

Rzeczpospolita
PolskaPolitechnika
WarszawskaUnia Europejska
Europejski Fundusz Społeczny

Rachunek Prawdopodobieństwa i Elementy Statystyki Matematycznej

Anna Dembińska

Wydział Matematyki i Nauk Informacyjnych

Wykład 13

Projekt „NERW 2 PW. Nauka – Edukacja – Rozwój – Współpraca”
współfinansowany jest ze środków Unii Europejskiej w ramach
Europejskiego Funduszu Społecznego.

Zadanie 10 pn. „Modyfikacja programów studiów na kierunkach
prowadzonych przez Wydział Matematyki i Nauk Informacyjnych”,
realizowane w ramach projektu „NERW 2 PW. Nauka – Edukacja –
Rozwój – Współpraca”, współfinansowanego ze środków Unii
Europejskiej w ramach Europejskiego Funduszu Społecznego.

13. ANALIZA ZGODNOŚCI

TESTOWANIE ZGODNOŚCI Z DOWOLNYM ROZKŁADEM

1. Test zgodności χ^2 -Pearsona

H_0 : badana próba losowa pochodzi z zadanego rozkładu (lub rodziny rozkładów)

H_1 : badana próba losowa nie pochodzi z zadanego rozkładu (lub rodziny rozkładów)

Statystyka testowa $\chi^2 = \sum_{j=1}^k \frac{(n_j - np_j^0)^2}{np_j^0}$, gdzie k - ilość klas; p_j^0 - prawdopodobieństwa teoretyczne wpadnięcia obserwacji do j -tej klasy przy założeniu prawdziwości H_0 (jeśli H_0 nie jest hipotezą prostą, to brakujące parametry rozkładu z H_0 wyznaczamy metodą największej wiarygodności), n_j - liczba obserwacji, które znalazły się w j -tej klasie, n - licznosc próby.

Zbiór krytyczny $W = \langle \chi_{1-\alpha, k-1-r}^2; +\infty \rangle$, gdzie r jest ilością parametrów szacowanych z próby, zaś $\chi_{1-\alpha, k-1-r}^2$ to kwantyl rzędu $1 - \alpha$ rozkładu chi-kwadrat o $k - 1 - r$ stopniach swobody.

Jeżeli wyznaczona wartość statystyki χ^2 należy do zbioru krytycznego W , to, na poziomie istotności α , H_0 odrzucamy.

UWAGI dotyczące testu zgodności χ^2 -Pearsona:

- Statystyka testowa χ^2 testu zgodności χ^2 -Pearsona w sytuacji, gdy hipoteza H_0 jest prawdziwa, ma jedynie w przybliżeniu rozkład χ^2 o $k - 1 - r$ stopniach swobody. Przybliżenie to uznajemy za dopuszczalne, gdy wszystkie $np_j^0 \geq 5$, a za dobre, gdy wszystkie $np_j^0 \geq 10$. Gdyby np_j^0 nie były duże, to wtedy przybliżenie rozkładem χ^2 nie działa i pozostaje symulacyjne wyznaczanie rozkładu statystyki testowej.
- Gdy testem zgodności χ^2 -Pearsona sprawdzamy zgodność z rozkładem ciągłym, to podczas jego dyskretyzacji, końce przedziałów klas wybieramy tak, by prawdopodobieństwa klas p_j^0 były przynajmniej w przybliżeniu równe i by był spełniony warunek, że wszystkie $np_j^0 \geq 5$.

Test zgodności χ^2 -Pearsona jest zaimplementowany w R, niestety jedynie dla prostych hipotez H_0 :

```
> chisq.test(x, p, simulate.p.value=FALSE, B=2000)
```

gdzie

- x to wektor z licznosciami poszczególnych klas,
- p to wektor z prawdopodobieństwami teoretycznymi p_j^0 poszczególnych klas,

- `simulate.p.value` musimy ustawić na `TRUE` jeśli chcemy symulacyjnie wyznaczyć wartość p-value testu, wtedy zostanie przeprowadzonych domyślnie `B=2000` losowań realizacji próbki losowej.

Przykład 13.1. 100 losowo wybranych studentów zapytano ile znają języków obcych. Otrzymane wyniki zapisano w poniższej tabeli:

liczba języków	liczba studentów
0	25
1	40
2	30
3	5

Czy na podstawie powyższych danych można uznać, że rozkład liczby języków obcych, którymi posługują się studenci, jest (a) rozkładem Poissona o średniej równej 1, (b) rozkładem Poissona? Przyjąć poziom istotności 0.01.

Rozwiązanie przykładu 13.1:

- (a) H_0 : badana próba losowa pochodzi z rozkładu Poissona o średniej 1,
 H_1 : badana próba losowa nie pochodzi z rozkładu Poissona o średniej 1.

W hipotezie H_0 jest rozkład Poissona o średniej równej 1. Korzystając z tego, że jeśli X ma rozkład Poissona z parametrem λ , to $EX = \lambda$, otrzymujemy $\lambda = EX = 1$.

Szukamy wektora prawdopodobieństw teoretycznych poszczególnych klas przy założeniu prawdziwości H_0 , tzn. wektora $(p_0^0, p_1^0, p_2^0, p_3^0)$, gdzie

$$p_0^0 = P(X = 0), p_1^0 = P(X = 1), p_2^0 = P(X = 2), p_3^0 = P(X \geq 3) = P(X > 2)$$

$$\text{ i } X \sim \text{Poiss}(\lambda = 1).$$

Przypomnijmy, że jeśli X ma rozkład Poissona z parametrem λ , to

- `dpois(x=k, lambda=λ)` podaje prawdopodobieństwo $P(X = k)$;
- `ppois(x=k, lambda=λ, lower.tail=FALSE)` podaje prawdopodobieństwo $P(X > k)$.

Zatem szukany wektor prawdopodobieństw $(p_0^0, p_1^0, p_2^0, p_3^0)$ otrzymamy pisząc

```
> prawdep0 <- c(dpois(c(0,1,2),lambda=1),
                ppois(2,lambda=1,lower.tail=FALSE))
```

Zauważamy, że możemy użyć testu zgodności χ^2 -Pearsona, bo wszystkie $np_j^0 \geq 5$ i tylko jedno np_j^0 nie spełnia warunku $np_j^0 \geq 10$. Rzeczywiście

```
> 100*prawdep0
```

daje następujące wartości: 36.78794, 36.78794, 18.39397, 8.03014.

Potrzebujemy także wektora z zaobserwowanymi licznosciami

```
> licznosci <- c(25,40,30,5)
Przeprowadzamy test zgodności  $\chi^2$ -Pearsona
> chisq.test(x=licznosci, p=prawdop0)
Odczytujemy p wartość:
```

$$p - value = 0,005787 < \alpha = 0,01 \Rightarrow \text{odrzucaamy } H_0,$$

gdzie $\alpha = 0,01$ to poziom istotności testu. Wyciągamy więc wniosek, że rozkład liczby języków obcych, którymi posługują się studenci, nie jest rozkładem Poissona o średniej równej 1.

(b) H_0 : badana próba losowa pochodzi z rozkładu Poissona,
 H_1 : badana próba losowa nie pochodzi z rozkładu Poissona.

Teraz H_0 jest hipotezą złożoną, więc musimy zacząć od oszacowania parametru λ , używając metody największej wiarygodności. Osoby znające tę metodę mogą pokazać, że estymatorem największej wiarygodności parametru λ rozkładu Poissona $Poiss(\lambda)$ jest

$$\hat{\lambda}_{NW} = \bar{X} = \frac{1}{n} \sum_{i=1}^n X_i.$$

Wartość tego estymatora wyznaczymy w R pisząc:

```
> dane <- c(rep(0,25),rep(1,40),rep(2,30),rep(3,5))
lub krócej
> dane <- rep(0:3,licznosci)
> estymator.lambdy <- mean(dane)
Możemy też wyznaczyć wartość tego estymatora automatycznie:
> library(MASS)
> estymator.lambdy <- fitdistr(dane,"Poisson")$estimate # 1.15
```

Wektor prawdopodobieństw $(p_0^0, p_1^0, p_2^0, p_3^0)$ liczymy w R pisząc

```
> prawdop0.b <- c(dpois(c(0,1,2),lambda=estymator.lambdy),
+ ppois(2,lambda=estymator.lambdy,lower.tail=FALSE))
```

Widzimy, że wszystkie $np_j^0 \geq 10$:

```
> 100*prawdop0.b # 31.66368 36.41323 20.93761 10.98549
```

wiec możemy użyć testu zgodności χ^2 -Pearsona. Wyliczamy statystykę testową tego testu:

```
> chisq.test(x=licznosci, p=prawdop0.b)
otrzymując
X-squared = 8.9394
```

Następnie wyznaczamy zbiór krytyczny $W = \langle \chi^2_{1-\alpha, k-1-r}; +\infty \rangle$, gdzie $\alpha = 0,01$, $k = 4$ i $r = 1$ (bo szacowaliśmy jeden parametr czyli λ). Zatem szukamy kwantyla

$$\chi^2_{1-\alpha, k-1-r} = \chi^2_{1-0,01; 4-1-1} = \chi^2_{0,99; 2},$$

i znajdujemy go korzystając z R:

```
> qchisq(0.99, 2)
```

Otrzymujemy $\chi^2_{0,99; 2} \approx 9.21$, co daje $W \approx \langle 9.21; +\infty \rangle$.

Widzimy, że

$$\chi^2 = 8.9394 \notin W \approx \langle 9.21; +\infty \rangle,$$

więc nie mamy podstaw do odrzucenia H_0 i stwierdzamy, że rozkład liczby języków obcych, którymi posługują się studenci, jest rozkładem Poissona.

BARDZO WAŻNA UWAGA: W punkcie (b) nie możemy bezpośrednio skorzystać z funkcji `chisq.test(x=licznosci, p=prawdop0.b)` odczytując p-value, bo ono odpowiada prostej hipotezie H_0 : badana próba losowa pochodzi z rozkładu Poissona o średniej równej 1,15. Aby przetestować złożoną hipotezę H_0 : badana próba losowa pochodzi z rozkładu Poissona, możemy jedynie wykorzystać obliczoną w ten sposób wartość statystyki testowej `X-squared = 8.9394`, ale zbiór krytyczny musimy wyznaczyć sami.

2. Test Kołmogorowa-Smirnowa

Test Kołmogorowa-Smirnowa jest przeznaczony do sprawdzania zgodności rozkładu, z którego pochodzi próba losowa, z dowolnym rozkładem ciągłym.

```
> ks.test(x=dane, y="nazwa.dystrybuanty.rozkladu",
          liczby.opisujace.parametry.rozkladu)
```

Niestety `ks.test` obsługuje test Kołmogorowa-Smirnowa jedynie dla prostej H_0 . W przypadku złożonej H_0 (np. H_0 : rozkład badanej cechy jest wykładniczy z nieznanym parametrem λ) pozostaje nam samemu wyznaczyć przybliżoną wartość *p-value* metodą symulacji.

SPRAWDZANIE NORMALNOŚCI ROZKŁADU

I. Metody graficzne sprawdzania normalności

1. Wykres skrzynkowy. Jeśli nie jest on symetryczny lub widać na nim dużo obserwacji odstających (dla rozkładu normalnego $\mathcal{N}(\mu, \sigma^2)$ średnio 7 obserwacji na 1000 znajduje się poza przedziałem

$$(Q_1 - 1.5IQR, Q_3 + 1.5IQR),$$

to uznajemy, że rozkład cechy w próbie znacznie odbiega od rozkładu normalnego. W przeciwnym przypadku, dane mogą, ale nie muszą,

pochodzić z rozkładu normalnego i trzeba to sprawdzać innymi metodami.

```
> boxplot(dane)
```

Niestety wykresy skrzynowe są wiarygodne jedynie dla dużej liczności prób (jeśli mamy mało danych, to nie warto na nie patrzeć podczas badania normalności rozkładu).

2. Histogram częstości i jądrowy estymator gęstości.

```
> hist(dane,freq=F)
```

```
> lines(density(dane))
```

Pierwsza z powyższych komend rysuje w R histogram częstości (nie liczności); druga nanosi na poprzednio wykonany rysunek jądrowy estymator gęstości.

3. Wykres kwantylowy (wykres normalności) (Q-Q plot: quantile versus quantile plot).

Niech x_1, x_2, \dots, x_n oznacza realizację próby losowej. Po jej uporządkowaniu (od obserwacji najmniejszej do największej) otrzymujemy tzw. *statystyki porządkowe* z próby, oznaczane $x_{1:n} \leq x_{2:n} \leq \dots \leq x_{n:n}$. Wykres kwantylowy to zbiór punktów o współrzędnych $(u_{(i-0,5)/n}, x_{i:n})$, gdzie $i = 1, 2, \dots, n$ zaś $u_{(i-0,5)/n}$ to kwantyl standardowego rozkładu normalnego rzędu $(i - 0, 5)/n$.

Jeśli próba losowa pochodzi z rozkładu normalnego $\mathcal{N}(\mu, \sigma^2)$, to wykres kwantylowy jest zbiorem punktów leżących mniej-więcej na prostej $y = \sigma x + \mu$.

```
> qqnorm(dane)
```

```
> qqline(dane)
```

Pierwsza z powyższych komend rysuje w R wykres kwantylowy; druga nanosi na ten wykres linię przechodzącą przez kwartyle.

II. Testy normalności

Wyróżnia się dwa rodzaje testów normalności:

- **testy uniwersalne:**

H_0 : rozkład, z którego pochodzi badana próba losowa, jest normalny

H_1 : rozkład, z którego pochodzi badana próba losowa, nie jest normalny

- **testy kierunkowe:** badamy pewną ustaloną własność rozkładu normalnego, np. sprawdzamy czy rozkład, z którego pochodzi badana próba losowa, jest rozkładem symetrycznym (test skośności) albo czy jest rozkładem o kurtozie równej zero (test kurtozy).

Test skośności:

H_0 : rozkład, z którego pochodzi badana próba losowa, jest symetryczny

H_1 : rozkład, z którego pochodzi badana próba losowa, nie jest symetryczny (czyli jest skośny)

UWAGA: Ściśle rzecz biorąc, nie testujemy tu normalności rozkładu, lecz jedynie jego symetrię.

Test kurtozy:

H_0 : rozkład, z którego pochodzi badana próba losowa, ma kurtozę równą zero

H_1 : rozkład, z którego pochodzi badana próba losowa, ma kurtozę różną od zera

UWAGA: Ściśle rzecz biorąc, nie testujemy tu normalności rozkładu, lecz jedynie sprawdzamy czy jego kurtoza wynosi zero. Przyjeliśmy tu definicję kurtozy zmiennej losowej, według której kurtoza rozkładu normalnego wynosi zero (a nie 3).

```
> install.packages("fBasics")
> library(fBasics)
> dagoTest(dane)
```

Testy uniwerslane

1. Test Shapiro-Wilka

Zaproponowany w 1965 r. jest to dziś uznawany za najlepszy test uniwersalny normalności rozkładu. Konstrukcja tego testu opiera się na wykresie kwantylowym. Dokładniej, wyznacza się linię, która jest możliwie najlepiej dopasowana do punktów tego wykresu (mówiąc precyzyjniej, wyznacza się tzw. *prostą regresji*) i następnie bada się stopień dopasowania tych punktów do owej prostej.

```
> shapiro.test(dane)
```

2. Test Andersona-Darlinga

```
> install.packages("nortest")
> library(nortest)
> ad.test(dane)
```

3. Test Craméra-von Misesa

```
> install.packages("nortest")  
> library(nortest)  
> cvm.test(dane)
```

W porównaniu z testem Craméra-von Misesa, test Andersona-Darlinga zwraca większą uwagę na ogony rozkładu.

4. Test Lilliefors'a

```
> install.packages("nortest")  
> library(nortest)  
> lillie.test(dane)
```

Test Lilliefors'a sprawuje się średnio gorzej niż test Andersona-Darlinga i test Craméra-von Misesa.