy w zbiorze danych znajduje się obserwacja o rozkładzie $\mathcal{N}(\mu, 1)$, to wykryjemy ją tylko wtedy, kiedy $\mu = (1 + \varepsilon)\sqrt{2 \ln(n)}$.

5.2 Model oparty o dane pochodzące z rozkładu Poissona

Postaramy się przeprowadzić weryfikację modelu typu „igła w stogu siana”, opartego o dane pochodzące z rozkładu Poissona. Niech zbiór danych, o który oprzemy późniejsze rozważania, składa się z niezależnych zmiennych losowych: $\underline{X}_1, \dots, \underline{X}_n$, gdzie $n \geq 2$. Przyjmijmy, że \underline{X}_i , dla $i \in \{1, \dots, n\}$, pochodzi z rozkładu Poissona o parametrze λ_i oraz, że statystyka testowa:

$$T(\underline{X}_i) = \underline{X}_i,$$

odpowiada i -tej hipotezie zerowej: $H_{0_i} : \lambda_i = \tilde{\lambda}$, gdzie $\tilde{\lambda} \in (0, \infty)$, jest ustaloną wartością. Wiemy, że przy założeniu prawdziwości hipotezy H_{0_i} , $T(\underline{X}_i) \sim \text{Poisson}(\tilde{\lambda})$. Wyznamy model typu „igła w stogu siana” dla tego przypadku:

$$\begin{aligned} H_0 : \quad & \forall i \in \{1, \dots, n\} : \lambda_i = \tilde{\lambda}, \\ H_1 : \quad & \exists i \sim \mathcal{U}\{1, \dots, n\} : \lambda_i > \tilde{\lambda}. \end{aligned} \tag{E1}$$

Odrzucamy globalną hipotezę zerową H_0 , na ustalonym poziomie istotności α , przy pomocy korekty Bonferroniego, gdy:

$$\max_{i \in \{1, \dots, n\}} T(X_i) \geq k,$$

gdzie wartość k definiujemy jako:

$$k := \min \left\{ m : \sum_{i=m}^{\infty} \frac{\tilde{\lambda}^i e^{-\tilde{\lambda}}}{i!} \leq \frac{\alpha}{n} \right\}.$$

Ponieważ rozkład Poissona jest rozkładem dyskretnym, określenie wartości k w taki sposób pozwala nam na znalezienie kwantyla rzędu $1 - \frac{\alpha}{n}$, lub wyższego, rozkładu $Poisson(\tilde{\lambda})$. W skutek tego, kontrolujemy prawdopodobieństwo popelnienia błędu I rodzaju na ustalonym poziomie α .

5.2.1 Asymptotyczna wartość kwantyla

Aby rozpocząć rozważania na temat efektywności korekty Bonferroniego oraz testu opartego o lemat Neymana-Pearsona, dla wyżej skonstruowanego modelu typu „igła w stogu siana”, musimy wyznaczyć optymalny próg wykrycia. Wiemy, że dystrybuenta rozkładu Poissona o parametrze λ jest określona i ma postać:

$$\mathbb{F}_\lambda(T') = P(t \leq T') = \sum_{i=0}^{T'} \frac{\lambda^i e^{-\lambda}}{i!},$$

dla określonej liczby $T' \in \{0, 1, 2, \dots\}$, oraz zmiennej losowej $t \sim Poisson(\lambda)$. Z tego wynika, że:

$$1 - \mathbb{F}_\lambda(T') = P(t > T') = \sum_{i=T'+1}^{\infty} \frac{\lambda^i e^{-\lambda}}{i!}.$$

Znajomość dystrybenty rozkładu Poissona pozwala nam, w stosunku do rozkładu standardowego normalnego, szybciej i łatwiej znaleźć przybliżenie kwantyla rzędu $1 - \frac{\alpha}{n}$ rozkładu $Poisson(\tilde{\lambda})$, gdzie $\alpha \in (0, 1)$, $n \in \mathbb{N}_+$. Niech $T' \in \{0, 1, 2, \dots\}$ i $t \sim Poisson(\tilde{\lambda})$, $T = T' + 1$. Przy pomocy wzoru Stirlinga wyznaczmy:

$$\begin{aligned} 1 - \mathbb{F}_{\tilde{\lambda}}(T') = P(t \geq T) &= \sum_{i=T}^{\infty} \frac{\tilde{\lambda}^i e^{-\tilde{\lambda}}}{i!} \\ &= \frac{\tilde{\lambda}^T e^{-\tilde{\lambda}}}{T!} \left(1 + \frac{\tilde{\lambda}}{T+1} + \frac{\tilde{\lambda}^2}{(T+1)(T+2)} + \dots \right) \\ &\approx \tilde{\lambda}^T e^{-\tilde{\lambda}} \frac{e^T}{T^T \sqrt{2\pi T}} \left(1 + \frac{\tilde{\lambda}}{T+1} + \frac{\tilde{\lambda}^2}{(T+1)(T+2)} + \dots \right). \end{aligned}$$

Oznaczmy $A_{T, \tilde{\lambda}} = \left(1 + \frac{\tilde{\lambda}}{T+1} + \frac{\tilde{\lambda}^2}{(T+1)(T+2)} + \dots \right)$. Oszacujmy $A_{T, \tilde{\lambda}}$, wiedząc, że:

$$1 \leq A_{T, \tilde{\lambda}} \leq \sum_{i=0}^{\infty} \frac{\tilde{\lambda}^i}{i!}.$$

Z tego wynika, że $A_{T,\tilde{\lambda}} \in [1, e^{\tilde{\lambda}}]$. Dla $n \rightarrow \infty$, wyznaczmy przybliżenie wartości T :

$$\begin{aligned} P(t \geq T) &\approx \frac{\alpha}{n} \\ \tilde{\lambda}^T e^{-\tilde{\lambda}} \frac{e^T}{T^T \sqrt{2\pi T}} A_{T,\tilde{\lambda}} &\approx \frac{\alpha}{n} \quad \backslash \ln(\cdot) \\ T \ln(\tilde{\lambda}) - \tilde{\lambda} + T - \left(T + \frac{1}{2}\right) \ln(T) - \frac{1}{2} \ln(2\pi) + \ln(A_{T,\tilde{\lambda}}) &\approx \ln(\alpha) - \ln(n), \end{aligned}$$

dla $n \rightarrow \infty$, jest to równoważne z:

$$\begin{aligned} T \ln(\tilde{\lambda}) + T - \left(T + \frac{1}{2}\right) \ln(T) &\approx -\ln(n) \\ T \left(\ln(\tilde{\lambda}) + 1\right) - \left(T + \frac{1}{2}\right) \ln(T) &\approx -\ln(n). \end{aligned}$$

Niech:

$$T \approx b_n \ln(n).$$

Z powyższego wynika, że:

$$\begin{aligned} -b_n \ln(n) \left(\ln(\tilde{\lambda}) + 1\right) + \left(b_n \ln(n) + \frac{1}{2}\right) \ln(b_n \ln(n)) &\approx \ln(n) \\ -b_n \ln(\tilde{\lambda}) - b_n + \left(b_n + \frac{1}{2\ln(n)}\right) (\ln(b_n) + \ln(\ln(n))) &\approx 1 \\ -b_n \ln(\tilde{\lambda}) - b_n + b_n \ln(b_n) + b_n \ln(\ln(n)) + \frac{\ln(b_n)}{2\ln(n)} + \frac{\ln(\ln(n))}{2\ln(n)} &\approx 1 \\ b_n \left(-\ln(\tilde{\lambda}) + \ln(\ln(n))\right) - b_n + b_n \ln(b_n) + \frac{\ln(b_n)}{2\ln(n)} + \frac{\ln(\ln(n))}{2\ln(n)} &\approx 1. \end{aligned}$$

Posługując się następującym faktem z analizy matematycznej:

$$d_n \approx e_n \iff \frac{d_n}{e_n} \xrightarrow{n \rightarrow \infty} 1,$$

Zdefiniujmy:

$$b_n := \frac{1}{\ln(\ln(n))}.$$

W związku z czym:

$$\begin{aligned} \frac{-\ln(\tilde{\lambda}) + \ln(\ln(n))}{\ln(\ln(n))} - \frac{1}{\ln(\ln(n))} - \frac{\ln(\ln(\ln(n)))}{\ln(\ln(n))} \\ - \frac{\ln(\ln(\ln(n)))}{2\ln(n)} + \frac{\ln(\ln(n))}{2\ln(n)} &\approx 1. \end{aligned}$$

Przy czym:

$$\begin{aligned} \frac{-\ln(\tilde{\lambda}) + \ln(\ln(n))}{\ln(\ln(n))} &\xrightarrow{n \rightarrow \infty} 1, \\ -\frac{1}{\ln(\ln(n))} &\xrightarrow{n \rightarrow \infty} 0, \\ -\frac{\ln(\ln(\ln(n)))}{\ln(\ln(n))} &\xrightarrow{n \rightarrow \infty} 0, \\ -\frac{\ln(\ln(\ln(n)))}{2\ln(n)} &\xrightarrow{n \rightarrow \infty} 0, \\ \frac{\ln(\ln(n))}{2\ln(n)} &\xrightarrow{n \rightarrow \infty} 0. \end{aligned}$$

Z powyższych rozważań wynika, że przybliżenie wartości T jest postaci:

$$T \approx \frac{\ln(n)}{\ln(\ln(n))}.$$

W porównaniu do progu wykrycia z poprzedniego podrozdziału, również w tym przypadku przybliżenie kwantyla rzędu $1 - \frac{\alpha}{n}$, dla dużych n , nie zależy od wartości ustalonego wcześniej poziomu istotności α oraz parametru $\tilde{\lambda}$.

5.2.2 Asymptotyczna moc korekty Bonferroniego

Znając aproksymację kwantyla rzędu $1 - \frac{\alpha}{n}$ rozkładu Poissona o parametrze: $\lambda = \tilde{\lambda}$, wyznaczmy asymptotyczną moc korekty Bonferroniego w oparciu o weryfikację modelu typu „igła w stogu siana”. Zanim do tego przejdziemy, sprawdźmy, czy korekta Bonferroniego kontroluje ryzyko błędu I rodzaju na ustalonym poziomie istotności $\alpha \in (0, 1)$. Kolejne rozważania oprzyjmy o model (E1). Przyjmując zmienną losową $t \sim \text{Poisson}(\tilde{\lambda})$, wyznaczmy p -wartość odpowiadającą statystyce $T(\underline{X}_i)$:

$$p_i = 1 - P(t \leq T(\underline{X}_i)) = 1 - \sum_{i=0}^{\lfloor T(\underline{X}_i) \rfloor} \frac{\tilde{\lambda}^i e^{-\tilde{\lambda}}}{i!}.$$

Z rozważań przeprowadzonych w rozdziale trzecim wynika, że p_i przyjmuje asymptotycznie rozkład $\mathcal{U}[0, 1]$. Oznaczmy:

$$k := \max \left\{ m : \sum_{i=m+1}^{\infty} \frac{\tilde{\lambda}^i e^{-\tilde{\lambda}}}{i!} > \frac{\alpha}{n} \right\},$$

$$T_n := \frac{\ln(n)}{\ln(\ln(n))},$$

takie, że:

$$k \approx T_n.$$

Rozpatrzmy i udowodnijmy poniższe twierdzenie.

Twierdzenie 5.2. *Korekta Bonferroniego kontroluje prawdopodobieństwo popełnienia błędu I rodzaju na poziomie $\alpha \in (0, 1)$, dla modelu typu „igła w stogu siana”, opartego o dane pochodzące z rozkładu Poissona.*

Dowód. Następujący dowód oprzyjmy o model (E1) oraz dane z nim związane. Przyjmijmy poziom istotności $\alpha \in (0, 1)$. Niech statystyka testowa: $T(\underline{X}_i) \sim \text{Poisson}(\tilde{\lambda})$, gdzie $i \in \{1, \dots, n\}$. Prawdopodobieństwo popełnienia błędu I rodzaju wynosi:

$$\begin{aligned} P\left(\min_{i \in \{1, \dots, n\}} p_i \leq \frac{\alpha}{n} \mid H_0\right) &= 1 - P\left(\bigcap_{i=1}^n \left\{p_i \geq \frac{\alpha}{n}\right\} \mid H_0\right) \\ &= 1 - \prod_{i=1}^n P\left(p_i \geq \frac{\alpha}{n} \mid H_0\right). \end{aligned}$$

Przy czym:

$$\prod_{i=1}^n P\left(p_i \geq \frac{\alpha}{n} \mid H_0\right) = \prod_{i=1}^n P(T(\underline{X}_i) < k \mid H_0) = \prod_{i=1}^n \mathbb{F}_{\tilde{\lambda}}(k) = \left(\mathbb{F}_{\tilde{\lambda}}(k)\right)^n,$$

gdzie:

$$\begin{aligned} 1 - \frac{\alpha}{n} &\leq \mathbb{F}_{\tilde{\lambda}}(k) \leq 1 - \frac{\alpha}{n} + \frac{\tilde{\lambda}^k e^{-\tilde{\lambda}}}{k!} \\ \left(1 - \frac{\alpha}{n}\right)^n &\leq \left(\mathbb{F}_{\tilde{\lambda}}(k)\right)^n \leq \left(1 - \frac{\alpha}{n} + \frac{\tilde{\lambda}^k e^{-\tilde{\lambda}}}{k!}\right)^n. \end{aligned}$$

Przy $n \rightarrow \infty$:

$$\exp(-\alpha) \leq \left(\mathbb{F}_{\tilde{\lambda}}(k)\right)^n \leq \exp\left(n\left(-\frac{\alpha}{n} + \frac{\tilde{\lambda}^k e^{-\tilde{\lambda}}}{k!}\right)\right),$$

gdzie:

$$\exp\left(n\left(-\frac{\alpha}{n} + \frac{\tilde{\lambda}^k e^{-\tilde{\lambda}}}{k!}\right)\right) = \exp(-\alpha) \exp\left(\frac{n\tilde{\lambda}^k e^{-\tilde{\lambda}}}{k!}\right) \xrightarrow{n \rightarrow \infty} \infty,$$

ponieważ:

$$\begin{aligned} \frac{n\tilde{\lambda}^k e^{-\tilde{\lambda}}}{k!} &\xrightarrow{n \rightarrow \infty} n\tilde{\lambda}^k e^{-\tilde{\lambda}} e^k k^{-k} (2\pi k)^{-\frac{1}{2}} \\ &\xrightarrow{n \rightarrow \infty} \frac{e^{-\tilde{\lambda}}}{\sqrt{2\pi}} \exp\left(\ln(n) + k \ln(\tilde{\lambda}) + k - k \ln(k) - \frac{1}{2} \ln(k)\right) \\ &\xrightarrow{n \rightarrow \infty} \frac{e^{-\tilde{\lambda}}}{\sqrt{2\pi}} \exp\left(\ln(n) + T_n \ln(\tilde{\lambda}) + T_n - T_n \ln(T_n) - \frac{1}{2} \ln(T_n)\right) \\ &\xrightarrow{n \rightarrow \infty} \frac{e^{-\tilde{\lambda}}}{\sqrt{2\pi}} \exp\left(-T_n \left(-\frac{\ln(n)}{T_n} - \ln(\tilde{\lambda}) - 1 + \ln(T_n) + \frac{\ln(T_n)}{2T_n}\right)\right) \\ &\xrightarrow{n \rightarrow \infty} \frac{e^{-\tilde{\lambda}}}{\sqrt{2\pi}} \exp\left(\frac{\ln(n)}{\ln(\ln(n))} \ln(\ln(\ln(n)))\right) \xrightarrow{n \rightarrow \infty} \infty. \end{aligned}$$

Z powyższego wynika, że, przy $n \rightarrow \infty$:

$$\begin{aligned} P\left(\min_{i \in \{1, \dots, n\}} p_i \leq \frac{\alpha}{n} \mid H_0\right) &= 1 - \prod_{i=1}^n P(T(\underline{X}_i) < k \mid H_0) \\ &\leq 1 - \left(1 - \frac{\alpha}{n}\right)^n \approx 1 - e^{-\alpha} \\ &\leq \alpha. \end{aligned}$$

■

Uwaga 5.1. Porównując powyższe modele, ryzyko popełnienia błędu I rodzaju, dla modelu opartego o dane pochodzące z rozkładu ciągłego, równe jest $1 - e^{-\alpha}$, natomiast, w przypadku modelu opartego o dane pochodzące z rozkładu dyskretnego, nie możemy określić dokładnej wartości prawdopodobieństwa popełnienia błędu I rodzaju (wiemy, że jest ono mniejsze lub równe $1 - e^{-\alpha}$).

Wykazaliśmy, że dla modelu (E1), korekta Bonferroniego także kontroluje błąd I rodzaju. Postaramy się wyznaczyć asymptotyczną moc rozważanej procedury. Ustalmy poziom istotności $\alpha \in (0, 1)$. Załóżmy, bez straty ogólności rozważań, że $\lambda_1 > \lambda$. Tak jak w poprzednim podrozdziale, rozważmy dwa przypadki:

1. Niech $\lambda_1 = (1 - \varepsilon) \frac{\ln(n)}{\ln(\ln(n))}$, gdzie $\varepsilon \in (0, 1)$. Wyznaczmy asymptotyczną moc korekty Bonferroniego, w oparciu o powyższy model:

$$\begin{aligned} P\left(\min_{i \in \{1, \dots, n\}} p_i \leq \frac{\alpha}{n} \mid H_1\right) &\leq P\left(p_1 \leq \frac{\alpha}{n} \mid H_1\right) + P\left(\min_{i \in \{2, \dots, n\}} p_i \leq \frac{\alpha}{n} \mid H_1\right) \\ &= P(\underline{X}_1 > k \mid H_1) + 1 - \prod_{i=2}^n P\left(p_i \geq \frac{\alpha}{n} \mid H_1\right) \\ &\leq P(\underline{X}_1 > T_n \mid H_1) + 1 - \left(1 - \frac{\alpha}{n}\right)^{n-1}. \end{aligned}$$

Przy czym:

$$\begin{aligned} P(\underline{X}_1 > k \mid H_1) &= P(\underline{X}_1 > T_n \mid H_1) \\ &= P\left(\frac{\underline{X}_1 - \lambda_1}{\sqrt{\lambda_1}} > \frac{T_n - \lambda_1}{\sqrt{\lambda_1}} \mid H_1\right) \\ &= P\left(\frac{\underline{X}_1 - \lambda_1}{\sqrt{\lambda_1}} > \frac{\varepsilon}{\sqrt{1 - \varepsilon}} \sqrt{T_n} \mid H_1\right). \end{aligned}$$

W oparciu o fakt, że rozkład $\mathcal{N}(\lambda, \lambda)$ jest doskonałym przybliżeniem rozkładu *Poisson*(λ), gdy $\lambda \rightarrow \infty$, wynika, że:

$$\frac{\underline{X}_1 - \lambda_1}{\sqrt{\lambda_1}} \xrightarrow{D} \mathcal{N}(0, 1).$$

W związku z:

$$\frac{\varepsilon}{\sqrt{1 - \varepsilon}} \sqrt{T_n} \xrightarrow{n \rightarrow \infty} \infty,$$

dochodzimy do:

$$P\left(\frac{\underline{X}_1 - \lambda_1}{\sqrt{\lambda_1}} > \frac{\varepsilon}{\sqrt{1 - \varepsilon}} \sqrt{T_n} \mid H_1\right) \xrightarrow{n \rightarrow \infty} 0.$$

Wobec powyższego, asymptotyczna moc korekty Bonferroniego wynosi w tym przypadku:

$$P\left(\min_{i \in \{1, \dots, n\}} p_i \leq \frac{\alpha}{n} \mid H_1\right) \leq 0 + 1 - e^{-\alpha} \leq \alpha.$$

2. Niech $\lambda_1 = (1 + \varepsilon) \frac{\ln(n)}{\ln(\ln(n))}$, gdzie $\varepsilon \in (0, 1)$. Wyznaczmy asymptotyczną moc korekty

Bonferroniego:

$$\begin{aligned}
P\left(\min_{i \in \{1, \dots, n\}} p_i \leq \frac{\alpha}{n} \mid H_1\right) &\geq P(p_1 \geq \frac{\alpha}{n} \mid H_1) \\
&= P(\underline{X}_1 > k \mid H_1) \\
&= P(\underline{X}_1 > T_n \mid H_1) \\
&= P\left(\frac{X_1 - \lambda_1}{\sqrt{\lambda_1}} > \frac{T_n - \lambda_1}{\sqrt{\lambda_1}} \mid H_1\right) \\
&= P\left(\frac{X_1 - \lambda_1}{\sqrt{\lambda_1}} > \frac{-\varepsilon}{\sqrt{1 + \varepsilon}} \sqrt{T_n} \mid H_1\right).
\end{aligned}$$

Na mocy faktów:

$$\begin{aligned}
\frac{X_1 - \lambda_1}{\sqrt{\lambda_1}} &\xrightarrow{D} \mathcal{N}(0, 1), \\
\frac{-\varepsilon}{\sqrt{1 + \varepsilon}} \sqrt{T_n} &\xrightarrow{n \rightarrow \infty} -\infty,
\end{aligned}$$

asymptotyczna moc procedury wynosi:

$$P\left(\min_{i \in \{1, \dots, n\}} p_i \leq \frac{\alpha}{n} \mid H_1\right) \geq P\left(\frac{X_1 - \lambda_1}{\sqrt{\lambda_1}} > \frac{-\varepsilon}{\sqrt{1 + \varepsilon}} \sqrt{T_n} \mid H_1\right) \xrightarrow{n \rightarrow \infty} 1.$$

Z powyższych podpunktów wynika, że korekta Bonferroniego wykryje obserwację pochodzącą z rozkładu $Poisson(\lambda)$ spośród obserwacji z rozkładu $Poisson(\lambda')$, jeśli λ będzie równa $(1 + \varepsilon) \frac{\ln(n)}{\ln(\ln(n))}$. W przypadku, kiedy $\lambda = (1 - \varepsilon) \frac{\ln(n)}{\ln(\ln(n))}$, test oparty o rozważaną procedurę osiąga asymptotycznie moc mniejszą niż α . W rezultacie, ryzyko popełnienia błędu II rodzaju jest zbyt duże, żeby test uznać za skuteczny. Dla modeli typu „igła w stogu siana”, opartych na danych pochodzących z rozkładu Poissona, korekta Bonferroniego nie jest dobrą procedurą weryfikacji. Przyczyną tego stwierdzenia jest fakt, że nie zawsze procedura ta osiąga wysoką moc.

6 Wnioski

Podsumujmy wszystkie rozważania jakie przeprowadziliśmy we wcześniejszych rozdziałach. Podstawy teorii testowania hipotez statystycznych pomogły nam zrozumieć ideę tworzenia testów, oraz nauczyły nas ich konstrukcji. Dzięki przeprowadzonym testom udowodniliśmy, że p -wartość pochodzi z rozkładu $\mathcal{U}[0, 1]$, kiedy statystyka testowa pochodzi z rozkładu ciągłego. Wykazaliśmy również, że jeśli statystyka testowa pochodzi z rozkładu Poissona, to jej przekształcenie osiąga asymptotycznie rozkład jednostajny na przedziale $[0, 1]$. Rozważania na temat wartości p pozwoliły nam wykazać, że korekta Bonferroniego kontroluje błąd polegający na odrzuceniu globalnej hipotezy zerowej w sytuacji, kiedy jest ona prawdziwa.

W konsekwencji, skonstruowaliśmy modele typu „igła w stogu siana” dla danych pochodzących z rozkładów normalnego oraz Poissona. Dla rozkładu normalnego wyznaczyliśmy asymptotyczną wartość kwantyla rzędu $1 - \frac{\alpha}{n}$ (gdzie: $\alpha \in (0, 1)$, $n \rightarrow \infty$), czyli asymptotyczny próg wykrycia, równą $\sqrt{2 \ln(n)}$. Wykazaliśmy, że korekta Bonferroniego osiąga asymptotycznie moc równą 1, w sytuacji, gdy parametr średniej, obserwacji różniący się rozkładem, jest równy: $\mu = (1 + \varepsilon) \sqrt{2 \ln(n)}$, gdzie $\varepsilon \in (0, 1)$. W przypadku kiedy ten sam parametr jest równy $(1 - \varepsilon) \sqrt{2 \ln(n)}$, pokazaliśmy, że korekta Bonferroniego osiąga moc bliską α . Za pomocą lematu Neymana-Pearsona pokazaliśmy, że nie możemy wyznaczyć testu o mocy większej niż α , gdy hipoteza alternatywna zakłada, że spośród $n - 1$ obserwacji o rozkładzie $\mathcal{N}(0, 1)$ znajduje się obserwacja o rozkładzie $\mathcal{N}((1 - \varepsilon) \sqrt{2 \ln(n)}, 1)$, gdzie $\varepsilon \in (0, 1)$.

Dla modelu opartego o n obserwacji z rozkładu Poissona, wyznaczyliśmy próg wykrycia równy $\frac{\ln(n)}{\ln(\ln(n))}$. Przeprowadziliśmy test modelu i wykazaliśmy, że, tak jak w poprzednim przypadku, korekta Bonferroniego przyjmuje asymptotycznie moc równą 1, kiedy parametr λ rozkładu odstającej obserwacji jest równy $(1 + \varepsilon)\frac{\ln(n)}{\ln(\ln(n))}$. Podobnie wykazaliśmy, że moc korekty Bonferroniego jest mniejsza od poziomu istotności α , gdy obserwacja różni się rozkładem, pochodzi z rozkładu *Poisson* $\left(\left(1 - \varepsilon\right)\frac{\ln(n)}{\ln(\ln(n))}\right)$.

Literatura

- [1] R. V. Hogg, J. W. McKean, A. T. Craig. *Introduction to Mathematical Statistics (6th Edition)*. Pearson Education, 2005.
- [2] N. L. Johnson, A. W. Kemp, S. Kotz. *Univariate Discrete Distributions(Third Edition)*. John Wiley & Sons, Inc. Hoboken, New Jersey. 2005.
- [3] A. Futschik, T. Taus, S. Zehetmayer. *An omnibus test for the global null hypothesis*. <https://journals.sagepub.com/doi/pdf/10.1177/0962280218768326>. Dostęp: 30.01.2023.
- [4] E. Candes. *STATS 300C: Theory of Statistics, Spring 2022. Lecture 2*. 31.03.2022. <https://candes.su.domains/teaching/stats300c/Lectures/Lecture02.pdf>. Dostęp: 30.01.2023.