

Rozdział 3

Wnioskowanie statystyczne

3.1 Populacja, próbka, parametry rozkładu i estymatory

W *populacji*, cecha X ma rozkład F . Nieznana liczba θ jest *parametrem* tego rozkładu (pewną charakterystyką liczbową). *Próbka*: X_1, \dots, X_n – wartości cechy X dla n elementów wylosowanych z populacji. Jeśli losujemy *ze zwracaniem* lub jeśli populacja jest duża, to X_1, \dots, X_n są *niezależnymi* zmiennymi losowymi o rozkładzie F . *Estymator* parametru θ jest to wielkość obliczona na podstawie próbki,

$$\hat{\theta} = \hat{\theta}(X_1, \dots, X_n),$$

która jest dostępnym oszacowaniem (przybliżeniem) nieznanej liczby θ .

- Cecha X jakościowa, o wartościach 0, 1 (tak, nie).

Parametr p – frakcja (procent) *populacji*, w której cecha ma wartość 1. Inaczej: p – frakcja elementów wyróżnionych w populacji, czyli „wskaźnik struktury”:

$$p = \mathbb{P}(X = 1).$$

Estymator:

$$\hat{p} = \frac{K}{n},$$

gdzie K oznacza liczbę elementów wyróżnionych w *próbce* (liczbę jedynek w ciągu X_1, \dots, X_n).

- **Cecha X jakościowa, o k wartościach w_1, \dots, w_k .**

Parametry – p_1, \dots, p_k , gdzie p_i – frakcja *populacji*, w której cecha X ma wartość w_i . Zmienna losowa X ma rozkład prawdopodobieństwa dany tabelką:

wynik	w_1	\dots	w_i	\dots	w_k
prawdopodobieństwo	p_1	\dots	p_i	\dots	p_k

gdzie $p_i = \mathbb{P}(X = w_i)$. Oczywiście, $p_1 + \dots + p_k = 1$.

Dla *próbki* n -elementowej, budujemy „tabelkę powtórzeń”:

wartość cechy X	w_1	\dots	w_i	\dots	w_k
liczba elementów próbki	N_1	\dots	N_i	\dots	N_k

gdzie

N_i = liczba elementów *próbki*, dla których cecha X ma wartość w_i .

Oczywiście, $N_1 + \dots + N_k = n$.

- **Cecha X ilościowa (o wartościach liczbowych).**

Parametry (np.): μ – wartość średnia cechy X w *populacji*; σ^2 – wariancja cechy X w *populacji*.

$$\mu = \mathbb{E}(X); \quad \sigma^2 = \text{Var}(X).$$

Estymatory:

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i; \quad S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2.$$

3.2 Przedziały ufności

3.2.1 DEFINICJA. Niech θ będzie nieznanym parametrem, X_1, \dots, X_n – obserwowaną próbką. Mówimy, że $[\underline{\theta}, \bar{\theta}]$ jest **przedziałem ufności** dla θ na poziomie $1 - \alpha$, jeśli $\underline{\theta} = \underline{\theta}(X_1, \dots, X_n)$ i $\bar{\theta} = \bar{\theta}(X_1, \dots, X_n)$ oraz

$$\mathbb{P} \left(\underline{\theta} \leq \theta \leq \bar{\theta} \right) \geq 1 - \alpha.$$

Cecha X ilościowa, ciągła, rozkład normalny.

Zakładamy, że cecha X ma w populacji rozkład $N(\mu, \sigma^2)$.

- ★ **Przedział ufności dla średniej μ , znana wariancja σ^2 .**

$$\left[\bar{X} - \frac{\sigma z}{\sqrt{n}}, \bar{X} + \frac{\sigma z}{\sqrt{n}} \right],$$

inaczej: $\mu = \bar{X} \pm \sigma z / \sqrt{n}$, gdzie $z = z_{1-\alpha/2}$ – kwantyl rozkładu $N(0, 1)$ (standardowego normalnego).

- ★ **Przedział ufności dla średniej μ , nieznana wariancja σ^2 .**

$$\left[\bar{X} - \frac{St}{\sqrt{n}}, \bar{X} + \frac{St}{\sqrt{n}} \right],$$

inaczej: $\mu = \bar{X} \pm St / \sqrt{n}$, gdzie $t = t_{1-\alpha/2}(n-1)$ – kwantyl rozkładu $t(n-1)$ (t-Studenta z $n-1$ stopniami swobody).

- ★ **Przedział ufności dla wariancji σ^2 .**

$$\left[\frac{(n-1)S^2}{c_2}, \frac{(n-1)S^2}{c_1} \right],$$

gdzie $c_1 = \chi_{\alpha/2}^2(n-1)$, i $c_2 = \chi_{1-\alpha/2}^2(n-1)$ są kwantylami rzędu, odpowiednio, $\alpha/2$ i $1 - \alpha/2$ rozkładu chi-kwadrat z $n-1$ stopniami swobody.

Cecha X jakościowa, o wartościach 0,1 (tak, nie).

- ★ **Przedział ufności dla wskaźnika struktury p .**

$$\left[\hat{p} - \frac{\sqrt{\hat{p}(1-\hat{p})}z}{\sqrt{n}}, \hat{p} + \frac{\sqrt{\hat{p}(1-\hat{p})}z}{\sqrt{n}} \right],$$

inaczej $p = \hat{p} \pm \sqrt{\hat{p}(1-\hat{p})}z / \sqrt{n}$, gdzie $z = z_{1-\alpha/2}$ – kwantyl rozkładu $N(0, 1)$ (standardowego normalnego).

Uwaga. To jest rozwiązanie przybliżone, które można stosować gdy n jest duże, zaś \hat{p} niezbyt bliskie 0 i 1, powiedzmy: $n\hat{p}(1-\hat{p}) \geq 9$.

3.3 Testy istotności

Hipoteza statystyczna – przypuszczenie na temat rozkładu prawdopodobieństwa, opisującego *populację*.

$$\begin{aligned}H_0 &: \text{hipoteza zerowa;} \\ H_1 &: \text{hipoteza alternatywna.}\end{aligned}$$

Test – procedura, która na podstawie danych (próbki X_1, \dots, X_n) prowadzi do decyzji

$$\begin{aligned}\text{albo} &\longrightarrow \text{odrzuć } H_0 \text{ (na korzyść } H_1); \\ \text{albo} &\longrightarrow \text{nie odrzucać } H_0.\end{aligned}$$

3.3.1 DEFINICJA. Test jest **na poziomie istotności α** , jeśli

$$\mathbb{P}_{H_0}(\text{odrzućmy } H_0) \leq \alpha.$$

\mathbb{P}_{H_0} – prawdopodobieństwo obliczone przy założeniu, że H_0 jest *prawdziwa*.

Najczęściej test ma postać:

$$\text{odrzućmy } H_0, \text{ jeśli } T > c,$$

gdzie $T = T(X_1, \dots, X_n)$ jest „statystyką testową” (obliczoną na podstawie próbki), zaś c nazywa się *poziomem krytycznym* testu (zazwyczaj odczytanym z odpowiednich tablic). Test jest na poziomie istotności α , jeśli $\mathbb{P}_{H_0}(T > c) \leq \alpha$.

Uwaga: Niekiedy odrzućmy H_0 , jeśli $T < c$.

p-value.

Przypuśćmy, że obliczona na podstawie danych $X_1 = x_1, \dots, X_n = x_n$ wartość statystyki testowej jest równa liczbie $t = T(x_1, \dots, x_n)$. Z tablic rozkładu zmiennej losowej T można odczytać wielkość

$$p = \mathbb{P}_{H_0}(T > t),$$

którą nazywamy *p-value*. Małe *p-value* świadczy przeciwko hipotezie zerowej:

$$\text{odrzućmy } H_0, \text{ jeśli } p < \alpha,$$

gdzie α jest założonym poziomem istotności.

Typowe zagadnienia, w których używa się testów istotności:

- **Porównanie z „normą”.** Rozważamy cechę X , która ma w populacji rozkład F . Mamy próbkę X_1, \dots, X_n . Testujemy hipotezę, że X ma „spodziewany” rozkład F_0 :

$$H_0 : F = F_0$$

(przeciw alternatywie $H_1 : F \neq F_0$). Rozważa się też nieco inne hipotezy.

- **Porównanie 2 populacji.** Badamy dwie populacje. Cecha X ma w pierwszej populacji rozkład F_1 , zaś w drugiej – rozkład F_2 . Mamy dwie próbki: X_{11}, \dots, X_{1n_1} – z pierwszej populacji, X_{21}, \dots, X_{2n_2} – z drugiej populacji. Badamy, czy rozkład cechy X jest w obu populacjach *jednakowy*. Testujemy hipotezę

$$H_0 : F_1 = F_2$$

(przeciw alternatywie $H_1 : F_1 \neq F_2$).

- **Porównanie k populacji.** Badamy k populacji. Cecha X ma w j -tej populacji rozkład F_j ($j = 1, \dots, k$). Mamy k próbek: X_{j1}, \dots, X_{jn_j} – jest próbka z j -tej populacji ($j = 1, \dots, k$). Badamy, czy rozkład cechy X jest we wszystkich populacjach *jednakowy*. Testujemy hipotezę

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

(przeciw alternatywie H_1 : nie wszystkie rozkłady są jednakowe).

Cecha X ilościowa, ciągła, rozkład normalny.

Zakładamy, że cecha X ma w populacji rozkład $N(\mu, \sigma^2)$.

Porównanie z „normą”.

- ★ **Test** $H_0 : \mu \leq \mu_0$ **przeciwko** $H_1 : \mu > \mu_0$, gdzie μ_0 jest ustaloną liczbą. Na poziomie istotności α , odrzucamy H_0 , gdy

$$\sqrt{n} \frac{\bar{X} - \mu_0}{S} > t, \quad t = t_{1-\alpha}(n-1).$$

Inaczej: Odrzucamy $H_0 : \mu \leq \mu_0$, jeśli $\bar{X} > \mu_0 + St/\sqrt{n}$. Czasami mówi się wtedy: „średnia \bar{X} jest *istotnie większa* od μ_0 ”.

- ★ **Test** $H_0 : \mu = \mu_0$ **przeciwko** $H_1 : \mu \neq \mu_0$, gdzie μ_0 jest ustaloną liczbą. Na poziomie istotności α , odrzucamy H_0 , gdy

$$\sqrt{n} \frac{|\bar{X} - \mu_0|}{S} > t, \quad t = t_{1-\alpha/2}(n-1).$$

Inaczej: Odrzucamy $H_0 : \mu = \mu_0$, jeśli $|\bar{X} - \mu_0| > St/\sqrt{n}$. Czasami mówi się wtedy: „średnia \bar{X} jest *istotnie różna* od μ_0 ”.

- ★ **Test** $H_0 : \sigma \leq \sigma_0$ **przeciwko** $H_1 : \sigma > \sigma_0$, gdzie σ_0 jest ustaloną liczbą. Odrzucamy H_0 , gdy

$$\frac{n-1}{\sigma_0^2} S^2 > c, \quad c = \chi_{1-\alpha}^2(n-1).$$

- ★ **Test** $H_0 : \sigma = \sigma_0$ **przeciwko** $H_1 : \sigma \neq \sigma_0$, gdzie σ_0 jest ustaloną liczbą. Odrzucamy H_0 , gdy

$$\frac{n-1}{\sigma_0^2} S^2 > c_2 \text{ lub } \frac{n-1}{\sigma_0^2} S^2 < c_1, \quad c_1 = \chi_{\alpha/2}^2(n-1), c_2 = \chi_{1-\alpha/2}^2(n-1).$$

Porównanie 2 populacji.

Mamy dwie próbki, wylosowane (niezależnie) z dwóch populacji:

X_{11}, \dots, X_{1n_1} – z rozkładu $N(\mu_1, \sigma^2)$,

X_{21}, \dots, X_{2n_2} – z rozkładu $N(\mu_2, \sigma^2)$.

Zakładamy równość wariancji w obu populacjach. Znaczenie symboli \bar{X}_1 , \bar{X}_2 , S_1^2 i S_2^2 jest oczywiste.

- ★ **Test** $H_0 : \mu_1 \leq \mu_2$ **przeciwko** $H_1 : \mu_1 > \mu_2$. *Zakładamy*, że $\sigma_1^2 = \sigma_2^2$. Testujemy więc hipotezę o wartościach oczekiwanych, nie kwestionując założenia o równości wariancji. Odrzucamy H_0 , gdy

$$\frac{\bar{X}_1 - \bar{X}_2}{\sqrt{(n_1-1)S_1^2 + (n_2-1)S_2^2}} \sqrt{\frac{n_1 n_2}{n_1 + n_2}} (n_1 + n_2 - 2) > t,$$

$$t = t_{1-\alpha}(n_1 + n_2 - 2).$$

Uwaga: Oczywiście jest modyfikacja procedury testowania, gdy testujemy $H_0 : \mu_1 = \mu_2$ przeciwko $H_1 : \mu_1 \neq \mu_2$.

- ★ **Test** $H_0 : \sigma_1^2 \leq \sigma_2^2$ **przeciwko** $H_1 : \sigma_1^2 > \sigma_2^2$. Odrzucamy H_0 , gdy

$$\frac{S_1^2}{S_2^2} > f, \quad f = F_{1-\alpha}(n_1-1, n_2-1).$$

Jeśli testujemy $H_0 : \sigma_1^2 = \sigma_2^2$ przeciw $H_1 : \sigma_1^2 \neq \sigma_2^2$, to test na poziomie istotności α odrzuca H_0 gdy

$$\frac{S_1^2}{S_2^2} > f_2 \text{ lub } \frac{S_1^2}{S_2^2} < f_1,$$

$f_1 = F_{\alpha/2}(n_1-1, n_2-1)$, $f_2 = F_{1-\alpha/2}(n_1-1, n_2-1)$ są kwantylami rozkładu F-Snedecora.

Porównanie r populacji: Analiza wariancji.

Rozważmy r niezależnych próbek:

$$\begin{array}{ll} \text{próbka 1:} & Y_{11}, \dots, Y_{1n_1} \text{ z rozkładu } N(\mu_1, \sigma^2); \\ \dots & \dots \\ \text{próbka } j: & Y_{j1}, \dots, Y_{jn_j} \text{ z rozkładu } N(\mu_j, \sigma^2); \\ \dots & \dots \\ \text{próbka } r: & Y_{r1}, \dots, Y_{rn_r} \text{ z rozkładu } N(\mu_r, \sigma^2). \end{array}$$

Zakładamy tu równość wariancji wszystkich rozkładów. Interesować nas będzie hipoteza

$$H_0 : \mu_1 = \dots = \mu_r.$$

Hipoteza ta sprowadza się do stwierdzenia, że wszystkie próbki pochodzą z tego samego rozkładu.

Oznaczenia:

$$\bar{Y}_j = \frac{1}{n_j} \sum_{i=1}^{n_j} Y_{ji}, \quad \bar{Y} = \frac{1}{n} \sum_{j=1}^r n_j \bar{Y}_j = \frac{1}{n} \sum_{j=1}^r \sum_{i=1}^{n_j} Y_{ji}$$

są to odpowiednio – średnia dla j -tej próbki i średnia globalna.

Niech

$$SST = \sum_{j=1}^r \sum_{i=1}^{n_j} (Y_{ji} - \bar{Y})^2,$$

$$SSB = \sum_{j=1}^r n_j (\bar{Y}_j - \bar{Y})^2, \quad SSW = \sum_{j=1}^r \sum_{i=1}^{n_j} (Y_{ji} - \bar{Y}_j)^2.$$

SSB jest sumą kwadratów *pomiędzy* próbkami (ang. “Sum of Squares, Between”), SSW jest sumą kwadratów *wewnątrz* próbek (ang. “Within”), zaś SST jest *całkowitą* sumą kwadratów (ang. “Total”).

Tożsamość analizy wariancji: $SST = SSB + SSW$.

Za statystykę testową przyjmujemy iloraz

$$F = \frac{MSB}{MSW} = \frac{SSB/(r-1)}{SSW/(n-r)}.$$

Test ANOVA: Hipotezę H_0 odrzucamy, jeśli

$$F > F_{1-\alpha}(r-1, n-r),$$

gdzie $F_{1-\alpha}(r-1, n-r)$ oznacza kwantyl rozkładu F-Snedocora z $r-1$ stopniami swobody w liczniku i $n-r$ w mianowniku.

Cecha X jakościowa, o wartościach 0,1 (tak,nie).

Porównanie z „normą”.

Niech p oznacza wskaźnik struktury, \hat{p} – jego estymator.

★ **Test $H_0 : p \leq p_0$ przeciw alternatywie $H_1 : p > p_0$.**

$$\text{odrzucaamy } H_0, \text{ jeśli } \sqrt{n} \frac{\hat{p} - p_0}{\sqrt{p_0(1 - p_0)}} > z_{1-\alpha},$$

gdzie $z_{1-\alpha}$ oznacza kwantyl rozkładu $N(0, 1)$. Inaczej: odrzucaamy H_0 , jeśli $\hat{p} > p_0 + z_{1-\alpha} \sqrt{p_0(1 - p_0)}/\sqrt{n}$ (\hat{p} jest „istotnie większe” od p_0).

★ **Test $H_0 : p = p_0$ przeciw alternatywie $H_1 : p \neq p_0$.**

$$\text{odrzucaamy } H_0, \text{ jeśli } \sqrt{n} \frac{|\hat{p} - p_0|}{\sqrt{p_0(1 - p_0)}} > z_{1-\alpha/2},$$

gdzie $z_{1-\alpha/2}$ oznacza kwantyl rozkładu $N(0, 1)$. Inaczej: odrzucaamy H_0 , jeśli $|\hat{p} - p_0| > z_{1-\alpha/2} \sqrt{p_0(1 - p_0)}/\sqrt{n}$ (\hat{p} jest „istotnie różne” od p_0). *Uwaga:* To są rozwiązania przybliżone. Można je stosować gdy, powiedzmy, $np_0 \geq 5$ i $n(1 - p_0) \geq 5$.

Porównanie 2 wskaźników struktury.

Niech p_1 i p_2 oznaczają wskaźniki struktury w dwóch populacjach. Pobieramy próbki rozmiarów n_1 i n_2 z obu populacji i obserwujemy w próbkach odpowiednio K_1 i K_2 elementów wyróżnionych. Niech $\hat{p}_1 = K_1/n_1$ i $\hat{p}_2 = K_2/n_2$ – będą estymatorami obu wskaźników struktury, zaś $\hat{p} = (K_1 + K_2)/(n_1 + n_2)$.

★ **Test $H_0 : p_1 \leq p_2$ przeciw alternatywie $H_1 : p_1 > p_2$.**

$$\text{odrzucaamy } H_0, \text{ jeśli } \sqrt{\frac{n_1 n_2}{n_1 + n_2}} \frac{\hat{p}_1 - \hat{p}_2}{\sqrt{\hat{p}(1 - \hat{p})}} > z_{1-\alpha},$$

gdzie $z_{1-\alpha}$ oznacza kwantyl rozkładu $N(0, 1)$. (\hat{p}_1 jest „istotnie większe” od \hat{p}_2).

Cecha k -wartościowa, porównanie z „normą”.

Jakościowa cecha X ma k wartości w_1, \dots, w_k .

Hipoteza H_0 : rozkład cechy X w populacji jest dany tabelką:

wartość	w_1	\dots	w_i	\dots	w_k
prawdopodobieństwo	p_1	\dots	p_i	\dots	p_k

gdzie $p_i = \mathbb{P}(X = w_i)$. Oczywiście, $p_1 + \dots + p_k = 1$.

Dane mają postać „tabelki powtórzeń”:

wartość cechy X	w_1	\dots	w_i	\dots	w_k
liczba elementów próbki	N_1	\dots	N_i	\dots	N_k

gdzie N_i – liczba elementów próbki, dla których $X = w_i$ (oczywiście, $N_1 + \dots + N_k = n$).

Test zgodności chi-kwadrat. Obliczamy statystykę „chi-kwadrat”:

$$\chi^2 = \sum_{i=1}^k \frac{(N_i - np_i)^2}{np_i}.$$

Test na poziomie istotności α (w przybliżeniu):

odrzucaamy H_0 jeśli $\chi^2 > c$,

gdzie $c = \chi^2_{1-\alpha}(k-1)$ jest kwantylem rzędu $1 - \alpha$ rozkładu chi-kwadrat z $k - 1$ stopniami swobody.

Ogólny schemat budowania statystyki chi-kwadrat:

$$\chi^2 = \sum \frac{(\text{wielkość obserwowana} - \text{wielkość oczekiwana})^2}{\text{wielkość oczekiwana}}.$$

Dwie cechy jakościowe, badanie niezależności.

Dla pojedynczego elementu obserwujemy *parę* cech (X, Y) , przy czym X ma możliwe wartości $1, \dots, r$, zaś Y – wartości $1, \dots, s$ (te wartości należy traktować jako umowne „etykiety”, kodujące cechy jakościowe, nie jako liczby). Rozkład jest opisany dwuwymiarową tabelką (p_{ij}) o wierszach $i = 1, \dots, r$ i kolumnach $j = 1, \dots, s$, gdzie

$$p_{ij} = \mathbb{P}(X = i, Y = j).$$

Rozkłady „brzegowe” cech X i Y , rozpatrywanych oddzielnie:

$$\mathbb{P}(X = i) = p_{i\bullet} = \sum_{j=1}^s p_{ij},$$

$$\mathbb{P}(Y = j) = p_{\bullet j} = \sum_{i=1}^r p_{ij}.$$

Jeśli obserwujemy cechy (X, Y) dla n elementów próbki, możemy zbudować dwuwymiarową tabelkę (N_{ij}) , gdzie

$$N_{ij} = \text{liczba elementów próbki, dla których } (X, Y) = (i, j).$$

Jest to *tablica kontyngencji*. Wielkości „brzegowe” w tej tabelce oznaczymy

$$N_{i\bullet} = \sum_{j=1}^s N_{ij}, \quad N_{\bullet j} = \sum_{i=1}^r N_{ij}.$$

Test niezależności chi-kwadrat. Rozpatrzmy hipotezę, która stwierdza *niezależność* zmiennych X i Y :

$$\mathbb{P}(X = i, Y = j) = \mathbb{P}(X = i)\mathbb{P}(Y = j),$$

czyli

$$H_0 : p_{ij} = p_{i\bullet}p_{\bullet j}, \quad (i = 1, \dots, r; j = 1, \dots, s).$$

Statystyka do testowania hipotezy o niezależności jest następująca:

$$\chi^2 = \sum_{i=1}^r \sum_{j=1}^s \frac{(N_{ij} - N_{i\bullet}N_{\bullet j}/n)^2}{N_{i\bullet}N_{\bullet j}/n}.$$

Test:

$$\text{odrzucamy } H_0 \text{ jeśli } \chi^2 > c,$$

gdzie $c = \chi^2_{1-\alpha}((r-1)(s-1))$ jest kwantylem rzędu $1 - \alpha$ rozkładu chi-kwadrat.

Porównanie r wskaźników struktury.

Niech p_1, \dots, p_k będą wskaźnikami struktury w r populacjach (mamy cechę jakościową X o wartościach $0, 1$). Pobieramy próbkę rozmiaru n_i z i -tej populacji ($i = 1, \dots, r$). Niech K_i będzie liczbą elementów wyróżnionych w próbce z i -tej populacji.

Hipoteza

$$H_0 : p_1 = \dots = p_r.$$

Statystyka testowa chi-kwadrat przybiera tu postać:

$$\chi^2 = \sum_{i=1}^r \frac{(K_i - n_i \hat{p})^2}{n_i \hat{p}(1 - \hat{p})},$$

gdzie $\hat{p} = \sum_i K_i / \sum_i n_i$.

Test:

$$\text{odrzucamy } H_0 \text{ jeśli } \chi^2 > c,$$

gdzie $c = \chi^2_{1-\alpha}(r-1)$ jest kwantylem rzędu $1 - \alpha$ rozkładu chi-kwadrat.