

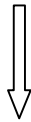
Statystyczne testy nieparametryczne

Testami nieparametrycznymi nazywamy testy służące do weryfikacji hipotez nieparametrycznych, tj hipotez nie dotyczących wartości nieznanymi parametrów populacji (choć czasem pojęcie to oznacza hipotezy nie zakładające rozkładu Normalnego dla populacji). Ze względu na różnorodność hipotez nieparametrycznych, klasę testów nieparametrycznych można podzielić na następujące podklasy:

- testy zgodności (z pewnym rozkładem teoretycznym), w tym testy normalności,
- testy jednorodności, czyli zgodności dwóch (lub więcej) rozkładów,
- testy niezależności,
- inne testy, w tym np. testy weryfikujące hipotezę, że próba ma charakter losowy.

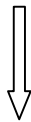
Badanie niezależności rozkładu dwu cech

- Cecha (X, Y) ma dwuwymiarowy, nieznaną rozkład



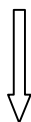
Test Chi-Kwadrat niezależności

- Cecha (X, Y) ma dwuwymiarowy rozkład ciągły



Współczynnik korelacji rangowej Spearmana
Współczynnik korelacji rangowej Kendalla

- Czy kolejność obserwacji w próbie jest losowa?



test serii

ZALEŻNOŚĆ MIĘDZY CECHAMI JAKOŚCIOWYMI /SKATEGORYZOWANYMI/

X, Y – cechy obserwowane

Próba: $(X_1, Y_1), \dots, (X_k, Y_m)$

klasy cechy X	klasy cechy Y			
	1	2	...	m
1	n_{11}	n_{12}	...	n_{1m}
2	n_{21}	n_{22}	...	n_{2m}
⋮	⋮	⋮	⋮	⋮
k	n_{k1}	n_{k2}	...	n_{km}

H_0 : Cechy X i Y są niezależne

H_1 : Cechy X i Y są zależne

Test Chi-kwadrat (c^2) niezależności

$$C_{emp}^2 = \sum_{i=1}^k \sum_{j=1}^m \frac{(n_{ij} - n_{ij}^t)^2}{n_{ij}^t}$$

n_{ij} – liczba obserwacji realizujących i-tą wartość cechy X i j-tą wartość Y

n_{ij}^t – teoretyczna liczba obserwacji realizujących i-tą wartość cechy X i j-tą wartość Y (wg. rozkładów brzegowych dla każdej z tych dwu cech)

$$n_{ij}^t = \frac{n_{i\bullet} \cdot n_{\bullet j}}{N}, \quad N = \sum_{i=1}^k \sum_{j=1}^m n_{ij},$$

$$n_{i\bullet} = \sum_{j=1}^m n_{ij}, \quad n_{\bullet j} = \sum_{i=1}^k n_{ij}$$

Jeśli $C_{emp}^2 > C_{kryt}^2$, to hipotezę H_0 odrzucamy.

$$C_{kryt}^2 = C_{a,v}^2, \text{ gdzie } v = (k-1) \cdot (m-1)$$

Przykład:

Badano dwie właściwości wędliny: związanie (słabo związana, związana, dobrze związana) oraz smakowitość (dostateczna, dobra, bardzo dobra). Analizę przeprowadzono dla 60 batonów wędliny. Wyniki były następujące:

X- smakowitość	Y₁-słabo związana	Y₂- związana	Y₃- dobrze związana	n_i
X₁- dostateczna	9	5	3	17
X₂-dobra	4	12	6	22
X₃-b. dobra	1	6	14	21
n._j	14	23	23	60

Hipoteza badawcza i statystyczna brzmi:

H_0 : Smakowitość wędliny (X) nie zależy od stopnia związania (Y) tj. cechy te są niezależne

Konstruujemy funkcję testową opartą na rozkładzie chi-kwadrat.

Dane (**liczebności**) **teoretyczne** N_{ij}^t (jeśli cechy są niezależne czyli $P(A \cap B) = P(A) \cdot P(B)$, to rozkłady brzegowe $P(A)$ i $P(B)$ wyznaczają prawdopodobieństwo podklas $P(A \cap B)$):

X- smakowitość	Y₁-słabo związana	Y₂- związana	Y₃- dobrze związana	n_i
X₁- dostateczna	4	6,5	6,5	17
X₂-dobra	5	8,5	8,5	22
X₃-b. dobra	5	8	8	21
n_j	14	23	23	60

$$n_{11}^t = \frac{17 \cdot 14}{60} = 3.97, \quad n_{12}^t = \frac{17 \cdot 23}{60} = 6.52, \dots$$

$$n_{33}^t = \frac{21 \cdot 23}{60} = 8.05$$

$$C_{emp}^2 = \frac{(9-3.97)^2}{3.97} + \frac{(5-6.52)^2}{6.52} + \dots + \frac{(14-8.05)^2}{8.05}$$

$$= 19.2$$

$$C_{kryt}^2 = C_{a,v}^2 = C_{0.05,4}^2 = 9.49$$

Ponieważ $C_{emp}^2 > C_{kryt}^2$, to hipotezę H_0 odrzucamy.

Wyniki pozwalają stwierdzić, że smakowitość badanej wędliny zależy od związania jej składników. Dla podniesienia walorów sensorycznych tej wędliny należy tak prowadzić proces technologiczny, aby uzyskać możliwie największe jej związanie

Można wyznaczyć także współczynnik kontyngencji P , który przyjmuje wartość zero, gdy występuje całkowita niezależność cech.

$$P = \sqrt{\frac{c^2}{N + c^2}}$$

Dla naszego przykładu

$$P = \sqrt{\frac{c^2}{N + c^2}} = \sqrt{\frac{19.2}{60 + 19.2}} = 0.492$$

co świadczy o dużej sile związku między rozważanymi cechami.

TESTY ZGODNOŚCI

Hipotezy tego typu dotyczą zgodności rozkładu empirycznego z rozkładem określonym przez hipotezę lub zgodności (jednorodności) rozkładów pewnej cechy w kilku populacjach bez określania, o jaki rozkład chodzi. Z tego też powodu testy służące do weryfikacji takich hipotez nazywamy **testami zgodności (jednorodności)**.

Do najczęściej stosowanych testów zgodności należą:

- χ^2 (chi-kwadrat) Pearsona
- λ (lambda) Kołmogorowa-Smirnowa
- w Shapiro-Wilka

Niech hipotezą zerową będzie przypuszczenie, że cecha X ma w populacji rozkład określony dystrybuantą $F_0(x)$:

$$H_0 : F(x) = F_0(x) \quad \text{wobec} \quad H_1 : F(x) \neq F_0(x)$$

Statystyka

$$c^2 = \sum_{j=1}^k \frac{(n_j - n_j^t)^2}{n_j^t}$$

przy prawdziwości H_0 ma asymptotyczny rozkład c^2 z liczbą stopni swobody $\nu = k - u - 1$.

Wielkość $n_j^t = np_j$ jest teoretyczną (to znaczy, obliczoną przy założeniu prawdziwości testowanej hipotezy H_0) liczebnością w j -tym przedziale, k jest liczbą przedziałów klasowych, a u liczbą parametrów populacyjnych, szacowanych z próby.

Wartość empiryczną statystyki

$$C_{emp}^2 = \sum_j \frac{(n_j - n_j^t)^2}{n_j^t}$$

porównujemy z wartością krytyczną odczytaną z tablic statystycznych

$$C_{a, v=k-u-1}^2$$

wnioskując analogicznie jak w pozostałych hipotezach.

Elementem kluczowym przy wykorzystaniu statystyki Chi-kwadrat jest wielkość

$$p_j^t = P(x \in (x_{1j}; x_{2j}))$$

która jest teoretycznym (to znaczy, obserwowanym przy założeniu prawdziwości testowanej hipotezy H_0) prawdopodobieństwem wystąpienia obserwacji w j -tym przedziale.

Przykład: Pracodawca przypuszcza, że liczba pracowników nieobecnych w różne dni tygodnia nie jest taka sama.

W celu sprawdzenia swojego przypuszczenia obserwował, przez pewien okres, liczby pracowników nieobecnych w kolejnych dniach tygodnia. Wyniki obserwacji zawiera tabela:

dzień tygodnia	liczba nieobecnych
poniedziałek	200
wtorek	160
środa	140
czwartek	140
piątek	100

Badaną cechą X jest dzień, w którym pracownik był nieobecny w pracy. Jest to cecha jakościowa o wartościach: poniedziałek, wtorek, ... , piątek.

Hipoteza badawcza, że absencja pracownika jest zależna od dnia tygodnia pracy, może być zapisana ‘przez negację’, to znaczy sugerujemy **brak preferencji w opuszczaniu dni.**

Zapis statystyczny tego przypuszczenia pracodawcy ma postać hipotezy:

H : cecha X ma rozkład:

Pon.	Wtk.	Śro.	Czw.	Ptk.
1/5	1/5	1/5	1/5	1/5

Do weryfikacji badanej hipotezy stosujemy test chi–kwadrat zgodności, przyjmując $\alpha = 0.05$.

Pomocnicze obliczenia funkcji testowej zawiera tabela:

	n_i	n_i teoret.
Pon	200	148
Wtk	160	148
Śro	140	148
Czw	140	148
Ptk	100	148
suma	740	740

Wartość statystyki jest wyznaczona według formuły:

$$C_{emp}^2 = \sum_j \frac{(n_j - n_j^t)^2}{n_j^t} =$$
$$\frac{(200 - 148)^2}{148} + \dots + \frac{(100 - 148)^2}{148} = 35,68$$

Ponieważ wartość krytyczna

$$C_{a,v=k-u-1}^2 = C_{0.05,5-0-1}^2 = C_{0.05,4}^2 = 9.49$$

zachodzi relacja $C_{emp}^2 > C_{0.05,4}^2$, czyli hipotezę o zgodności z określonym rozkładem odrzucamy.

Oznacza to, że przypuszczenie pracodawcy o nierównomiernym rozkładzie absencji w zakładzie pracy można uznać za uzasadnione.

Test c2 zgodności (jednorodności) kilku rozkładów

Obserwujemy tę samą cechę w kilku populacjach. Interesuje nas odpowiedź na pytanie, czy rozkłady te są takie same (**co pociąga za sobą równość wszystkich parametrów**).

Jeżeli dystrybuantę danej cechy w i -tej populacji oznaczymy jako F_i , to hipoteza zerowa ma postać:

$$H_0 : F_1 = F_2 = \dots = F_k$$

Zastosowanie testu χ^2 wymaga zestawienia próby w postaci tabeli dwukierunkowej.

W jednym kierunku umieszczamy poziomy danej cechy, w drugim populacje.

numer populacji	klasy cechy X			
	X_1	X_2	\dots	X_r
1	n_{11}	n_{12}	\dots	n_{1r}
2	n_{21}	n_{22}	\dots	n_{2r}
\vdots	\vdots	\vdots	\vdots	\vdots
k	n_{k1}	n_{k2}	\dots	n_{kr}

Statystyka testowa ma postać:

$$C_{emp}^2 = \sum_{i=1}^k \sum_{j=1}^r \frac{(n_{ij} - n_{ij}^t)^2}{n_{ij}^t}$$

gdzie $n_{ij}^t = \frac{n_{i\bullet} \cdot n_{\bullet j}}{n}$

n_{ij} – oznacza liczbę obserwacji reprezentujących i -tą populację i j -tą klasę cechy X .

Z indeksem górnym t , jest to odpowiednia **liczebność teoretyczna**.

Przy prawdziwości H_0 statystyka ta ma rozkład χ^2 Pearsona z liczbą stopni swobody $\nu = (k-1)(r-1)$.

Wnioskowanie przebiega analogicznie jak przy innych hipotezach.

przykład:

We wszystkich 10 sklepach pewnej sieci sklepów jest takie samo zapotrzebowanie na mleko o różnej zawartości tłuszczu.

sklep	0,5%	2,0%	3,2%
1	34	36	28
2	48	42	46
3	15	18	10
4	61	45	51
5	37	29	46
6	18	16	25
7	39	28	35
8	42	18	31
9	41	38	29
10	19	26	16

Wartość statystyki C_{emp}^2 wynosi 23,8. Powoduje to, iż krytyczny poziom istotności (P-value) wynosi 16%. Nie ma podstaw, aby na poziomie istotności 5% odrzucić hipotezę o takim samym rozkładzie preferencji odnośnie zawartości tłuszczu we wszystkich 10 sklepach