

# Testy adaptacyjne dla problemu $k$ prób

Grzegorz Wyłupek

Instytut Matematyczny  
Polskiej Akademii Nauk  
Oddział Wrocław

# Problem testowania

Założmy, że  $X_{l1}, \dots, X_{ln_l}$ ,  $l = 1, \dots, k$ , są  $k$  niezależnymi próbkami losowymi, gdzie  $X_{lj}$  pochodzi z populacji o ciągłej dystrybucji  $F_l$ .  
Rozważmy problem testowania

$$\mathcal{H} : F_1 = \dots = F_k$$

przeciwko globalnej alternatywie

$$\mathcal{A} : F_i \neq F_j \text{ dla pewnych } 1 \leq i < j \leq k.$$

# Notacja

Niech  $N = \sum_{l=1}^k n_l$ ,  $p_l = n_l/N$ . Dla uproszczenia, będziemy identyfikować  $X_{11}, \dots, X_{1n_1}, \dots, X_{k1}, \dots, X_{kn_k}$  z

$X_1, \dots, X_{n_1}, \dots, X_{n_1+1}, \dots, X_{n_1+n_2}, \dots, X_{n_1+\dots+n_{k-1}+1}, \dots, X_N$ .

Niech  $R_{n_1+\dots+n_{l-1}+1}, \dots, R_{n_1+\dots+n_l}$  będzie wektorem rang  $X_{l1}, \dots, X_{ln_l}$  w połączonej próbie  $X_1, \dots, X_N$ .

Oznaczmy przez  $\mathbf{Z}_l$  zbiór indeksów odpowiadających  $l$ -tej próbie

$\mathbf{Z}_l = \{n_1 + \dots + n_{l-1} + 1, \dots, n_1 + \dots + n_l\}$  i połączmy

$\mathbf{Z} = \bigcup_{l=1}^k \mathbf{Z}_l = \{1, \dots, N\}$ .

Dla  $l = 1, \dots, k$ ;  $i = 1, \dots, N$  zdefiniujemy

$$c_{Ni}(l) = \sqrt{\frac{n_l(N - n_l)}{N}} \begin{cases} n_l^{-1} & \text{gdy } i \in \mathbf{Z}_l, \\ -(N - n_l)^{-1} & \text{gdy } i \in \mathbf{Z} - \mathbf{Z}_l, \end{cases}$$

$$W_r^{[2]}(l) = \sum_{j=1}^r \left\{ \sum_{i=1}^N c_{Ni}(l) b_j \left( \frac{R_i - 0.5}{N} \right) \right\}^2 - \text{gładka 2-próbkowa statystyka,}$$

gdzie  $b_0 \equiv 1, b_1, b_2, \dots$  jest bazą wielomianów Legendre'a na  $(0, 1)$ .

$$W_r = \sum_{l=1}^k (1 - p_l) W_r^{[2]}(l) = \sum_{j=1}^r \left[ \sum_{l=1}^k (1 - p_l) \left\{ \sum_{i=1}^N c_{Ni}(l) b_j \left( \frac{R_i - 0.5}{N} \right) \right\}^2 \right],$$

$W_r$  - gładka  $k$ -próbkowa statystyka.

$W_1$  -  $k$ -próbkowa statystyka Kruskala-Wallisa (1952);

$W_2$  - suma  $k$ -próbkowych statystyk Kruskala-Wallisa i Mooda;

$W_4$  -  $k$ -próbkowa statystyka Boosa (1986).

$W_r$  jest kombinacją statystyk  $W_r^{[2]}(l)$

dla podproblemów  $(\mathbf{Z}_l, \mathbf{Z} - \mathbf{Z}_l)$ ,  $l = 1, \dots, k$ .

$$W_r^{[2]}(l) = \sum_{j=1}^r \left\{ \sum_{i=1}^N c_{Ni}(l) b_j \left( \frac{R_i - 0.5}{N} \right) \right\}^2.$$

Łączny rozkład w próbach odpowiadających  $(\mathbf{Z}_l, \mathbf{Z} - \mathbf{Z}_l)$  oraz  $(\mathbf{Z}_{l'}, \mathbf{Z} - \mathbf{Z}_{l'})$ ,  $l \neq l'$ , może być różny. Dlatego będziemy szukać optymalnego  $r$ , opierając się na danych, w każdym podproblemie dwóch prób  $(\mathbf{Z}_l, \mathbf{Z} - \mathbf{Z}_l)$  oddzielnie.

$$S(l) = \min\{r : W_r^{[2]}(l) - r \log N \geq W_j^{[2]}(l) - j \log N, 1 \leq r, j \leq d(N)\},$$

$d(N)$  - liczba modeli na liście.

$S(l)$  - uproszczona reguła Schwarzera (1978) BIC;

patrz Janic-Wróblewska i Ledwina (2000).

Alternatywna reguła wyboru modelu:

$A(l)$  - uproszczona reguła Akaike (1973) AIC

$$A(l) = \min\{r : W_r^{[2]}(l) - 2r \geq W_j^{[2]}(l) - 2j, 1 \leq r, j \leq d(N)\}.$$

Wybór między regułami  $S(l)$  i  $A(l)$ , lub równoważnie pomiędzy karami, powinien być oparty na danych - Ingłot i Ledwina (2006).

Mamy  $W_r^{[2]}(l) = \sum_{j=1}^r \{\hat{B}_j(l)\}^2$ ,

gdzie  $\hat{B}_j(l) = \sum_{i=1}^N c_{Ni}(l) b_j([R_i - 0.5]/N)$ .

Przy ustalonym  $l$  i pod warunkiem  $\mathcal{H}$ ,

$\hat{B}_1(l), \hat{B}_2(l), \dots, \hat{B}_{d(N)}(l)$  są asymptotycznie nieskorelowanymi zmiennymi losowymi o rozkładzie normalnym  $N(0, 1)$ .

Nasza decyzja o karze będzie oparta  
na zmiennych  $\hat{B}_j(l)$ ,  $j = 1, \dots, d(N)$ .

Dla każdego  $d(N) \rightarrow \infty$  i  $c \geq 0$ , pod warunkiem  $\mathcal{A}$ , mamy

$$P\left(\max_{1 \leq j \leq d(N)} |\hat{B}_j(l)| \leq \sqrt{c \log N}\right) \rightarrow 0, \text{ gdy } N \rightarrow \infty.$$

Ponadto, jeżeli  $\mathcal{H}$  jest prawdziwa, to dla każdego  $\gamma \in (0, 1/9]$ ,  
 $d(N) = o([N/\log N]^\gamma)$ , i  $c > 2\gamma$  otrzymujemy

$$P\left(\max_{1 \leq j \leq d(N)} |\hat{B}_j(l)| \leq \sqrt{c \log N}\right) \rightarrow 1, \text{ gdy } N \rightarrow \infty.$$

Przy ustalonym  $c > 2\gamma$ , zdefiniujemy zmienną losową

$$I_N(c, l) = \mathbf{1}\left(\max_{1 \leq j \leq d(N)} |\hat{B}_j(l)| \leq \sqrt{c \log N}\right),$$

gdzie  $\mathbf{1}(\mathcal{E})$  jest indykátorem zdarzenia  $\mathcal{E}$ .

Na mocy powyższego,  $I_N(c, l)$  może służyć jako przełącznik pomiędzy karą Schwarz'a i karą Akaike.

Zdefiniujemy

$$\Pi_r(c, l) = (r \log N)[I_N(c, l)] + (2r)[1 - I_N(c, l)]$$

i nową regułą wyboru modelu

$$T(l) = \min\{r : W_r^{[2]}(l) - \Pi_r(c, l) \geq W_j^{[2]}(l) - \Pi_j(c, l), 1 \leq r, j \leq d(N)\}.$$



Nowy test adaptacyjny odrzuca  $\mathcal{H}$  dla dużych wartości statystyki

$$W_T = \sum_{l=1}^k (1 - p_l) W_{T(l)}^{[2]}(l).$$

Powyższa formuła definiuje całą klasę statystyk indeksowanych parametrem  $c \in [0, +\infty]$ .

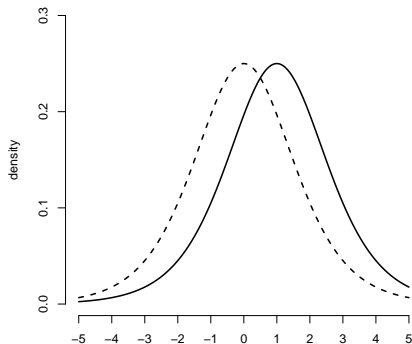
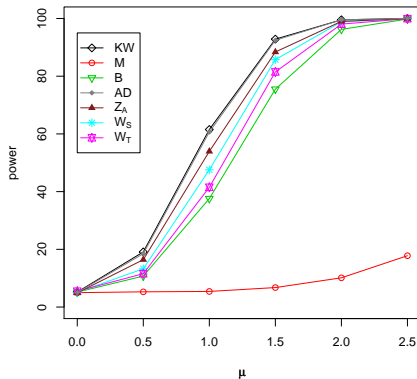
Wybór  $c$  w nowej karze  $\Pi_r(c, l)$  i odpowiadającej statystyce  $W_T$  ma decydujący wpływ na zachowanie mocy testu adaptacyjnego. Pozwala to na konstrukcję testów typu omnibus jak i wyspecjalizowanych wariantów. Wyniki asymptotyczne zachodzą dla szerokiego zakresu  $c$  i  $d(N)$ .

Niech  $\chi_{k-1}^2$  będzie zmienną losową o rozkładzie chi-kwadrat z  $k - 1$  stopniami swobody. Ponadto niech  $\xrightarrow{\mathcal{P}}$  i  $\xrightarrow{\mathcal{D}}$  oznaczają odpowiednio zbieżność według prawdopodobieństwa i według rozkładu.

**Twierdzenie 1.** *Założmy, że  $\mathcal{H}$  jest prawdziwa,  $d(N) = o([N/\log N]^\gamma)$  i  $0 < \gamma \leq 1/9$ . Ponadto niech  $c > 2\gamma$ . Wtedy  $T(l) \xrightarrow{\mathcal{P}} 1$ ,  $l = 1, \dots, k$ , i  $W_T \xrightarrow{\mathcal{D}} \chi_{k-1}^2$ , gdy  $N \rightarrow \infty$ .*

**Twierdzenie 2.** *Założmy, że  $d(N)$  jest jak w twierdzeniu 1 oraz  $d(N) \rightarrow \infty$ , gdy  $N \rightarrow \infty$ . Niech  $c > 2\gamma$ . Wtedy test adaptacyjny odrzucający  $\mathcal{H}$  dla dużych wartości  $W_T$  jest zgodny dla każdej alternatywy z  $\mathcal{A}$ .*

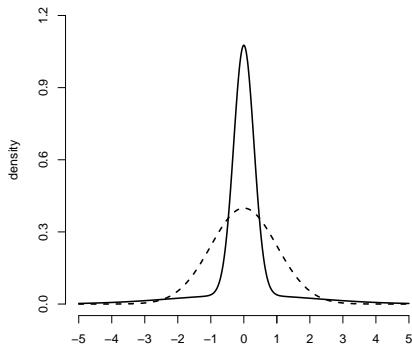
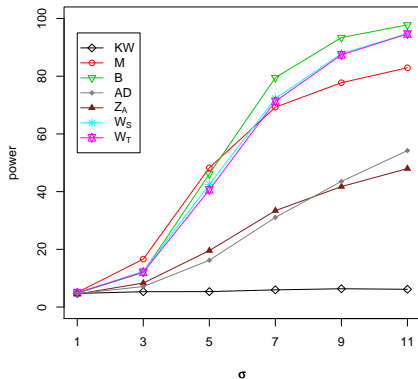
$\mathcal{H}$  : Logistic(0)/Logistic(0)/Logistic(0)     $\mathcal{A}(\mu)$  : Logistic(0)/Logistic( $\mu$ )/Logistic(0)



Rysunek 1. Lewy panel: porównanie mocy  $n_1 = n_2 = n_3 = 30$ ,  $\alpha = .05$ ,  $c = 2.3$ ,  $d(N) = 12$ . Oparte na  $10^4$  powtórzeń MC. Moce przemnożone przez 100. Prawy panel: - - - Logistic(0), — Logistic(1).

$\mathcal{H} : N(0, 1)/N(0, 1)/N(0, 1)$

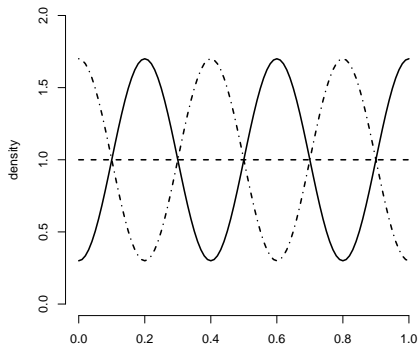
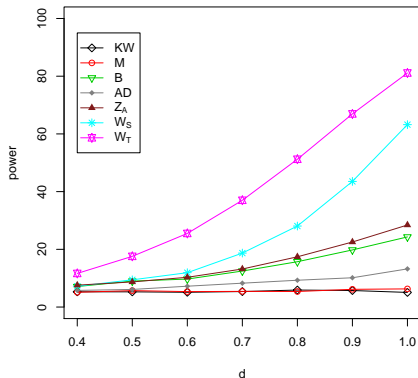
$\mathcal{A}(\sigma) : N(0, 1)/SC(\sigma)/N(0, 1)$



Rysunek 2. Lewy panel: porównanie mocy  $n_1 = n_2 = n_3 = 30$ ,  $\alpha = .05$ ,  $c = 2.3$ ,  $d(N) = 12$ . Oparte na  $10^4$  powtórzeń MC. Moce przemnożone przez 100. Prawy panel: - - -  $N(0, 1)$ , —  $SC(7)$ .

$\mathcal{H} : U(0, 1)/U(0, 1)/U(0, 1)$

$\mathcal{A}(d) : U(0, 1)/C_5(d)/C_5(-d)$



Rysunek 3. Lewy panel: porównanie mocy  $n_1 = n_2 = n_3 = 30$ ,  $\alpha = .05$ ,  $c = 2.3$ ,  $d(N) = 12$ . Oparte na  $10^4$  powtórzeń MC. Moce przemnożone przez 100. Prawy panel: - - -  $U(0, 1)$ , —  $C_5(0.8)$ , - · - · -  $C_5(-0.8)$ .