

Badanie zgodności z określonym rozkładem

H_0 : Cecha X ma rozkład F

F jest dowolnym rozkładem prawdopodobieństwa

Test chi–kwadrat zgodności

F jest rozkładem ciągłym

Test Kołmogorowa

F jest rozkładem normalnym

Test Shapiro–Wilka

Test Chi–kwadrat zgodności (poziom istotności α)

Klasa	Liczebność
1	n_1
2	n_2
\vdots	\vdots
k	n_k

Statystyka testowa

$$\chi_{\text{emp}}^2 = \sum_{i=1}^k \frac{(n_i - n_i^t)^2}{n_i^t}$$

$$n_i^t = Np_i^t, \quad N = \sum_{i=1}^k n_i,$$

$$p_i^t = P_F\{X \text{ przyjęła wartość z klasy } i\}$$

Wartość krytyczna $\chi^2(\alpha; k - u - 1)$ (u jest liczbą nieznanymi parametrów hipotetycznego rozkładu F)

Wniosek. Jeżeli $\chi_{\text{emp}}^2 > \chi^2(\alpha; k - u - 1)$, to hipotezę H_0 odrzucamy

Przykład. Pracodawca przypuszcza, że liczba pracowników nieobecnych w różne dni tygodnia nie jest taka sama. W tym celu w ciągu pewnego okresu czasu zebrał następujące dane

Dzień	n_i
Poniedziałek	200
Wtorek	160
Środa	140
Czwartek	140
Piątek	100

Populacja:

Cecha X :

dzień nieobecności pracownika

Założenie:

cecha przyjmuje wartości będące nazwami dni tygodnia (cecha jakościowa)

Formalizacja:

Liczbę pracowników nieobecnych w kolejne dni tygodnia można przedstawić jako odesetek załogi. Odesetki te można interpretować jako prawdopodobieństwo nieobecności pracownika w danym dniu tygodnia. Jeżeli ilość pracowników nieobecnych w kolejne dni tygodnia jest „mniej więcej” taka sama, to można ten fakt sformalizować jako identyczne prawdopodobieństwo nieobecności pracownika w poszczególne dni tygodnia. Tak więc, weryfikowana będzie hipoteza

$$H_0 : X \text{ ma rozkład } \frac{\text{Pon}}{1/5} \quad \frac{\text{Wtk}}{1/5} \quad \frac{\text{Śro}}{1/5} \quad \frac{\text{Czw}}{1/5} \quad \frac{\text{Ptk}}{1/5}$$

Technika statystyczna:

test chi–kwadrat zgodności
poziom istotności $\alpha = 0.05$

Obliczenia

Dzień	n_i	p_i^t	n_i^t	$(n_i - n_i^t)^2 / n_i^t$
Poniedziałek	200	1/5	148	$\frac{(200-148)^2}{148} = 18.270$
Wtorek	160	1/5	148	$\frac{(160-148)^2}{148} = 0.973$
Środa	140	1/5	148	$\frac{(140-148)^2}{148} = 0.432$
Czwartek	140	1/5	148	$\frac{(140-148)^2}{148} = 0.432$
Piątek	100	1/5	148	$\frac{(100-148)^2}{148} = 15.676$
	740			$\chi_{\text{emp}}^2 = 35.676$

Wartość krytyczna

$$\chi^2(\alpha; k - u - 1) = \chi^2(0.05; 5 - 0 - 1) = 9.4877$$

Odpowiedź: hipotezę odrzucamy

Wniosek:

Odrzucamy hipotezę o równomiernym rozkładzie nieobecności w tygodniu. Zatem przypuszczenie pracodawcy można uznać za uzasadnione

Przykład. Na pewnej uczelni badano strukturę miesięcznych dochodów (na głowę) w rodzinach studentów. W tym celu wylosowano grupę 192 studentów i zanotowano miesięczne dochody w ich rodzinach. Uzyskano następujące wyniki (w setkach złotych):

x_i	x_{i+1}	n_i
poniżej 6		6
6	7	11
7	8	18
8	9	27
9	10	32
10	11	35
11	12	24
12	13	20
13	14	13
powyżej 14		6

Czy można, że rozkładów dochodów w rodzinach studenckich jest normalny?

Populacja:

studenci pewnej uczelni

Cecha X :

miesięczne dochody na głowę w rodzinach studentów

Założenie:

cecha ciągła

Formalizacja:

H_0 : Cecha X ma rozkład normalny $N(\mu, \sigma^2)$

Technika statystyczna:

test chi–kwadrat zgodności
poziom istotności $\alpha = 0.05$

Obliczenia

Szereg ma $k = 10$ klas

Do całkowitego określenia hipotetycznego rozkładu
brakuje dwóch parametrów, czyli $u = 2$

Wartość krytyczna

$$\chi^2(\alpha; k - u - 1) = \chi^2(0.05; 10 - 2 - 1) = 14.0671$$

Wyznaczenie wartości statystyki χ_{emp}^2

Wyznaczenie prawdopodobieństw teoretycznych

$$p_i^t = P\{x_i < X < x_{i+1}\} = F\left(\frac{x_{i+1} - \mu}{\sigma}\right) - F\left(\frac{x_i - \mu}{\sigma}\right)$$

Z próby wyznaczamy $\bar{x} = 10.09$, $s^2 = 4.81$

$$p_i^t = F\left(\frac{x_{i+1} - \bar{x}}{s}\right) - F\left(\frac{x_i - \bar{x}}{s}\right) = F(z_{i+1}) - F(z_i)$$

x_i	x_{i+1}	z_i	z_{i+1}	$F(z_i)$	$F(z_{i+1})$	p_i^t
poniżej 6		$-\infty$	-1.82	0.0000	0.0344	0.0344
6	7	-1.82	-1.36	0.0344	0.0869	0.0525
7	8	-1.36	-0.91	0.0869	0.1814	0.0943
8	9	-0.91	-0.45	0.1814	0.3264	0.1450
9	10	-0.45	0.00	0.3264	0.5000	0.1736
10	11	0.00	0.45	0.5000	0.6736	0.1736
11	12	0.45	0.91	0.6736	0.8186	0.1450
12	13	0.91	1.36	0.8186	0.9131	0.0943
13	14	1.36	1.82	0.9131	0.9656	0.0525
powyżej 14		1.82	∞	0.9656	1.0000	0.0344

Wyznaczenie wartości statystyki testowej

x_i	x_{i+1}	n_i	p_i^t	n_i^t	$(n_i - n_i^t)^2 / n_i^t$
poniżej 6		6	0.0344	6.53	0.36
6	7	11	0.0525	10.18	0.07
7	8	18	0.0943	18.05	0.00
8	9	27	0.1450	27.84	0.02
9	10	32	0.1736	33.41	0.06
10	11	35	0.1736	33.41	0.08
11	12	24	0.1450	27.84	0.53
12	13	20	0.0943	18.05	0.21
13	14	13	0.0525	10.18	0.78
powyżej 14		6	0.0344	6.53	0.03
		192			2.14

Odpowiedź. Nie odrzucamy hipotezy

Wniosek. Możemy uznać, że miesięczne dochody na głowę w rodzinach studentów mają rozkład normalny $N(10.09, 4.81)$

H_0 : Cecha X ma rozkład ciągły F

Test Kołmogorowa zgodności (poziom istotności α)

Próba X_1, \dots, X_n

Próbkę X_1, \dots, X_n porządkujemy niemalejąco

$$X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n)}$$

Tak uporządkowane obserwacje nazywamy statystykami pozycyjnymi

Statystyka testowa

$$D_n = \max_{1 \leq i \leq n} \left\{ \max \left\{ \left| F(X_{(i)}) - \frac{i-1}{n} \right|, \left| \frac{i}{n} - F(X_{(i)}) \right| \right\} \right\}$$

Wartość krytyczna testu Kołmogorowa $D(\alpha; n)$

Jeżeli $D_n > D(\alpha; n)$, to hipotezę H_0 odrzucamy

Przykład. Obserwujemy cechę X

H_0 : X ma rozkład jednostajny na przedziale $(0, 1)$

Test Kołmogorowa zgodności ($\alpha = 0.05$)

Dystrybuanta rozkładu jednostajnego

$$F(x) = \begin{cases} 0, & \text{jeżeli } x \leq 0, \\ x, & \text{jeżeli } 0 < x \leq 1, \\ 1, & \text{jeżeli } x \geq 1. \end{cases}$$

Próba:

0.62799, 0.06409, 0.70428, 0.98945, 0.59274, 0.49986,
0.31380, 0.02692, 0.63097, 0.95539, 0.84967, 0.86341,
0.80649, 0.73106, 0.99956, 0.86344, 0.00744, 0.28036,
0.66399, 0.21092.

Wartość statystyki testowej $D_{20} = 0.24274$

Wartość krytyczna $D(0.05; 20) = 0.29408$

Ponieważ $D_{20} < D_{0.05,20}$, więc nie odrzucamy weryfikowanej hipotezy

Wyznaczenie wartości statystyki testowej

i	$F(X_{(i)})$	i/n	$(i-1)/n$	$ F(X_{(i)}) - i/n $	$ F(X_{(i)}) - (i-1)/n $
1	0.00744	0.05	0.00	0.04256	0.00744
2	0.02692	0.10	0.05	0.07308	0.02308
3	0.06409	0.15	0.10	0.08591	0.03591
4	0.21092	0.20	0.15	0.01092	0.06092
5	0.28036	0.25	0.20	0.03036	0.08036
6	0.31380	0.30	0.25	0.01380	0.06380
7	0.49986	0.35	0.30	0.14986	0.19986
8	0.59274	0.40	0.35	0.19274	0.24274
9	0.62799	0.45	0.40	0.17799	0.22799
10	0.63097	0.50	0.45	0.13097	0.18097
11	0.66399	0.55	0.50	0.11399	0.16399
12	0.70428	0.60	0.55	0.10428	0.15428
13	0.73106	0.65	0.60	0.08106	0.13106
14	0.80649	0.70	0.65	0.10649	0.15649
15	0.84967	0.75	0.70	0.09967	0.14967
16	0.86341	0.80	0.75	0.06341	0.11341
17	0.86344	0.85	0.80	0.01344	0.06344
18	0.95539	0.90	0.85	0.05539	0.10539
19	0.98945	0.95	0.90	0.03945	0.08945
20	0.99956	1.00	0.95	0.00044	0.04956

H_0 : Cecha X ma rozkład normalny

Test Shapiro–Wilka (poziom istotności α)

Próbkę X_1, \dots, X_n porządkujemy niemalejąco

$$X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n)}$$

Tak uporządkowane obserwacje nazywamy statystykami pozycyjnymi

Statystyka testowa

$$W = \frac{\left(\sum_{i=1}^{[n/2]} a_{i:n} (X_{(n-i+1)} - X_{(i)}) \right)^2}{\text{var}X}$$

$a_{i:n}$ są tablicowanymi współczynnikami

$$[n/2] = \begin{cases} n/2, & \text{dla } n \text{ parzystych} \\ (n-1)/2, & \text{dla } n \text{ nieparzystych} \end{cases}$$

Wartość krytyczna testu Shapiro–Wilka $W_n(\alpha)$

Jeżeli $W \leq W_n(\alpha)$, to hipotezę H_0 odrzucamy.

Przykład. Z cechy X pobrano próbę ($n = 19$): 12.4, 14.2, 14.9, 15.6, 16.1, 17.3, 17.9, 18.2, 18.6, 19.3, 19.7, 20.4, 21.9, 22.8, 23.7, 25.2, 25.9, 27.4.

H_0 : Cecha X ma rozkład normalny

Test Shapiro–Wilka ($\alpha = 0.05$)

$$\bar{x} = 19.3842, \text{ var}x = 730.57$$

Licznik statystyki W

i	$x_{(19-i+1)} - x_{(i)}$	$a_{i:19}$	$a_{i:19}(x_{(19-i+1)} - x_{(i)})$
1	27.4–12.4=15.0	0.4808	7.21200
2	25.9–14.2=11.7	0.3232	3.78144
3	25.2–14.9=10.3	0.2561	2.63783
4	23.7–15.6= 8.1	0.2059	1.66779
5	22.8–16.1= 6.7	0.1641	1.09947
6	21.9–16.8= 5.1	0.1271	0.64821
7	20.4–17.3= 3.1	0.0932	0.28892
8	19.7–17.9= 1.8	0.0612	0.11016
9	19.3–18.2= 1.1	0.0303	0.03333
			17.47915

Statystyka testowa $W = \frac{305.52}{730.57} = 0.418$

Wartość krytyczna $W_{19}(0.05) = 0.901$

Ponieważ $W < W_{19}(0.05)$, więc hipotezę odrzucamy