

STUDIA METODOLOGICZNE

Jacek WESOŁOWSKI
Jakub TARCZYŃSKI

Podstawy matematyczne technik imputacyjnych

Streszczenie. *W artykule przedstawiono podstawy metodologii imputacyjnej (w tym metodologii wielokrotnej imputacji), koncentrując się na wyjaśnieniu matematycznej strony zagadnień. Analizowano sytuację, gdy obserwacje tworzące pierwotną próbkę są niezależnymi zmiennymi losowymi o jednakowym rozkładzie, a braki odpowiedzi pojawiają się losowo w sposób niezależny od obserwacji. W szczególności wskazano na problemy pojawiające się, gdy w imputacji wielokrotnej stosowany jest standardowy estymator Rubina wariancji estymatora wielokrotnej imputacji i wskazano na możliwe ulepszenie tego popularnego estymatora. Punktem wyjścia analiz jest sytuacja, gdy za pojawianie się braków odpowiedzi odpowiada mechanizm deterministyczny.*

Słowa kluczowe: imputacja, imputacja wielokrotna, estymator imputacyjny, estymator Rubina, imputacja średnią, imputacja typu hot-deck, imputacja regresyjna.

Słowem imputacja określa się zespół metod uzupełniania braków w próbce w taki sposób, aby zastosowanie estymatorów standardowych, tzn. takich, które byłyby użyte, gdyby w próbce nie było braków obserwacji, prowadziło do efektywnej procedury estymacyjnej.

Literatura dotycząca technik imputacyjnych jest bardzo obszerna. Składają się na nią setki specjalistycznych artykułów w czasopismach statystycznych, artykuły przeglądowe, np. Andridge i Little (2010), Donders, van der Heijden, Stijnen i Moons (2006) czy Norazian Ramli, Yahaya, Ramli i Yusof (2013) oraz kilka większych publikacji monograficznych, takich jak: Rubin (1987), Little i Rubin

(2002), de Waal, Pannekoek i Scholtus (2011) czy Garson (2012) lub van Buuren (2012). Literatura ta najczęściej dotyczy opisu konkretnych metod, ich zastosowań i porównań. Techniki te bywają dostępne w pakietach statystycznych (Horton i Lipsitz, 2001; van Buuren i Groothuis-Oudshoorn, 2011; Misztal, 2012). We wszystkich tych publikacjach stosunkowo mały nacisk kładziony jest na podstawy teoretyczne.

Okazuje się, że w standardowym modelu statystycznym niezależnych obserwacji o jednakowych rozkładach daje się precyzyjnie opisać (w formie twierdzeń matematycznych) podstawy teoretyczne technik imputacyjnych. Przedstawiane opracowanie jest próbą takiego właśnie systematycznego i matematycznego uściślenia uniwersalnych paradygmatów imputacji. W rozdziale I rozważana jest uproszczona sytuacja, gdy zbiór braków obserwacji jest ustalony. Przy tym założeniu wyprowadzono wzory ogólne na estymator imputacyjny średniej i jego parametry oraz zastosowano je w czterech podstawowych schematach imputacyjnych: imputacji średnią, dwóch procedurach typu hot-deck oraz imputacji regresyjnej. W rozdziale II analizowano imputację wielokrotną w modelu z rozdziału I, a następnie opracowany schemat ogólny wykorzystano w sytuacji, gdy w każdym kroku procedury stosowana jest imputacja typu hot-deck. Używany w literaturze termin „imputacja hot-deck”, w tym opracowaniu — nieco rozszerzony — odnosi się do sytuacji, gdy elementy imputowane są zgodnie z pewną procedurą losową, która zależna jest od obserwowanej części próbki.

W rozdziałach III i IV rozważany jest bardziej realistyczny losowy model pojawiania się braków odpowiedzi, tzn. taki, w którym każda obserwacja może być obecna w próbkę lub nie, przy czym zdarzenie to ma charakter losowy, czyli można mu przypisać pewne prawdopodobieństwo. Poprzez zastosowanie techniki warunkowania, rezultaty otrzymane w rozdziałach III i IV można stosunkowo prosto wywieść z wyników przedstawionych w rozdziałach I i II. W takim podejściu istotną rolę odgrywa przyjęte w artykule założenie o niezależności statystycznej obserwacji i braków odpowiedzi. W literaturze mówi się w takiej sytuacji o brakach odpowiedzi typu MCAR (*missing completely at random*), np. Little i Rubin (2002). To podejście jest ważne, bo choć często nie jest właściwym modelem dla całej populacji, to okazuje się dobrym przybliżeniem sytuacji rzeczywistej w odpowiednio wybranych podpopulacjach (chodzi o klasy elementów podobnych pod względem skłonności do udzielania odpowiedzi).

Należy zwrócić uwagę na, prezentowane w opracowaniu, ulepszenie klasycznego estymatora Rubina służącego do estymacji wariancji estymatora imputacyjnego w przypadku imputacji wielokrotnej. Zasadnicza część artykułu ma matematyczną strukturę twierdzeń i dowodów. Najciekawsze wyniki matematyczne (w rozdziałach II i IV) zilustrowano eksperymentami numerycznymi. Konkluzje przedstawiono w rozdziale V.

I. DETERMINISTYCZNY ZBIÓR BRAKÓW OBSERWACJI
— IMPUTACJA JEDNOKROTNA

I.0. Ogólny schemat imputacyjny

Zmienne losowe X_1, \dots, X_n stanowią pierwotną próbkę, $\mathbf{X} = (X_1, \dots, X_n)$. Niech $\hat{\theta} = g(\mathbf{X})$ oznacza estymator parametru θ , gdzie g jest odpowiednią funkcją n -argumentową. Rozważać będziemy sytuację, gdy obserwowane są tylko zmienne X_i , $i \in R \subsetneq \{1, \dots, n\}$, czyli R jest zbiorem indeksów obserwowanych elementów próbki pierwotnej. W konsekwencji pierwotny wektor obserwacji dzieli się naturalnie na dwa wektory: wektor odpowiedzi obserwowanych $\mathbf{X}_R = (X_i, i \in R)$ oraz wektor odpowiedzi nieobserwowanych $\mathbf{X}_{R^c} = (X_i, i \in R^c = \{1, \dots, n\} \setminus R)$. Taką niepełną próbkę \mathbf{X}_R można uzupełnić przyjmując $\tilde{\mathbf{X}} = (\tilde{X}_1, \dots, \tilde{X}_n)$, gdzie $\tilde{X}_i = X_i$ dla $i \in R$ są zmiennymi obserwowanymi oraz \tilde{X}_i , $i \in R^c$ są tzw. zmiennymi imputowanymi, które mogą zależeć (i często zależą) od zmiennych obserwowanych. Często zmienne są imputowane pewnymi funkcjami tych zmiennych, ale zdarzają się też sytuacje, w których nie ma bezpośredniej zależności funkcyjnej, a zmienne obserwowane wpływają jedynie na rozkład prawdopodobieństwa zmiennych imputowanych.

Definicja I.0.1. Imputacyjną wersją estymatora $\hat{\theta} = g(\mathbf{X})$ parametru θ nazywamy estymator postaci

$$\tilde{\theta}_{imp} = g(\tilde{\mathbf{X}}).$$

W skrócie mówimy, że $\tilde{\theta}_{imp}$ jest estymatorem imputacyjnym parametru θ .

Zakładamy, że zmienne X_1, \dots, X_n są niezależne i mają jednakowe rozkłady, takie jak pewna zmienna losowa X o wartości oczekiwanej $E X = \mu$ i wariancji $\text{Var } X = \sigma^2$, które są nieznanne.

Do estymacji średniej $\theta = \mu$ standardowo stosowany jest estymator

$$\hat{\mu} = g(\mathbf{X}) = \bar{X} = \frac{1}{n} \sum_{i=1}^n X_i.$$

Jest to estymator nieobciążony, czyli $E \hat{\mu} = \mu$, a jego wariancja wynosi $\text{Var } \hat{\mu} = \frac{\sigma^2}{n}$.

Do estymacji parametru σ^2 standardowo stosowany jest estymator

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

Jest to estymator nieobciążony, czyli $E S^2 = \sigma^2$.

Zatem nieobciążony estymator $d^2(\hat{\mu})$ wariancji $\text{Var } \hat{\mu}$ estymatora $\hat{\mu}$ ma postać $d^2(\hat{\mu}) = \frac{1}{n} S^2$.

Z centralnego twierdzenia granicznego (i tzw. twierdzenia Śluckiego) wynika, że $\frac{\hat{\mu} - \mu}{S} \sqrt{n}$ ma asymptotycznie rozkład standardowy normalny, co pozwala na wnioskowanie o przybliżonej precyzji estymatorów, rozumianej jako przedział ufności na zadanym poziomie ufności.

Zadanie imputacji pojawia się, gdy obserwowane są tylko zmienne losowe $X_i, i \in R \subsetneq \{1, \dots, n\}$. Polega ono na uzupełnieniu obserwowanej części próbki o część nieobserwowaną $\tilde{X}_i, i \in R^c$. Niech r oznacza liczbę obserwowanych elementów próbki, czyli $r = \#(R)$. Dla potrzeb rozdziałów I i II zakładamy, że podzbiór $R \subsetneq \{1, \dots, n\}$, a zatem i jego liczebność r jest ustalona. W kolejnych rozdziałach odejdziemy od tego założenia.

Nowa próbka ma postać $\tilde{\mathbf{X}} = (\tilde{X}_1, \dots, \tilde{X}_n)$, przy czym zakłada się, że

$$E \tilde{X}_i = \tilde{\mu}_i = \tilde{\mu} \quad \text{oraz} \quad \text{Var } \tilde{X}_i = \tilde{\sigma}_i^2 = \tilde{\sigma}^2, \quad \text{jeśli } i \in R^c$$

oraz

$$\text{Corr}(\tilde{X}_i, \tilde{X}_j) = \tilde{\rho}_{ij} = \tilde{\rho}, \quad \text{jeśli } i, j \in R^c, i \neq j,$$

$$\text{Corr}(X_i, \tilde{X}_j) = \rho_{ij} = \rho, \quad \text{jeśli } i \in R, j \in R^c.$$

Zajmiemy się teraz, zgodnie z Definicją I.0.1, imputacyjnym estymatorem średniej μ .

TWIERDZENIE I.0.1. *Imputacyjny estymator średniej μ ma postać*

$$\hat{\mu}_{imp} = \frac{r \bar{X}_R + (n-r) \tilde{\bar{X}}}{n}, \quad (\text{I.0.1})$$

gdzie $\bar{X}_R = \frac{1}{r} \sum_{i \in R} X_i$ oraz $\tilde{\bar{X}} = \frac{1}{n-r} \sum_{i \in R^c} \tilde{X}_i$.

Jego wartość oczekiwana wynosi

$$E \hat{\mu}_{Imp} = \frac{r}{n} \mu + \frac{n-r}{n} \tilde{\mu}, \quad (I.0.2)$$

a jego obciążenie wynosi

$$B \hat{\mu}_{Imp} = \frac{n-r}{n} (\tilde{\mu} - \mu),$$

więc estymator $\hat{\mu}_{Imp}$ jest nieobciążony tylko gdy $\tilde{\mu} = \mu$.

Dowód. Z postaci funkcji g występującej w definicji estymatora $\hat{\mu}$ wynika, że

$$\hat{\mu}_{Imp} = g(\tilde{X}) = \frac{1}{n} \left(\sum_{i \in R} X_i + \sum_{i \in R^c} \tilde{X}_i \right),$$

co jest równoważne z (I.0.1).

Wzór (I.0.1) implikuje

$$E \hat{\mu}_{Imp} = \frac{1}{n} \left(r E \bar{X}_R + (n-r) E \bar{\tilde{X}} \right).$$

Wzór (I.0.2) wynika z tej równości, ponieważ $E \bar{X}_R = \mu$, a $E \bar{\tilde{X}} = \tilde{\mu}$.

Wzór na obciążenie jest natychmiastową konsekwencją definicji obciążenia $B \hat{\mu}_{Imp} = E \hat{\mu}_{Imp} - \mu$ i wzoru (I.0.2). ■

TWIERDZENIE I.0.2. *Wariancja imputacyjnego estymatora średniej $\hat{\mu}_{Imp}$ ma postać*

$$\text{Var } \hat{\mu}_{Imp} = \frac{r\sigma^2 + (n-r)\tilde{\sigma}^2 + (n-r)(n-r-1)\tilde{\rho}\tilde{\sigma}^2 + 2r(n-r)\rho\sigma\tilde{\sigma}}{n^2} \quad (I.0.3)$$

Dowód. Z postaci estymatora $\hat{\mu}_{Imp}$ i podstawowych własności wariancji otrzymujemy

$$\text{Var } \hat{\mu}_{Imp} = \frac{1}{n^2} \left(r^2 \text{Var } \bar{X}_R + (n-r)^2 \text{Var } \bar{\tilde{X}} + 2r(n-r) \text{Cov} \left(\bar{X}_R, \bar{\tilde{X}} \right) \right). \quad (I.0.4)$$

$$\text{Ale } \text{Var } \bar{X}_R = \frac{\sigma^2}{r}.$$

Z kolei

$$\begin{aligned}\text{Var}\bar{\tilde{X}} &= \frac{1}{(n-r)^2} \left(\sum_{i \in R^c} \text{Var}\tilde{X}_i + \sum_{i, j \in R^c, i \neq j} \text{Cov}(\tilde{X}_i, \tilde{X}_j) \right) = \\ &= \frac{1}{(n-r)^2} \left((n-r)\tilde{\sigma}^2 + (n-r)(n-r-1)\tilde{\rho}\tilde{\sigma}^2 \right).\end{aligned}$$

A zatem

$$\text{Var}\bar{\tilde{X}} = \frac{1 + (n-r-1)\tilde{\rho}}{n-r} \tilde{\sigma}^2. \quad (\text{I.0.5})$$

Natomiast

$$\text{Cov}(\bar{X}_R, \bar{\tilde{X}}) = \frac{1}{r(n-r)} \sum_{j \in R} \sum_{i \in R^c} \text{Cov}(X_j, \tilde{X}_i) = \rho\sigma\tilde{\sigma}. \quad (\text{I.0.6})$$

Wstawiając te trzy wyrażenia do wzoru (I.0.4) otrzymujemy

$$\text{Var}\hat{\mu}_{Imp} = \frac{1}{n^2} \left(r\sigma^2 + (n-r)(1 + (n-r-1)\tilde{\rho})\tilde{\sigma}^2 + 2r(n-r)\rho\sigma\tilde{\sigma} \right)$$

Ten wzór jest równoważny wzorowi (I.0.3).■

Imputacyjna estymacja wariancji polega na wykorzystaniu standardowego estymatora wariancji do próbki imputowanej \tilde{X} , czyli zgodnie z definicją I.0.1

$$S_{Imp}^2 = S^2(\tilde{X}) = \frac{1}{n-1} \left(\sum_{j \in R} (X_j - \hat{\mu}_{Imp})^2 + \sum_{i \in R^c} (\tilde{X}_i - \hat{\mu}_{Imp})^2 \right). \quad (\text{I.0.7})$$

TWIERDZENIE I.0.3. *Imputacyjny estymator wariancji S_{Imp}^2 ma postać*

$$S_{Imp}^2 = \frac{1}{n-1} \left((r-1)S_R^2 + (n-r-1)\tilde{S}^2 + \frac{r(n-r)}{n} (\bar{X}_R - \bar{\tilde{X}})^2 \right), \quad (\text{I.0.8})$$

gdzie $S_R^2 = \frac{1}{r-1} \sum_{j \in R} (X_j - \bar{X}_R)^2$ oraz $\tilde{S}^2 = \frac{1}{n-r-1} \sum_{i \in R^c} (\tilde{X}_i - \bar{\tilde{X}})^2$.

Dowód. Ze wzoru (I.0.7) i definicji $\hat{\mu}_{Imp}$ mamy

$$S_{Imp}^2 = \frac{1}{n-1} \left(\sum_{i \in R} \left(X_i - \frac{1}{n} \left(r\bar{X}_R + (n-r)\tilde{X} \right) \right)^2 + \sum_{j \in R^c} \left(\tilde{X}_j - \tilde{X} + \frac{r}{n} \left(\tilde{X} - \bar{X}_R \right) \right)^2 \right).$$

Wykonując kwadratowanie pod obiema sumami i korzystając z faktu, że $\sum_{i \in R} (X_i - \bar{X}_R) = 0$ oraz $\sum_{i \in R^c} (\tilde{X}_i - \tilde{X}) = 0$ otrzymujemy

$$\begin{aligned} S_{Imp}^2 &= \frac{1}{n-1} \left(\sum_{i \in R} (X_i - \bar{X}_R)^2 + r \frac{(n-r)^2}{n^2} (\bar{X}_R - \tilde{X})^2 + \sum_{i \in R^c} (\tilde{X}_i - \tilde{X})^2 + (n-r) \frac{r^2}{n^2} (\bar{X}_R - \tilde{X})^2 \right) = \\ &= \frac{1}{n-1} \left((r-1)S_R^2 + (n-r-1)\tilde{S}^2 + \left(r \frac{(n-r)^2}{n^2} + (n-r) \frac{r^2}{n^2} \right) (\bar{X}_R - \tilde{X})^2 \right), \end{aligned}$$

co kończy dowód. ■

TWIERDZENIE I.0.4. *Wartość oczekiwana imputacyjnego estymatora wariancji S_{Imp}^2 ma postać*

$$\begin{aligned} ES_{Imp}^2 &= \frac{r}{n} \sigma^2 + \frac{n-r}{n} \tilde{\sigma}^2 - \frac{(n-r)(n-r-1)}{n(n-1)} \tilde{\rho} \tilde{\sigma}^2 - \\ &+ 2 \frac{r(n-r)}{n(n-1)} \rho \sigma \tilde{\sigma} + \frac{r(n-r)}{n(n-1)} (\mu - \tilde{\mu})^2. \end{aligned} \tag{I.0.9}$$

Dowód. Mamy $ES_R^2 = \sigma^2$ oraz $E\tilde{S}^2 = \sigma^2(1 - \tilde{\rho})$. Pierwszy fakt jest powszechnie znany, a drugi wynika z następującego rachunku

$$\begin{aligned} E\tilde{S}^2 &= E \left(\frac{1}{n-r-1} \sum_{i \in R^c} \tilde{X}_i^2 - \frac{1}{(n-r-1)(n-r)} \left(\sum_{i \in R^c} \tilde{X}_i \right)^2 \right) = \frac{1}{n-r-1} \sum_{i \in R^c} E\tilde{X}_i^2 - \\ &+ \frac{1}{(n-r-1)(n-r)} \left(\sum_{i \in R^c} E\tilde{X}_i^2 + \sum_{i, j \in R^c, i \neq j} E\tilde{X}_i \tilde{X}_j \right) = \frac{n-r}{n-r-1} (\tilde{\sigma}^2 + \tilde{\mu}^2) - \\ &+ \frac{1}{(n-r-1)(n-r)} \left((n-r)(\tilde{\sigma}^2 + \tilde{\mu}^2) + (n-r)(n-r-1)(\tilde{\rho}\tilde{\sigma}^2 + \tilde{\mu}^2) \right) = \tilde{\sigma}^2(1 - \tilde{\rho}). \end{aligned}$$

Co więcej

$$\begin{aligned} E\left(\bar{X}_R - \tilde{X}\right)^2 &= E\bar{X}_R^2 + E\tilde{X}^2 - 2E\bar{X}_R\tilde{X} = \text{Var}\bar{X}_R + \left(E\bar{X}_R\right)^2 + \\ &+ \text{Var}\tilde{X} + \left(E\tilde{X}\right)^2 - 2\text{Cov}\left(\bar{X}_R, \tilde{X}\right) - 2E\bar{X}_R E\tilde{X}. \end{aligned}$$

Wykorzystując wzory (I.0.5) oraz (I.0.6) otrzymujemy

$$E\left(\bar{X}_R - \tilde{X}\right)^2 = \frac{\sigma^2}{r} + \frac{\tilde{\sigma}^2}{n-r} + \frac{n-r-1}{n-r}\tilde{\rho}\tilde{\sigma}^2 - 2\rho\sigma\tilde{\sigma} + (\mu - \tilde{\mu})^2.$$

Ostatecznie ze wzoru (I.0.8) i powyżej uzyskanych tożsamości wynika, że

$$\begin{aligned} (n-1)E S_{Imp}^2 &= (r-1)\sigma^2 + (n-r-1)\tilde{\sigma}^2(1-\tilde{\rho}) + \\ &+ \frac{r(n-r)}{n}\left(\frac{\sigma^2}{r} + \frac{\tilde{\sigma}^2}{n-r} + \frac{n-r-1}{n-r}\tilde{\rho}\tilde{\sigma}^2 - 2\rho\sigma\tilde{\sigma} + (\mu - \tilde{\mu})^2\right), \end{aligned}$$

co po uproszczeniu daje wzór (I.0.9). ■

1.1. Imputacja średnią

Technika ta polega na przypisaniu każdemu nieobserwowanemu elementowi próbki średniej z części obserwowanej próbki, czyli

$$\tilde{X}_i = \frac{1}{r} \sum_{j \in R} X_j = \bar{X}_R, \quad i \in R^c. \quad (\text{I.1.1})$$

TWIERDZENIE I.1.1. *Imputacyjny estymator wartości oczekiwanej w przypadku imputacji średnią ma postać*

$$\hat{\mu}_{Imp} = \bar{X}_R. \quad (\text{I.1.2})$$

Jego wartość oczekiwana oraz wariancja wynoszą, odpowiednio

$$E\hat{\mu}_{Imp} = \mu \quad \text{oraz} \quad \text{Var}\hat{\mu}_{Imp} = \frac{\sigma^2}{r}. \quad (\text{I.1.3})$$

Dowód. Ze wzoru (I.1.1) wynika, że $\tilde{X} = \bar{X}_R$. Zatem wzór (I.0.1) prowadzi natychmiast do (I.1.2). Ponieważ $\tilde{\mu} = E\tilde{X}_i = E\bar{X}_R = \mu$, więc (I.0.2) daje pierwszy ze wzorów (I.1.3).

Dodatkowo z (I.1.1) wynika, że $\tilde{\sigma}^2 = \text{Var } \tilde{X}_i = \text{Var } \bar{X}_R = \frac{\sigma^2}{r}$ oraz

$$\tilde{\rho} = \text{Corr}(\tilde{X}_i, \tilde{X}_j) = \text{Corr}(\bar{X}_R, \bar{X}_R) = 1,$$

$$\rho = \text{Corr}(X_i, \tilde{X}_j) = \text{Corr}(X_i, \bar{X}_R) = \sqrt{r} \frac{\text{Cov}(X_i, \bar{X}_R)}{\sigma^2} = \sqrt{r} \frac{\frac{\sigma^2}{\sqrt{r}}}{\sigma^2} = \frac{1}{\sqrt{r}}.$$

Zatem zgodnie ze wzorem (I.0.3)

$$\text{Var } \hat{\mu}_{Imp} = \frac{r\sigma^2 + (n-r)\frac{\sigma^2}{r} + (n-r)(n-r-1)\frac{\sigma^2}{r} + 2r(n-r)\frac{1}{\sqrt{r}}\sigma\frac{\sigma}{\sqrt{r}}}{n^2}.$$

Po uproszczeniu otrzymujemy drugi ze wzorów (I.1.3). ■

Oczywiście wzory (I.1.3) wynikają wprost z postaci (I.1.2) estymatora imputacyjnego. W dowodzie, wyprowadzając je odpowiednio ze wzorów (I.0.2) oraz (I.0.3), chcieliśmy zaznaczyć uniwersalność proponowanego podejścia.

TWIERDZENIE I.1.2. *Imputacyjny estymator wariancji w przypadku imputacji średnią ma postać*

$$S_{Imp}^2 = \frac{r-1}{n-1} S_R^2, \quad (\text{I.1.4})$$

a jego wartość oczekiwana wynosi

$$\text{ES}_{Imp}^2 = \frac{r-1}{n-1} \sigma^2. \quad (\text{I.1.5})$$

Nieobciążony estymator wariancji estymatora $\hat{\mu}_{Imp}$ ma postać

$$v^2(\hat{\mu}_{Imp}) = \frac{n-1}{r(r-1)} S_{Imp}^2. \quad (\text{I.1.6})$$

Dowód. Zgodnie ze wzorem (I.1.1) mamy $\tilde{S}^2 = 0$ oraz $\tilde{X} = \bar{X}_R$. Zatem (I.1.4) wynika bezpośrednio ze wzoru (I.0.8).

Natomiast zgodnie ze wzorem (I.0.9) na wartość oczekiwaną imputacyjnego estymatora wariancji mamy

$$ES_{imp}^2 = \frac{r}{n} \sigma^2 + \frac{n-r}{n} \frac{\sigma^2}{r} - \frac{(n-r)(n-r-1)}{n(n-1)} \frac{\sigma^2}{r} - 2 \frac{r(n-r)}{n(n-1)} \frac{1}{\sqrt{r}} \sigma \frac{\sigma}{\sqrt{r}},$$

skąd po uproszczeniach otrzymuje się wzór (I.1.5).

Wzór (I.1.6) wynika z porównania wzoru (I.1.5) i drugiego ze wzorów (I.1.3).■

Wzór (I.1.5) wynika również bezpośrednio z postaci (I.1.4), ponieważ S_R^2 jest nieobciążonym estymatorem wariancji σ^2 . Wyprowadziliśmy go jednak ze wzoru (I.0.9), aby, podobnie jak w przypadku wzorów (I.0.2) i (I.0.3) (użytych do wyprowadzenia (I.1.3)), wskazać na jego uniwersalność.

Zauważmy, że ze wzoru (I.1.5) wynika, iż estymator S_{imp}^2 jest obciążonym estymatorem wariancji.

1.2. Imputacja typu hot-deck

1.2.1. Losowanie spośród respondentów

Technika ta polega na przypisaniu każdemu nieobserwowanemu elementowi próbki elementu wylosowanego spośród elementów obserwowanych. W najprostszej sytuacji jest to losowanie proste ze zwracaniem.

Niech K_i będzie numerem elementu wylosowanym dla i -tego nierespondenta, $i \in R^c$. Wtedy zmienne losowe K_i , $i \in R^c$, są niezależne i mają ten sam rozkład $P(K_i = j) = \frac{1}{r}$, $j \in R$, $i \in R^c$ oraz wektory losowe $\mathbf{K} = (K_i, i \in R^c)$ i $\mathbf{X}_R = (X_j, j \in R)$ są niezależne.

W konsekwencji

$$\tilde{X}_i = X_{K_i}, \quad i \in R^c. \quad (\text{I.2.1})$$

Uwaga. W przypadku obserwacji dwuwymiarowych $((X_j, Y_j), j \in R, Y_i, i \in R^c)$ stosowana bywa metoda wyboru najbliższego sąsiada. Polega ona na ustaleniu dla każdego $i \in R^c$ takiego $K_i \in R$, dla którego odległość $|Y_{K_i} - Y_i|$ jest najmniejsza spośród $|Y_j - Y_i|$, $j \in R$. Przy założeniu, że wektory losowe $\mathbf{w}_i := (X_i, Y_i)$, $i = 1, \dots, n$, są niezależne i mają jednakowe rozkłady (takie jak rozkład pary (X, Y)), zmienne losowe K_i , $i \in R^c$, mają rozkład jednostajny tak jak w opisanym wyżej modelu, natomiast wektory losowe \mathbf{K} i \mathbf{X}_R nie są niezależne. Zwracamy

uwagę, że założenie niezależności dotyczy par, tzn. wektory losowe \mathbf{w}_i , $i = 1, \dots, n$, są łącznie niezależne, natomiast dla każdego ustalonego $i \in \{1, \dots, n\}$ zmienne losowe X_i oraz Y_i (będące składowymi wektora \mathbf{w}_i) nie muszą być niezależne; najlepiej, gdy są np. dobrze skorelowane. Za takim podejściem przemawia następujące uzasadnienie: jeżeli obserwowane wartości Y_j i Y_i , gdzie $i \in R^c$, $j \in R$, są bliskie, to przy znacznej korelacji zmiennych X i Y można oczekiwać, że nieobserwowana wartość X_i jest bliska obserwowanej wartości X_j .

TWIERDZENIE 1.2.1. *Imputacyjny estymator wartości oczekiwanej w przypadku imputacji typu hot-deck polegającej na losowaniu respondentów ma postać*

$$\hat{\mu}_{Imp} = \frac{1}{n} \left(\sum_{j \in R} X_j + \sum_{i \in R^c} X_{K_i} \right). \quad (I.2.2)$$

Jego wartość oczekiwana oraz wariancja wynoszą, odpowiednio

$$E \hat{\mu}_{Imp} = \mu \quad \text{oraz} \quad \text{Var} \hat{\mu}_{Imp} = \frac{\sigma^2}{r} \left(1 + \frac{(n-r)(r-1)}{n^2} \right). \quad (I.2.3)$$

Dowód. Zgodnie z (I.2.1) mamy $\bar{\tilde{X}} = \frac{1}{n-r} \sum_{i \in R^c} X_{K_i}$, a zatem wzór (I.0.1) implikuje (I.2.2).

Zauważmy, że dla dowolnego $i \in R^c$

$$E(X_{K_i} | \mathbf{K}) = E X = \mu \quad \text{oraz} \quad \text{Var}(X_{K_i} | \mathbf{K}) = \text{Var}(X) = \sigma^2$$

i w konsekwencji

$$\hat{\mu} = E \bar{\tilde{X}} = E X_{K_i} = E E(X_{K_i} | \mathbf{K}) = \mu \quad (I.2.4)$$

oraz

$$\tilde{\sigma}^2 = \text{Var} X_{K_i} = \text{Var} E(X_{K_i} | \mathbf{K}) + E \text{Var}(X_{K_i} | \mathbf{K}) = \sigma^2. \quad (I.2.5)$$

Ze wzoru (I.2.4) wynika, że $E \bar{\tilde{X}} = \mu$, zgodnie ze wzorem (I.0.2), mamy więc pierwszy ze wzorów (I.2.3), czyli estymator $\hat{\mu}_{Imp}$ jest nieobciążonym estymatorem średniej.

Podobnie dla dowolnych $i, j \in R^c$ oraz $l \in R$ mamy

$$\text{Cov}(X_{K_i}, X_{K_j} | \mathbf{K}) = \sigma^2 \mathbf{I}(K_i = K_j) \quad \text{oraz} \quad \text{Cov}(X_l, X_{K_j} | \mathbf{K}) = \sigma^2 \mathbf{I}(K_j = l).$$

Stąd wynika, że dla dowolnych $i, j \in R^c$

$$\begin{aligned} \tilde{\rho} &= \text{Corr}(\tilde{X}_i, \tilde{X}_j) = \frac{\text{Cov}(X_{K_i}, X_{K_j})}{\sigma^2} = \\ &= \frac{\text{E Cov}(X_{K_i}, X_{K_j} | \mathbf{K}) + \text{Cov}(\text{E}(X_{K_i} | \mathbf{K}), \text{E}(X_{K_j} | \mathbf{K}))}{\sigma^2}. \end{aligned}$$

Ponieważ $\text{Cov}(\text{E}(X_{K_i} | \mathbf{K}), \text{E}(X_{K_j} | \mathbf{K})) = 0$, z powyższego wzoru otrzymujemy

$$\tilde{\rho} = \text{P}(K_i = K_j) = \frac{1}{r}. \quad (\text{I.2.6})$$

Z kolei dla dowolnych $j \in R$ oraz $i \in R^c$

$$\begin{aligned} \rho &= \text{Corr}(X_j, \tilde{X}_i) = \frac{\text{Cov}(X_j, X_{K_i})}{\sigma^2} = \\ &= \frac{\text{E Cov}(X_j, X_{K_i} | \mathbf{K}) + \text{Cov}(\text{E}(X_j | \mathbf{K}), \text{E}(X_{K_i} | \mathbf{K}))}{\sigma^2}. \end{aligned}$$

Ponieważ $\text{Cov}(\text{E}(X_j | \mathbf{K}), \text{E}(X_{K_i} | \mathbf{K})) = 0$, z powyższego wzoru otrzymujemy

$$\rho = \text{P}(K_i = j) = \frac{1}{r}. \quad (\text{I.2.7})$$

Wykorzystując wzory (I.2.4—I.2.7), zgodnie ze wzorem (I.0.3), mamy

$$\text{Var} \hat{\mu}_{imp} = \frac{r\sigma^2 + (n-r)\sigma^2 + (n-r)(n-r-1)\frac{\sigma^2}{r} + 2r(n-r)\frac{\sigma^2}{r}}{n^2}.$$

Po uproszczeniach otrzymujemy drugi ze wzorów (I.2.3). ■

TWIERDZENIE 1.2.2. *Imputacyjny estymator wariancji w przypadku imputacji typu hot-deck polegającej na losowaniu respondentów ma postać*

$$S_{imp}^2 = \frac{1}{n-1} \left((r-1)S_R^2 + \sum_{i \in R^c} \left(X_{K_i} - \frac{1}{n-r} \sum_{j \in R^c} X_{K_j} \right)^2 + \frac{r(n-r)}{n} \left(\bar{X}_R - \frac{1}{n-r} \sum_{j \in R^c} X_{K_j} \right)^2 \right),$$

a jego wartość oczekiwana wynosi

$$E S_{imp}^2 = \frac{r-1}{r} \frac{n(n-1)+r}{n(n-1)} \sigma^2. \quad (I.2.8)$$

Nieobciążony imputacyjny estymator wariancji estymatora $\hat{\mu}_{imp}$ ma postać

$$v^2(\hat{\mu}_{imp}) = \frac{n-1}{n(r-1)} \frac{n(n-1)+r+r(n-r)}{n(n-1)+r} S_{imp}^2. \quad (I.2.9)$$

Dowód. Wzór na S_{imp}^2 jest natychmiastową konsekwencją wzoru (I.2.1) oraz ogólnego wzoru na estymator imputacyjny wariancji (I.0.8). Z kolei wstawiając wzory (I.2.4)—(I.2.7) do wzoru (I.0.9) otrzymujemy

$$\begin{aligned} E S_{imp}^2 &= \frac{r}{n} \sigma^2 + \frac{n-r}{n} \sigma^2 - \frac{(n-r)(n-r-1)}{n(n-1)} \frac{\sigma^2}{r} - 2 \frac{r(n-r)}{n(n-1)} \frac{\sigma^2}{r} = \\ &= \left(1 - \frac{(n-r)(n+r-1)}{n(n-1)r} \right) \sigma^2. \end{aligned}$$

Ostatnie wyrażenie jest równoważne ze wzorem (I.2.8).

Wzór (I.2.9) wynika z porównania wzoru (I.2.8) oraz drugiego ze wzorów (I.2.3).■

Zatem dla dużych wartości r mamy $E S_{imp}^2 \approx \sigma^2$, czyli estymator wariancji jest w przybliżeniu nieobciążony.

1.2.2. Losowanie z rozkładu normalnego

Zakładamy, że próbka pochodzi z rozkładu normalnego, czyli $X_i, i \in R$, są niezależnymi zmiennymi losowymi o rozkładzie normalnym $N(\mu, \sigma^2)$. Zmienne imputowane generowane są w sposób niezależny z rozkładu $N(\bar{X}_R, S_R^2)$, tzn. $\tilde{X}_j | \mathbf{X}_R \sim N(\bar{X}_R, S_R^2), j \in R^c$, oraz $\tilde{X}_j, j \in R^c$, są warunkowo niezależne pod warunkiem \mathbf{X}_R . Zatem

$$E \tilde{X}_j = E E(\tilde{X}_j | \mathbf{X}_R) = E \bar{X}_R = \mu$$

oraz

$$\text{Var } \tilde{X}_j = \text{Var} E(\tilde{X}_j | \mathbf{X}_R) + E \text{Var}(\tilde{X}_j | \mathbf{X}_R) = \text{Var } \bar{X}_R + E S_R^2 = \sigma^2 \frac{1+r}{r}.$$

Ponieważ

$$\text{Cov}(X_i, \tilde{X}_j | \mathbf{X}_R) = 0 \quad \text{oraz} \quad \text{Cov}(\tilde{X}_i, \tilde{X}_j | \mathbf{X}_R) = 0.$$

więc

$$\text{Cov}(X_i, \tilde{X}_j) = \text{Cov}(E(X_i | \mathbf{X}_R), E(\tilde{X}_j | \mathbf{X}_R)) = \text{Cov}(X_i, \bar{X}_R) = \frac{\sigma^2}{r}$$

oraz

$$\text{Cov}(\tilde{X}_i, \tilde{X}_j | \mathbf{X}_R) = \text{Cov}(E(\tilde{X}_i | \mathbf{X}_R), E(\tilde{X}_j | \mathbf{X}_R)) = \text{Var} \bar{X}_R = \frac{\sigma^2}{r}.$$

TWIERDZENIE 1.2.3. *Dla imputacyjnego estymatora wartości oczekiwanej w przypadku imputacji typu hot-deck polegającej na losowaniu z rozkładu normalnego wartość oczekiwana oraz wariancja wynoszą odpowiednio*

$$E\hat{\mu}_{Imp} = \mu \quad \text{oraz} \quad \text{Var} \hat{\mu}_{Imp} = \frac{\sigma^2}{r} \left(1 + \frac{(n-r)r}{n^2} \right). \quad (I.2.10)$$

Natomiast

$$ES_{Imp}^2 = \left(1 - \frac{n-r}{n(n-1)} \right) \sigma^2. \quad (I.2.11)$$

Nieobciążony estymator wariancji estymatora $\hat{\mu}_{Imp}$ ma postać

$$\hat{v}^2(\hat{\mu}_{Imp}) = \frac{n-1}{rn} \frac{n^2 + (n-r)r}{n(n-1) - (n-r)} S_{Imp}^2. \quad (I.2.12)$$

Dowód. Wzory (I.2.10) są odpowiednio natychmiastową konsekwencją wzorów (I.0.2) i (I.0.3), natomiast wzór (I.2.11) wynika wprost ze wzoru (I.0.9). Wzór (I.2.12) jest z kolei konsekwencją wzoru (I.2.11) i drugiego ze wzorów (I.2.10). ■

1.3. Imputacja regresyjna

Zakładamy, że pary (X_i, Y_i) , $i = 1, 2, \dots, n$, tworzące próbkę pierwotną mają jednakowe rozkłady, takie jak para (X, Y) i są niezależne. Próbka obserwowana ma postać (X_i, Y_i) , $i \in R$, Y_j , $j \in R^c$ czyli zmienna Y -owa jest obserwowana

w całej próbkę pierwotnej, natomiast zmienna X -owa — jedynie dla respondentów. Jeśli znany jest rozkład warunkowy $X|Y$, to można wartość zmiennej losowej \tilde{X}_j generować z rozkładu warunkowego $X_j|Y_j, j \in R^c$. Wtedy własności statystyczne próbki imputowanej \tilde{X} nie różnią się od własności próbki pierwotnej X . Jednak najczęściej rozkład warunkowy może być jedynie estymowany na podstawie obserwacji $(X_i, Y_i), i \in R$. W takiej sytuacji precyzyjny opis własności estymatorów imputacyjnych jest trudny, w szczególności zależy od metody estymacji rozkładu warunkowego.

Zamiast estymacji rozkładu warunkowego można wykorzystać model regresyjny, który jest konstrukcją teoretyczną znajdującą zastosowanie w takiej sytuacji.

Niech $EX = \mu, EY = \nu, \text{Var}(X) = \sigma^2, \text{Var}(Y) = \tau^2$. Niech χ będzie współczynnikiem korelacji zmiennych X i Y . Niech λ będzie współczynnikiem regresji X względem Y , tzn. $\lambda = \chi \frac{\sigma}{\tau}$.

W najprostszej sytuacji, gdy znane są średnia ν oraz współczynnik regresji λ , technika imputacji regresyjnej polega na przypisaniu nierespondentom wartości predyktora regresyjnego, tzn.

$$\tilde{X}_j = \lambda(Y_j - \nu) + \bar{X}_R, j \in R^c.$$

TWIERDZENIE 1.3.1. *Imputacyjny estymator wartości oczekiwanej w przypadku imputacji regresyjnej ma postać*

$$\hat{\mu}_{Imp} = \bar{X}_R + \frac{n-r}{n} \lambda (\bar{Y}_{R^c} - \nu), \tag{I.3.1}$$

gdzie $\bar{Y}_{R^c} = \frac{1}{n-r} \sum_{j \in R^c} Y_j$.

Jego wartość oczekiwana oraz wariancja wynoszą odpowiednio

$$E \hat{\mu}_{Imp} = \mu \quad \text{oraz} \quad \text{Var} \hat{\mu}_{Imp} = \sigma^2 \left(\frac{1}{r} + \chi^2 \frac{n-r}{n^2} \right). \tag{I.3.2}$$

Dowód. Ponieważ

$$\tilde{X} = \frac{1}{n-r} \sum_{j \in R^c} \tilde{X}_j = \bar{X}_R + \lambda (\bar{Y}_{R^c} - \nu), \tag{I.3.3}$$

więc wzór (I.3.1) wynika wprost ze wzoru (I.0.1).

Z kolei

$$E(\bar{Y}_{R^c} - \nu) = 0.$$

a zatem pierwszy ze wzorów (I.3.2) wynika natychmiast z (I.3.1).

Ponieważ Y_j oraz \bar{X}_R są niezależne dla $j \in R^c$, to

$$\tilde{\sigma}^2 = \text{Var} \tilde{X}_j = \text{Var}(\lambda(Y_j - \nu)) + \text{Var}(\bar{X}_R) = \sigma^2 \left(\chi^2 + \frac{1}{r} \right).$$

Dla $i \in R, j \in R^c$

$$\text{Cov}(X_i, \tilde{X}_j) = \text{Cov}(X_i, \bar{X}_R) = \frac{\sigma^2}{r},$$

więc

$$\rho = \frac{\frac{\sigma^2}{r}}{\sigma^2 \sqrt{\chi^2 + \frac{1}{r}}} = \frac{1}{\sqrt{r(1+r\chi^2)}}.$$

Dla $i, j \in R^c$

$$\text{Cov}(\tilde{X}_i, \tilde{X}_j) = \text{Var}(\bar{X}_R) = \frac{\sigma^2}{r},$$

więc

$$\tilde{\varrho} = \frac{\frac{\sigma^2}{r}}{\sigma^2 \left(\chi^2 + \frac{1}{r} \right)} = \frac{1}{1+r\chi^2}$$

Ponieważ $\tilde{\varrho} \tilde{\sigma}^2 = \frac{\sigma^2}{r}$ oraz $\rho \sigma \tilde{\sigma} = \frac{\sigma^2}{r}$, więc wstawiając otrzymane wyżej wyrażenia na $\tilde{\sigma}^2$, ρ oraz $\tilde{\varrho}$ do wzoru (I.0.3) mamy

$$\begin{aligned} \text{Var} \hat{\mu}_{mp} &= \frac{r\sigma^2 + (n-r)\sigma^2 \left(\chi^2 + \frac{1}{r} \right) + (n-r)(n-r-1) \frac{\sigma^2}{r} + 2r(n-r) \frac{\sigma^2}{r}}{n^2} = \\ &= \frac{\sigma^2 \left(\frac{n^2}{r} + (n-r)\chi^2 \right)}{n^2}, \end{aligned}$$

co jest równoważne z drugim ze wzorów (I.3.2). ■

Uwaga. W odróżnieniu od wcześniej rozważanych przypadków w imputacji regresyjnej wartości imputowane zależą od współczynnika regresji λ oraz od wartości oczekiwanej ν . Najczęściej nie są one znane. Wtedy we wzorach na zmienne imputowane, na estymator imputacyjny średniej oraz na estymator imputacyjny wariancji (poniżej) zamiast λ oraz ν należy wstawić ich estymatory otrzymane na podstawie obserwowanej części próbki $(X_i, Y_i), i \in R, Y_j, j \in R^c$.

TWIERDZENIE I.3.2. *Imputacyjny estymator wariancji w przypadku imputacji regresyjnej ma postać*

$$S_{imp}^2 = \frac{1}{n-1} \left((r-1)S_R^2 + \lambda^2 \left((n-r-1)S_{R^c, Y}^2 + \frac{(n-r)r}{n} (\bar{Y}_{R^c} - \nu)^2 \right) \right), \quad (I.3.4)$$

a jego wartość oczekiwana wynosi

$$E S_{imp}^2 = \sigma^2 \left(\frac{r-1}{n-1} + \left(1 - \frac{r}{n} \right) \chi^2 \right). \quad (I.3.5)$$

Nieobciążony imputacyjny estymator wariancji estymatora $\hat{\mu}_{imp}$ ma postać

$$\hat{\nu}^2(\hat{\mu}_{imp}) = \frac{\frac{1}{r} + \frac{n-r}{n^2} \chi^2}{\frac{r-1}{n-1} + \frac{n-r}{n} \chi^2} S_{imp}^2. \quad (I.3.6)$$

Dowód. Ze wzoru (I.3.1) i definicji \tilde{S}^2 mamy

$$\tilde{S}^2 = \lambda^2 S_{R^c, Y}^2, \text{ gdzie } S_{R^c, Y}^2 = \frac{1}{n-r-1} \sum_{j \in R^c} (Y_j - \bar{Y}_{R^c})^2.$$

Ponownie wykorzystując wzór (I.3.1) otrzymujemy

$$\left(\bar{X}_R - \bar{\tilde{X}} \right)^2 = \lambda^2 (\bar{Y}_{R^c} - \nu)^2.$$

Wstawiając oba powyższe wyrażenia na \tilde{S}^2 oraz na $\left(\bar{X}_R - \bar{\tilde{X}} \right)^2$ do wzoru (I.0.8) otrzymujemy (I.3.4).

Natomiast wstawiając wyrażenia

$$\tilde{\sigma}^2 = \sigma^2 \left(\chi^2 + \frac{1}{r} \right), \quad \tilde{\varrho} \tilde{\sigma}^2 = \frac{\sigma^2}{r} \quad \text{oraz} \quad \rho \sigma \tilde{\sigma} = \frac{\sigma^2}{r}$$

do wzoru (I.0.9) na wartość oczekiwaną S_{Imp}^2 , po łatwych uproszczeniach, otrzymujemy (I.3.5).

Wzór (I.3.6) wynika z porównania wzoru (I.3.5) oraz drugiego ze wzorów (I.3.2). ■

II. DETERMINISTYCZNY ZBIÓR BRAKÓW OBSERWACJI — IMPUTACJA WIELOKROTNA

II.0. Ogólny schemat imputacji wielokrotnej

Imputacja wielokrotna została wprowadzona przez Rubina w monografii z 1987 r. Poniżej przytaczamy podstawowe zasady tej metody, zwane zasadami Rubina (*Rubin's rules*), przedstawione np. na stronach 39 i 40 monografii Carpenter i Kenward (2013).

Imputacja wielokrotna polega na generowaniu kilku próbek imputacyjnych

$$\tilde{X}^{(l)} = \left(X_i, i \in R, \tilde{X}_j^{(l)}, j \in R^c \right), \quad l = 1, \dots, m,$$

na podstawie obserwowanej części próbki $X_R = (X_i, i \in R)$.

Definicja II.0.1. *Estymatorem wielokrotnej imputacji Rubina parametru θ (odpowiadającym estymatorowi $\hat{\theta} = g(X)$) nazywamy średnią z estymatorów imputacyjnych dla próbek $\tilde{X}^{(l)}$, $l = 1, \dots, m$, czyli*

$$\hat{\theta}_{MImp} = \frac{1}{m} \sum_{l=1}^m \hat{\theta}_{Imp}^{(l)}.$$

Estymator Rubina wariancji estymatora $\hat{\theta}_{MImp}$ ma postać

$$\hat{v}_{MImp}^2 = \bar{U}_m + \frac{m+1}{m} B_m,$$

gdzie

$$\bar{U}_m := \frac{1}{m} \sum_{l=1}^m \hat{v}_{Imp}^{(l),2}$$

oznacza średnią estymatorów wariancji estymatorów imputacyjnych, natomiast

$$B_m := \frac{1}{m-1} \sum_{l=1}^m \left(\hat{\theta}_{Imp}^{(l)} - \hat{\theta}_{MImp} \right)^2 \quad (\text{II.0.1})$$

jest wariancją empiryczną estymatorów imputacyjnych.

Poniżej interesować się będziemy estymacją średniej μ w modelu opisanym w poprzednim rozdziale. Zakładać będziemy, że $\tilde{X}^{(l)}$, $l=1, \dots, m$, są warunkowo niezależne pod warunkiem \mathbf{X}_R . Oczywiście jest to równoważne temu, że imputowane części próbek $(\tilde{X}_j^{(l)}, j \in R^c)$, $l=1, \dots, m$, są warunkowo niezależne pod warunkiem \mathbf{X}_R . Zakłada się również, że $E \tilde{X}_j^{(l)} = \tilde{\mu}$, $\text{Var} \tilde{X}_j^{(l)} = \tilde{\sigma}^2$ nie zależą od $j \in R^c$ oraz od $l=1, \dots, m$, $\text{Corr}(\tilde{X}_i^{(l)}, \tilde{X}_j^{(l)}) = \tilde{\rho}$ nie zależy od $i, j \in R^c$ oraz od $l=1, \dots, m$, a $\text{Corr}(X_i, \tilde{X}_j^{(l)}) = \rho$ nie zależy od $i \in R, j \in R^c$ oraz $l=1, \dots, m$.

Dodatkowo zakłada się warunek liniowości regresji

$$E(\tilde{X}_j^{(l)} | \mathbf{X}_R) = \alpha \bar{X}_R + \beta$$

dla dowolnych $j \in R^c$ oraz $l=1, \dots, m$. Współczynniki α i β można wyliczyć korzystając z zależności

$$\tilde{\mu} = EE(\tilde{X}_j^{(l)} | \mathbf{X}_R) = \alpha \mu + \beta$$

oraz

$$\tilde{\mu}\mu + \rho\sigma\tilde{\sigma} = E\tilde{X}_j^{(l)}X_i = EX_iE(\tilde{X}_j^{(l)} | \mathbf{X}_R) = EX_i(\alpha \bar{X}_R + \beta) = \alpha \left(\frac{\sigma^2}{r} + \mu^2 \right) + \beta\mu.$$

Rozwiązując powyższy układ równań względem α i β otrzymujemy

$$E(\tilde{X}_j^{(l)} | \mathbf{X}_R) = r\rho \frac{\tilde{\sigma}}{\sigma} (\bar{X}_R - \mu) + \tilde{\mu} \quad (\text{II.0.2})$$

dla dowolnych $j \in R^c$ oraz $l=1, \dots, m$.

Na podstawie próbek imputacyjnych konstruowane są estymatory

$$\hat{\mu}_{Imp}^{(l)} = \frac{1}{n} \left[r \bar{X}_R + (n-r) \bar{\tilde{X}}^{(l)} \right], \quad l=1, \dots, m,$$

a estymator finalny jest ich średnią

$$\hat{\mu}_{MImp} = \frac{1}{m} \sum_{l=1}^m \hat{\mu}_{Imp}^{(l)} = \frac{r}{n} \bar{X}_R + \frac{n-r}{nm} \sum_{l=1}^m \tilde{X}^{(l)}. \quad (\text{II.0.3})$$

TWIERDZENIE II.0.1. *Wartość oczekiwana estymatora wielokrotnej imputacji dla średniej wynosi*

$$E\hat{\mu}_{MImp} = \frac{r}{n} \mu + \frac{n-r}{n} \tilde{\mu}. \quad (\text{II.0.4})$$

Jego wariancja ma postać

$$\begin{aligned} \text{Var}(\hat{\mu}_{MImp}) &= \\ &= \frac{mr\sigma^2 + (n-r)(1+(n-r-1)\tilde{\rho})\tilde{\sigma}^2 + 2mr(n-r)\varrho\sigma\tilde{\sigma} + (m-1)(n-r)^2r\varrho^2\tilde{\sigma}^2}{mn^2}. \end{aligned} \quad (\text{II.0.5})$$

Dowód. Wzór (II.0.4) wynika wprost z faktu, że dla każdego z estymatorów $\hat{\mu}_{Imp}^{(l)}$, $l = 1, \dots, m$, obowiązuje wzór (I.0.2).

Znajdziemy teraz wariancję estymatora wielokrotnej imputacji średniej

$$\text{Var} \hat{\mu}_{MImp} = \frac{1}{m^2} \left(\sum_{l=1}^m \text{Var} \hat{\mu}_{Imp}^{(l)} + \sum_{k \neq l}^m \text{Cov} \left(\hat{\mu}_{Imp}^{(k)}, \hat{\mu}_{Imp}^{(l)} \right) \right).$$

Ale zgodnie ze wzorem (I.0.3)

$$\text{Var} \hat{\mu}_{Imp}^{(l)} = \frac{1}{n^2} \left(r\sigma^2 + (n-r)(1+(n-r-1)\tilde{\rho})\tilde{\sigma}^2 + 2r(n-r)\varrho\sigma\tilde{\sigma} \right).$$

Z kolei

$$\begin{aligned} \text{Cov} \left(\hat{\mu}_{Imp}^{(k)}, \hat{\mu}_{Imp}^{(l)} \right) &= \frac{1}{n^2} \left(r^2 \text{Var} \left(\bar{X}_R \right) + (n-r)^2 \text{Cov} \left(\tilde{X}^{(k)}, \tilde{X}^{(l)} \right) + \right. \\ &\quad \left. + r(n-r) \left(\text{Cov} \left(\bar{X}_R, \tilde{X}^{(k)} \right) + \text{Cov} \left(\bar{X}_R, \tilde{X}^{(l)} \right) \right) \right). \end{aligned}$$

Wykorzystując wzór (I.0.6) otrzymujemy

$$\text{Cov} \left(\hat{\mu}_{Imp}^{(k)}, \hat{\mu}_{Imp}^{(l)} \right) = \frac{1}{n^2} \left(r\sigma^2 + (n-r)^2 \text{Cov} \left(\tilde{X}^{(k)}, \tilde{X}^{(l)} \right) + 2r(n-r)\varrho\sigma\tilde{\sigma} \right).$$

Ale

$$\text{Cov}\left(\widetilde{X}^{(k)}, \widetilde{X}^{(l)}\right) = \text{ECov}\left(\widetilde{X}^{(k)}, \widetilde{X}^{(l)} \mid \mathbf{X}_R\right) + \text{Cov}\left(\text{E}\left(\widetilde{X}^{(k)} \mid \mathbf{X}_R\right), \text{E}\left(\widetilde{X}^{(l)} \mid \mathbf{X}_R\right)\right).$$

Z warunkowej niezależności wynika, że pierwszy składnik po prawej stronie jest równy zeru, natomiast

$$\text{Cov}\left(\text{E}\left(\widetilde{X}^{(k)} \mid \mathbf{X}_R\right), \text{E}\left(\widetilde{X}^{(l)} \mid \mathbf{X}_R\right)\right) = \frac{1}{(n-r)^2} \sum_{i,j \in R^c} \text{Cov}\left(\text{E}\left(\widetilde{X}_i^{(k)} \mid \mathbf{X}_R\right), \text{E}\left(\widetilde{X}_j^{(l)} \mid \mathbf{X}_R\right)\right).$$

Z liniowości regresji (II.0.2) otrzymujemy dla dowolnych $i, j \in R^c$

$$\text{Cov}\left(\text{E}\left(\widetilde{X}_i^{(k)} \mid \mathbf{X}_R\right), \text{E}\left(\widetilde{X}_j^{(l)} \mid \mathbf{X}_R\right)\right) = r^2 \varrho^2 \frac{\widetilde{\sigma}^2}{\sigma^2} \text{Var}\left(\overline{X}_R\right) = r\varrho^2 \widetilde{\sigma}^2.$$

Zatem

$$\text{Cov}\left(\hat{\mu}_{Imp}^{(k)}, \hat{\mu}_{Imp}^{(l)}\right) = \frac{1}{n^2} \left(r\sigma^2 + (n-r)^2 r\varrho^2 \widetilde{\sigma}^2 + 2r(n-r)\varrho\sigma\widetilde{\sigma} \right). \quad (\text{II.0.6})$$

Ostatecznie

$$\begin{aligned} \text{Var}\hat{\mu}_{MImp} &= \frac{1}{mn^2} \left(r\sigma^2 + (n-r) \left(1 + (n-r-1)\widetilde{\rho} \right) \widetilde{\sigma}^2 + 2r(n-r)\varrho + \sigma\widetilde{\sigma} \right) \\ &\quad + \frac{m-1}{mn^2} \left(r\sigma^2 + (n-r)^2 r\varrho^2 \widetilde{\sigma}^2 + 2r(n-r)\varrho\sigma\widetilde{\sigma} \right), \end{aligned}$$

co po łatwych uproszczeniach pozwala otrzymać wzór (II.0.5).■

Zgodnie z drugą częścią Definicji II.0.1 estymatorem Rubina wariancji estymatora $\hat{\mu}_{MImp}$ nazywamy statystykę określoną wzorem

$$\hat{v}_{MImp}^2 = \overline{U}_m + \frac{m+1}{m} B_m, \quad (\text{II.0.7})$$

gdzie

$$\overline{U}_m := \frac{1}{mn} \sum_{l=1}^m \left(S_{Imp}^{(l)} \right)^2, \quad (\text{II.0.8})$$

przy czym

$$\left(S_{Imp}^{(l)}\right)^2 = \frac{1}{n-1} \sum_{i=1}^n \left(\tilde{X}_i^{(l)} - \hat{\mu}_{Imp}^{(l)}\right)^2, \quad l=1, \dots, m,$$

natomiast

$$B_m := \frac{1}{m-1} \sum_{l=1}^m \left(\hat{\mu}_{Imp}^{(l)} - \hat{\mu}_{MImp}\right)^2. \quad (\text{II.0.9})$$

Estymator (II.0.7), zwany estymatorem Rubina, na ogół nie jest nieobciążony.

TWIERDZENIE II.0.2. *Wartość oczekiwana estymatora Rubina wynosi*

$$\begin{aligned} E\hat{v}_{MImp}^2 = & \frac{r\sigma^2 + (n-r)\tilde{\sigma}^2 - \frac{(n-r)\left[(n-r-1)\tilde{\rho}\tilde{\sigma}^2 + 2r\rho\sigma\tilde{\sigma} - r(\mu - \tilde{\mu})^2\right]}{n-1}}{n^2} + \\ & + \frac{m+1}{m} \frac{(n-r)\left[1 + (n-r-1)\tilde{\varrho} - (n-r)r\varrho^2\right]}{n^2} \tilde{\sigma}^2. \end{aligned} \quad (\text{II.0.10})$$

Dowód. Zauważmy najpierw, że

$$E\bar{U}_m = \frac{1}{mn} \sum_{l=1}^m E\left(S_{Imp}^{(l)}\right)^2 = \frac{1}{n} E\left(S_{Imp}^{(1)}\right)^2.$$

A zatem zgodnie ze wzorem (I.0.9) otrzymujemy

$$E\bar{U}_m = \frac{1}{n^2} \left(r\sigma^2 + (n-r)\tilde{\sigma}^2 - \frac{(n-r)\left[(n-r-1)\tilde{\rho}\tilde{\sigma}^2 + 2r\rho\sigma\tilde{\sigma} - r(\mu - \tilde{\mu})^2\right]}{n-1} \right), \quad (\text{II.0.11})$$

czyli pierwszy składnik we wzorze (II.0.10).

Z kolei, jeśli $E\hat{\mu}_{Imp}^{(l)} = E\hat{\mu}_{MImp} = v, l=1, \dots, m$, to

$$\begin{aligned} \sum_{l=1}^m E\left(\hat{\mu}_{Imp}^{(l)} - \hat{\mu}_{MImp}\right)^2 &= \sum_{l=1}^m E\left(\hat{\mu}_{Imp}^{(l)} - v\right)^2 - 2E\left(\hat{\mu}_{MImp} - v\right) \sum_{l=1}^m \left(\hat{\mu}_{Imp}^{(l)} - v\right) + \\ &+ mE\left(\hat{\mu}_{MImp} - v\right)^2 = m\left(\text{Var}\hat{\mu}_{Imp}^{(1)} - \text{Var}\hat{\mu}_{MImp}\right). \end{aligned}$$

Ale

$$\begin{aligned} \text{Var } \hat{\mu}_{MImp} &= \frac{1}{m^2} \left(\sum_{l=1}^m \text{Var } \hat{\mu}_{Imp}^{(l)} + \sum_{l,k=1, l \neq k}^m \text{Cov}(\hat{\mu}_{Imp}^{(l)}, \hat{\mu}_{Imp}^{(k)}) \right) = \\ &= \frac{\text{Var } \hat{\mu}_{Imp}^{(1)} + (m-1) \text{Cov}(\hat{\mu}_{Imp}^{(1)}, \hat{\mu}_{Imp}^{(2)})}{m}. \end{aligned}$$

W konsekwencji

$$\begin{aligned} \sum_{l=1}^m \text{E}(\hat{\mu}_{Imp}^{(l)} - \hat{\mu}_{MImp})^2 &= m \left(\text{Var } \hat{\mu}_{Imp}^{(1)} - \frac{\text{Var } \hat{\mu}_{Imp}^{(1)} + (m-1) \text{Cov}(\hat{\mu}_{Imp}^{(1)}, \hat{\mu}_{Imp}^{(2)})}{m} \right) = \\ &= (m-1) (\text{Var } \hat{\mu}_{Imp}^{(1)} + \text{Cov}(\hat{\mu}_{Imp}^{(1)}, \hat{\mu}_{Imp}^{(2)})). \end{aligned}$$

Zatem

$$\text{E} B_m = \frac{1}{m-1} \sum_{l=1}^m \text{E}(\hat{\mu}_{Imp}^{(l)} - \hat{\mu}_{MImp})^2 = \text{Var } \hat{\mu}_{Imp}^{(1)} + \text{Cov}(\hat{\mu}_{Imp}^{(1)}, \hat{\mu}_{Imp}^{(2)}).$$

Stosując w powyższej równości wzory (I.0.3) i (II.0.6) otrzymujemy

$$\text{E} B_m = \frac{(n-r)[1 + (n-r-1)\tilde{\rho} - (n-r)r\tilde{\rho}^2]}{n^2} \tilde{\sigma}^2, \quad (\text{II.0.12})$$

czyli drugą część wzoru (II.0.10), bez mnożnika $\frac{m+1}{m}$. ■

II.1. Imputacja wielokrotna typu hot-deck

II.1.1. Losowanie spośród respondentów

Każdą z próbek imputacyjnych $\tilde{X}^{(l)}$, $l = 1, \dots, m$, tworzymy stosując metodę hot-deck opisaną w podrozdziale I.2. Niech $\mathbf{K}^{(l)} = (K_j^{(l)}, j \in R^c)$ będzie wektorem numerów elementów wylosowanych dla kolejnych nierespondentów przy tworzeniu próbki $\tilde{X}^{(l)}$. Zakładamy, że wektory losowe $\mathbf{K}^{(l)}$, $l = 1, \dots, m$, oraz \mathbf{X}_R są niezależne. Wtedy próbki imputacyjne $\tilde{X}^{(l)}$, $l = 1, \dots, m$, są warunkowo niezależne pod warunkiem \mathbf{X}_R .

Zgodnie z teorią rozwiniętą dla jednokrotnej imputacji hot-deck mamy

$$\tilde{\mu} = \mu, \quad \tilde{\sigma}^2 = \sigma^2, \quad \tilde{\rho} = \rho = \frac{1}{r}. \quad (\text{II.1.1})$$

Co więcej, z niezależności $\mathbf{K}^{(l)}$ i \mathbf{X}_R wynika, że dla $j \in R^c$

$$E(\tilde{X}_j^{(l)} | \mathbf{X}_R) = E(X_{K_j^{(l)}} | \mathbf{X}_R) = \bar{X}_R,$$

czyli $\alpha = 1$ oraz $\beta = 0$ w ogólnym wzorze na liniową regresję z podrozdziału II.0.

TWIERDZENIE II.1.1. *Estymator wielokrotnej imputacji ma postać*

$$\hat{\mu}_{MImp} = \frac{r}{n} \bar{X}_R + \frac{1}{nm} \sum_{l=1}^m \sum_{j \in R^c} X_{K_j^{(l)}}.$$

Jego wartość oczekiwana wynosi

$$E\hat{\mu}_{MImp} = \mu,$$

a jego wariancja ma postać

$$\text{Var} \hat{\mu}_{MImp} = \frac{\sigma^2}{r} \left(1 + \frac{(n-r)(r-1)}{mn^2} \right). \quad (\text{II.1.2})$$

Dowód. Wzór na postać estymatora wielokrotnej imputacji wynika wprost z ogólnej postaci (II.0.3) oraz wzoru (I.2.2) z twierdzenia I.2.1, natomiast wzór na jego wartość oczekiwaną jest konsekwencją wzoru (II.0.4) oraz równości $\tilde{\mu} = \mu$.

Zgodnie ze wzorem (II.0.5) mamy

$$\begin{aligned} \text{Var} \hat{\mu}_{MImp} &= \frac{mr\sigma^2 + (n-r) \left(1 + \frac{n-r-1}{r} \right) \sigma^2 + 2m(n-r)\sigma^2 + (m-1) \frac{(n-r)^2}{r} \sigma^2}{mn^2} = \\ &= \frac{\sigma^2}{n^2} \left(r + 2(n-r) + \frac{n^2}{r} - 2n + r \right) + \frac{\sigma^2}{mn^2} \left(n-r + \frac{(n-r)^2}{r} - \frac{n-r}{r} - \frac{(n-r)^2}{r} \right) = \\ &= \frac{\sigma^2}{r} + \frac{\sigma^2(n-r)(r-1)}{mn^2r}, \end{aligned}$$

co daje wzór (II.1.2). ■

TWIERDZENIE II.1.2. *Obciążenie estymatora Rubina wariancji \hat{v}_{MImp}^2 wynosi*

$$B\hat{v}_{MImp}^2 = E\hat{v}_{MImp}^2 - \text{Var}\hat{\mu}_{MImp} = -\frac{(n-r)[n(n-r+1)+2(r-1)]}{n^2(n-1)r}\sigma^2. \quad (\text{II.1.3})$$

Estymator nieobciążony wariancji estymatora $\hat{\mu}_{MImp}$ ma postać

$$\hat{v}_*^2 = \frac{n}{r}\bar{U}_m + \left(\frac{1}{m} + \frac{n(n+r-1)}{(n-1)r(r-1)}\right)B_m.$$

Dowód. Wstawiając do (II.0.11) wartości podane w (II.1.1) otrzymujemy

$$E\bar{U}_m = \frac{\sigma^2}{n^2}\left(n - \frac{(n-r)(n+r-1)}{r(n-1)}\right).$$

Z kolei wstawiając do (II.0.12) te same wartości z (II.1.1) otrzymujemy

$$EB_m = \sigma^2 \frac{(n-1)(r-1)}{rn^2}. \quad (\text{II.1.4})$$

Zatem

$$E\hat{v}_{MImp}^2 = E\bar{U}_m + EB_m + \frac{1}{m}EB_m = \frac{\sigma^2(r-1)(2n(n-1)+2r-nr)}{rn^2(n-1)} + \frac{\sigma^2(r-1)(n-r)}{mn^2r}.$$

Obciążenie wynosi więc

$$E\hat{v}_{MImp}^2 - \text{Var}\hat{\mu}_{MImp} = \frac{\sigma^2(r-1)(2n(n-1)+2r-nr)}{rn^2(n-1)} - \frac{\sigma^2}{r},$$

co daje po uproszczeniach wzór (II.1.3).

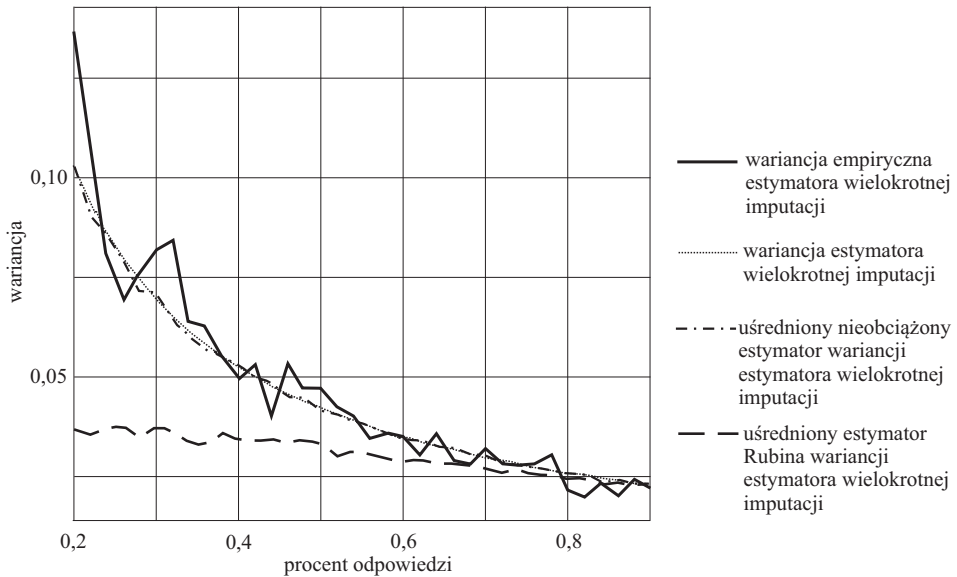
Zauważmy, że

$$\frac{n}{r}E\bar{U}_m + \frac{n(n+r-1)}{(n-1)r(r-1)}EB_m = \frac{\sigma^2}{r}.$$

Nieobciążoność estymatora \hat{v}_*^2 wynika z powyższego wzoru oraz (II.1.4) i wzoru (II.1.2) na wariancję estymatora $\hat{\mu}_{MImp}$. ■

Symulacyjna ilustracja wyników zawartych w twierdzeniach II.1.1 i II.1.2 przedstawiona jest na wyk. 1. Szczegółowy opis procedury symulacyjnej znajduje się w aneksie.

Wykr. 1. WYNIKI ESTYMACJI WARIANCJI ESTYMATORA ŚREDNIEJ
DLA WIELOKROTNEJ IMPUTACJI (HOT-DECK)



Źródło: obliczenia własne (zobacz aneks).

Wniosek II.1.3. *Asymptotycznie względne obciążenie estymatora Rubina \hat{v}_{MImp}^2 wynosi*

$$\lim_{n \rightarrow \infty} \frac{E\hat{v}_{MImp}^2 - \text{Var}\hat{\mu}_{MImp}}{\text{Var}\hat{\mu}_{MImp}} = -1. \quad (\text{II.1.5})$$

Dowód. Ze wzoru (II.1.2) mamy

$$\lim_{n \rightarrow \infty} \text{Var}\hat{\mu}_{MImp} = \frac{\sigma^2}{r}.$$

Z kolei wzór (II.1.3) implikuje

$$\lim_{n \rightarrow \infty} (E\hat{v}_{MImp}^2 - \text{Var}\hat{\mu}_{MImp}) = -\frac{\sigma^2}{r}. \blacksquare$$

Uwaga. Oczywiście istnieją też inne niż \hat{v}_*^2 nieobciążone estymatory wariancji estymatora wielokrotnej imputacji $\hat{\mu}_{MImp}$. W szczególności estymatory mające postać kombinacji liniowej statystyk \bar{U}_m i B_m . Analiza tej klasy estymatorów jest równoważna rozważaniu estymatorów wariancji $\hat{\mu}_{MImp}$ postaci

$$\hat{v}_{(\alpha,\beta)}^2 = \alpha \bar{U}_m + \beta B_m + \frac{1}{m} B_m,$$

przy warunku

$$\beta = \frac{1}{n-r} \left(\frac{n^2}{r-1} - \frac{n(n-1)+r}{(n-r)(n-1)} \alpha \right) \quad (\text{II.1.6})$$

zapewniającym nieobciążoność. Istotnym zagadnieniem jest znalezienie optymalnych wartości α , β , tzn. takich, przy których wariancja estymatora $\hat{v}_{(\alpha,\beta)}^2$ osiąga minimum. Rozwiązanie analityczne jest trudne, wymaga znajomości trzecich i czwartych momentów. Możliwe jest natomiast numeryczne badanie wariancji estymatora $\hat{v}_{(\alpha,\beta)}^2$ jako funkcji zmiennej α , przy ograniczeniu (II.1.6).

II.1.2. Losowanie z rozkładu normalnego

Korzystamy ze wzorów ogólnych, przy czym

$$\tilde{\mu} = \mu, \quad \tilde{\sigma}^2 = \sigma^2 \frac{1+r}{r}, \quad \varrho = \frac{1}{\sqrt{r(r+1)}}, \quad \tilde{\varrho} = \frac{1}{r+1}.$$

TWIERDZENIE II.1.4. *Wartość oczekiwana estymatora imputacji wielokrotnej przy losowaniu z rozkładu normalnego wynosi*

$$E \hat{\mu}_{MImp} = \mu,$$

a jego wariancja ma postać

$$\text{Var} \hat{\mu}_{MImp} = \frac{\sigma^2}{r} \left(1 + \frac{(n-r)r}{mn^2} \right). \quad (\text{II.1.7})$$

Dowód. Powyższe wzory wynikają wprost ze wzorów (II.0.4) i (II.0.5) z twierdzenia II.1.1. ■

TWIERDZENIE II.1.5. *Obciążenie estymatora Rubina wariancji \hat{v}_{MImp}^2 wynosi*

$$B\hat{v}_{MImp}^2 = E\hat{v}_{MImp}^2 - \text{Var}\hat{\mu}_{MImp} = -\frac{(n-r)[(n-1)(n-r)-1]}{n^2(n-1)r}\sigma^2. \quad (\text{II.1.8})$$

Estymator nieobciążony wariancji estymatora $\hat{\mu}_{MImp}$ ma postać

$$\hat{v}_*^2 = \frac{n}{r}\bar{U}_m + \left(\frac{1}{m} + \frac{n}{(n-1)r}\right)B_m.$$

Dowód. Ze wzorów (II.0.11) oraz (II.0.12) wynika, że

$$E\bar{U}_m = \frac{\sigma^2}{n^2}\left(n - \frac{n-r}{n-1}\right) \quad \text{oraz} \quad EB_m = \frac{\sigma^2(n-r)}{n^2}.$$

Zatem

$$E\hat{v}_{MImp}^2 - \text{Var}\hat{\mu}_{MImp} = \frac{\sigma^2}{n^2}\left(n + \frac{(n-r)(r-2)}{r(n-1)}\right) + \frac{m+1}{m}\frac{\sigma^2(n-r)}{n^2} - \frac{\sigma^2}{r}\left(1 + \frac{(n-r)r}{mn^2}\right)$$

skąd po prostych przekształceniach otrzymujemy wzór (II.1.8).

Ponieważ

$$\text{Var}\hat{\mu}_{MImp} - E\frac{1}{m}B_m = \frac{\sigma^2}{r},$$

więc wystarczy znaleźć α i β takie, że $\alpha E\bar{U}_m + \beta EB_m = \frac{\sigma^2}{r}$. Przyjmujemy

$$\alpha = \frac{n}{r}.$$

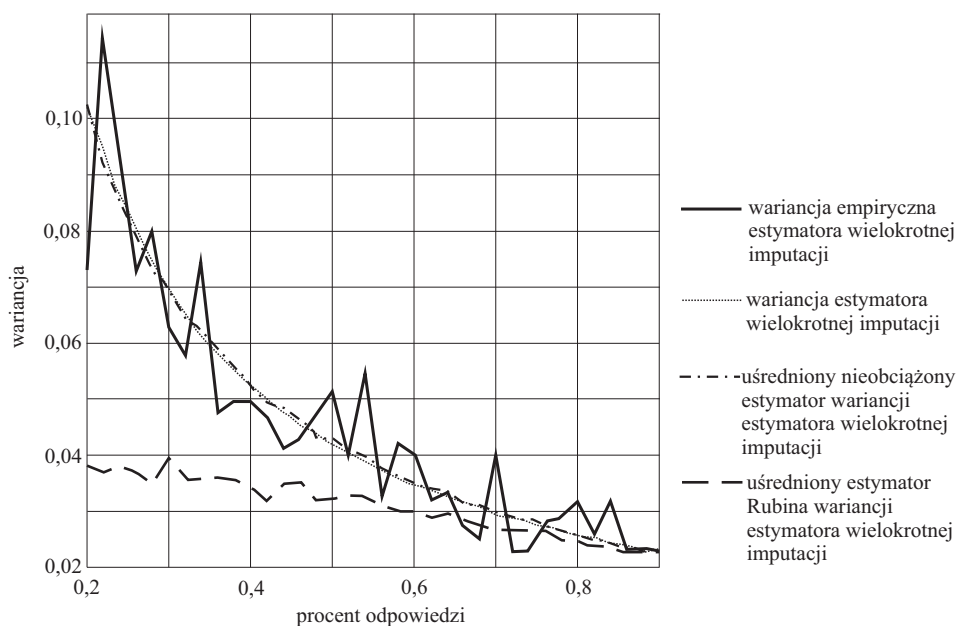
Wtedy nietrudno wyliczyć, że $\beta = -\frac{n(r-2)}{(n-1)r^2}$. ■

Podobnie jak we wniosku II.1.3 i w tym przypadku asymptotyczne obciążenie wynosi

$$\lim_{n \rightarrow \infty} \frac{E\hat{v}_{MImp}^2 - \text{Var}\hat{\mu}_{MImp}}{\text{Var}\hat{\mu}_{MImp}} = -1.$$

Symulacyjna ilustracja wyników zawartych w twierdzeniach II.1.4 i II.1.5 przedstawiona jest na wykr. 2. Szczegółowy opis procedury symulacyjnej znajduje się w aneksie.

**Wykr. 2. WYNIKI ESTYMACJI WARIANCJI ESTYMATORA ŚREDNIEJ
DLA WIELOKROTNEJ IMPUTACJI (HOT-DECK) Z ROZKŁADU NORMALNEGO**



Źródło: jak przy wykr. 1.

III. LOSOWY ZBIÓR BRAKÓW OBSERWACJI — IMPUTACJA JEDNOKROTNA

III.0. Ogólny schemat imputacyjny

Niech wektor losowy $\mathbf{X} = (X_1, \dots, X_n)$ oznacza pierwotną próbkę, natomiast $\mathbf{J} = (J_1, \dots, J_n)$ niech oznacza wektor losowy, w którym składowa $J_i = 1$, jeśli zmienna X_i jest obserwowana oraz $J_i = 0$, jeśli zmienna X_i nie jest obserwowana. W konsekwencji zbiór $\mathbf{R} = \{i \in \{1, \dots, n\} : J_i = 1\}$ jednostek, dla których obserwacja jest dostępna jest losowy.

Przy założeniu, że mechanizm pojawiania się braków odpowiedzi nie jest zależny od wartości badanych zmiennych (MCAR, *missing completely at random*), czyli że wektory losowe \mathbf{X} i \mathbf{J} są (statystycznie) niezależne, do analizy własności rozkładu warunkowego \mathbf{X} pod warunkiem \mathbf{J} można wygodnie stosować metodologię rozwiniętą w rozdziałach I i II, czyli w sytuacji, gdy zbiór \mathbf{R} jest deterministyczny. Ostateczne wyniki otrzymuje się poprzez uśrednianie

rezultatów dla rozkładu warunkowego względem rozkładu wektora \mathbf{J} . Szczegółowo jest to opisane dla imputacji jednokrotnej w bieżącym rozdziale i dla imputacji wielokrotnej w rozdziale kolejnym.

Konstrukcja wektora \mathbf{J} wymaga przeanalizowania pewnej subtelności technicznej. Losowość braków obserwacji modelowana jest pierwotnie przez zmienne losowe (indykatory odpowiedzi) $\delta_1, \dots, \delta_n$, przy czym $\delta_i = 1$, jeśli obserwacja X_i jest dostępna, a $\delta_i = 0$, jeśli obserwacja X_i nie jest dostępna (brak odpowiedzi). Zakładamy, że zmienne $\delta_1, \dots, \delta_n$ są niezależne oraz że wektory \mathbf{X} i $\boldsymbol{\delta} = (\delta_1, \dots, \delta_n)$ są również niezależne. Wprowadzamy oznaczenie p_i na prawdopodobieństwo i -tej odpowiedzi, $P(\delta_i = 1) = p_i$, dla każdego $i = 1, \dots, n$. Niech

$|\boldsymbol{\delta}| = \sum_{i=1}^n \delta_i$ oznacza liczbę elementów obserwowanych w próbie. Jeśli w danej

próbie nie ma żadnej dostępnej obserwacji, tzn. $|\boldsymbol{\delta}| = 0$, to badamy kolejną próbkę i postępujemy tak do momentu, gdy znajdziemy próbkę, dla której $|\boldsymbol{\delta}| > 0$.

Dopiero na zakończenie takiej procedury wprowadzamy wektor ostatecznych braków odpowiedzi \mathbf{J} — jeśli obserwacja X_i jest dostępna, przyjmujemy $J_i = 1$, a jeśli obserwacja X_i nie jest dostępna, przyjmujemy $J_i = 0$, $i = 1, \dots, n$. Czyli zamiast oryginalnego wektora $\boldsymbol{\delta}$ o składowych niezależnych, braki odpowiedzi opisane są za pomocą wektora $\mathbf{J} = (J_1, \dots, J_n)$, którego składowe nie są niezależne. Tym niemniej wektory losowe \mathbf{X} i \mathbf{J} są niezależne. Zauważmy, że

$$P(J_1 = \varepsilon_1, \dots, J_n = \varepsilon_n) = P(\delta_1 = \varepsilon_1, \dots, \delta_n = \varepsilon_n \mid |\boldsymbol{\delta}| > 0) = \frac{\prod_{i=1}^n p_i^{\varepsilon_i} (1-p_i)^{1-\varepsilon_i}}{1 - \prod_{i=1}^n (1-p_i)},$$

dla $\varepsilon_i \in \{0,1\}$, $i = 1, \dots, n$, takich, że $\sum_{i=1}^n \varepsilon_i > 0$. Wtedy $|\mathbf{J}| = |\boldsymbol{\delta}| > 0$ oraz

$$P(|\mathbf{J}| = k) = \frac{\sum_{1 \leq i_1, \dots, i_k \leq n} \prod_{l=1}^k p_{i_l} \prod_{j \notin \{i_1, \dots, i_k\}} (1-p_j)}{1 - \prod_{i=1}^n (1-p_i)}, \quad k = 1, \dots, n.$$

W szczególnym przypadku, gdy prawdopodobieństwo odpowiedzi jest stałe, tzn. $p_i = p$, $i = 1, \dots, n$,

$$P(J_1 = \varepsilon_1, \dots, J_n = \varepsilon_n) = \frac{p^{\sum_{i=1}^n \varepsilon_i} (1-p)^{n-\sum_{i=1}^n \varepsilon_i}}{1 - (1-p)^n}.$$

Wtedy zmienna losowa $|\delta|$ ma rozkład dwumianowy, $b(n, p)$, tzn.

$$P(|\delta| = k) = \binom{n}{k} p^k (1-p)^{n-k}, \quad k = 0, 1, \dots, n,$$

a zmienna losowa $|\mathbf{J}|$ ma rozkład dwumianowy ucięty w zerze, $b_+(n, p)$, tzn.

$$P(|\mathbf{J}| = k) = \binom{n}{k} \frac{p^k (1-p)^{n-k}}{q_n}, \quad k = 1, \dots, n,$$

gdzie

$$q_n = 1 - (1-p)^n.$$

Gdy prawdopodobieństwo braków odpowiedzi jest jednakowe i wynosi p , wtedy

$$E g(|\mathbf{J}|) = \frac{1}{q_n} \sum_{i=1}^n g(i) \binom{n}{i} p^i (1-p)^{n-i},$$

w szczególności

$$E|\mathbf{J}| = \frac{np}{q_n} \quad \text{oraz} \quad E \frac{1}{|\mathbf{J}|} = \frac{1}{q_n} \sum_{i=1}^n \frac{1}{i} \binom{n}{i} p^i (1-p)^{n-i}.$$

Uwaga. W pracy Szablowski, Wesołowski i Wieczorkowski (1996) pokazano, że $E \frac{1}{|\mathbf{J}|} \cong \frac{1}{np}$. Dokładniejsze przybliżenie można znaleźć np. w pracach Marciniak i Wesołowski (1999) albo Rempała (2004).

Zmienne imputacyjne oznaczamy symbolem \tilde{X}_i , $i = 1, \dots, n$. Estymator imputacyjny średniej ma postać

$$\hat{\mu}_{Imp} = \frac{1}{n} \sum_{i=1}^n (X_i J_i + \tilde{X}_i (1 - J_i)) = \frac{1}{n} \left(\sum_{j \in R} X_j + \sum_{j \in R^c} \tilde{X}_j \right). \quad (\text{III.0.1})$$

Natomiast imputacyjny estymator wariancji ma postać

$$S_{Imp}^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i J_i + \tilde{X}_i (1 - J_i) - \hat{\mu}_{Imp})^2. \quad (\text{III.0.2})$$

III.1. Imputacja średnią

Technika ta polega na przypisaniu każdemu nieobserwowanemu elementowi próbki średniej z obserwowanej części próbki, czyli

$$\tilde{X}_i = \frac{1}{|\mathbf{J}|} \sum_{j=1}^n X_j J_j =: \bar{X}_J. \quad (\text{III.1.1})$$

TWIERDZENIE III.1.1. *Imputacyjny estymator wartości oczekiwanej w przypadku imputacji średnią ma postać*

$$\hat{\mu}_{Imp} = \bar{X}_J. \quad (\text{III.1.2})$$

Jego wartość oczekiwana oraz wariancja wynoszą odpowiednio

$$E\hat{\mu}_{Imp} = \mu \quad \text{oraz} \quad \text{Var}\hat{\mu}_{Imp} = \sigma^2 E \frac{1}{|\mathbf{J}|}. \quad (\text{III.1.3})$$

Dowód. Ze wzorów (III.0.1) i (III.0.2) wynika, że

$$\hat{\mu}_{Imp} = \frac{1}{n} \sum_{i=1}^n (J_i X_i + (1 - J_i) \bar{X}_J) = \bar{X}_J + \frac{1}{n} \sum_{i=1}^n J_i X_i - \frac{1}{n} \bar{X}_J \sum_{i=1}^n J_i,$$

co prowadzi bezpośrednio do wzoru (III.1.2).

Z pierwszego ze wzorów (I.1.3) mamy

$$E(\hat{\mu}_{Imp} | \mathbf{J}) = \mu,$$

a zatem

$$E\hat{\mu}_{Imp} = EE(\hat{\mu}_{Imp} | \mathbf{J}) = \mu.$$

Z kolei

$$\text{Var}\hat{\mu}_{Imp} = \text{Var}E(\hat{\mu}_{Imp} | \mathbf{J}) + E\text{Var}(\hat{\mu}_{Imp} | \mathbf{J}) = E\text{Var}(\hat{\mu}_{Imp} | \mathbf{J}),$$

ponieważ $E(\hat{\mu}_{Imp} | \mathbf{J})$ nie jest losowa. Z drugiego ze wzorów (I.1.3) mamy

$$\text{Var}(\hat{\mu}_{Imp} | \mathbf{J}) = \frac{\sigma^2}{|\mathbf{J}|}.$$

Czyli zachodzi drugi ze wzorów (III.1.3). ■

TWIERDZENIE III.1.2. *Imputacyjny estymator wariancji w przypadku imputacji średnią ma postać*

$$S_{imp}^2 = \frac{1}{n-1} \sum_{i=1}^n J_i (X_i - \bar{X}_J)^2, \quad (III.1.4)$$

a jego wartość oczekiwana wynosi

$$ES_{imp}^2 = \frac{E|\mathbf{J}|-1}{n-1} \sigma^2. \quad (III.1.5)$$

Nieobciążony imputacyjny estymator wariancji estymatora $\hat{\mu}_{imp}$ ma postać

$$v^2(\hat{\mu}_{imp}) = \frac{(n-1)E \frac{1}{|\mathbf{J}|}}{E|\mathbf{J}|-1} S_{imp}^2. \quad (III.1.6)$$

W przypadku $p_i = p, i = 1, \dots, n$

$$ES_{imp}^2 = \frac{np - q_n}{q_n(n-1)} \sigma^2 \quad \text{oraz} \quad v^2(\hat{\mu}_{imp}) = \frac{q_n(n-1)E \frac{1}{|\mathbf{J}|}}{np - q_n} S_{imp}^2. \quad (III.1.7)$$

Dowód. Zgodnie ze wzorami (III.0.2) i (III.1.2) mamy

$$S_{imp}^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i J_i + \bar{X}_J (1 - J_i) - \bar{X}_J)^2 = \frac{1}{n-1} \sum_{i=1}^n J_i (X_i - \bar{X}_J)^2.$$

Zgodnie ze wzorem (I.1.5)

$$E(S_{imp}^2 | \mathbf{J}) = \frac{|\mathbf{J}|-1}{n-1} \sigma^2.$$

Zatem równość $ES_{imp}^2 = EE(S_{imp}^2 | \mathbf{J})$ pociąga wzór (III.1.5). ■

Zauważmy, że ze wzoru (III.1.5) wynika, iż estymator S_{imp}^2 jest obciążonym estymatorem wariancji.

III.2. Imputacja typu hot-deck

III.2.1. Losowanie spośród respondentów

Niech K_i będzie numerem elementu wylosowanym dla i -tej jednostki, gdy $i \notin R$. Wtedy zmienne losowe K_i $i \notin R$, są warunkowo niezależne pod warunkiem \mathbf{J} oraz $P(K_i = j | \mathbf{J}) = \frac{1}{|\mathbf{J}|}$, $j \in R$. Co więcej, wektory losowe $\mathbf{K}_{R^c} = (K_i, i \in R^c)$

i $\mathbf{X}_R = (X_j, j \in R)$ są warunkowo niezależne pod warunkiem \mathbf{J} .

W konsekwencji

$$\tilde{X}_i = X_{K_i}, \quad i \in \{1, \dots, n\}. \quad (\text{III.2.1})$$

TWIERDZENIE III.2.1. *Imputacyjny estymator wartości oczekiwanej w przypadku imputacji typu hot-deck polegającej na losowaniu respondentów ma postać*

$$\hat{\mu}_{Imp} = \frac{1}{n} \left(\sum_{j \in R} X_j + \sum_{i \in R^c} X_{K_i} \right). \quad (\text{III.2.2})$$

Jego wartość oczekiwana oraz wariancja wynoszą odpowiednio

$$E\hat{\mu}_{Imp} = \mu \quad \text{oraz} \quad \text{Var}\hat{\mu}_{Imp} = \frac{\sigma^2}{n} \left(1 + \frac{1}{n} + (n-1)E\left[\frac{1}{|\mathbf{J}|} - \frac{E|\mathbf{J}|}{n}\right] \right). \quad (\text{III.2.3})$$

W przypadku $p_i = p$, $i = 1, \dots, n$

$$\text{Var}\hat{\mu}_{Imp} = \frac{\sigma^2}{n} \left(1 + \frac{1}{n} + (n-1)E\left[\frac{1}{|\mathbf{J}|} - \frac{p}{q_n}\right] \right). \quad (\text{III.2.4})$$

Dowód. Z twierdzenia I.2.1 wynika, że $E(\hat{\mu}_{Imp} | \mathbf{J}) = \mu$. W konsekwencji, warunkując względem \mathbf{J} otrzymujemy pierwszy ze wzorów (III.2.3).

Ze wzoru (I.2.3) otrzymujemy

$$\text{Var}(\hat{\mu}_{Imp} | \mathbf{J}) = \frac{\sigma^2}{n} \left(1 + \frac{1}{n} + \frac{n-1}{|\mathbf{J}|} - \frac{|\mathbf{J}|}{n} \right).$$

Ponieważ warunkowa wartość oczekiwana $E(\hat{\mu}_{Imp}|\mathbf{J}) = \mu$ jest nielosowa, więc

$$\text{Var} \hat{\mu}_{Imp} = E \text{Var}(\hat{\mu}_{Imp}|\mathbf{J}),$$

skąd wynika drugi ze wzorów (III.2.3). Wzór (III.2.4) jest natychmiastową konsekwencją (III.2.3) oraz wzoru na wartość oczekiwaną w rozkładzie $b_+(n, p)$. ■

TWIERDZENIE III.2.2. *Imputacyjny estymator wariancji w przypadku imputacji typu hot-deck polegającej na losowaniu respondentów ma postać*

$$S_{Imp}^2 = \frac{(|\mathbf{J}|-1)S_R^2 + (n-|\mathbf{J}|-1)S_{R^c}^2 + \frac{|\mathbf{J}|(n-|\mathbf{J}|)}{n}(\bar{X}_R - \bar{X}_{R^c})^2}{n-1},$$

gdzie $S_{R^c}^2 = \frac{1}{n-|\mathbf{J}|-1} \sum_{i \in R^c} (X_{K_i} - \bar{X}_{R^c})^2$ oraz $\bar{X}_{R^c} = \frac{1}{n-|\mathbf{J}|} \sum_{j \in R^c} X_{K_j}$.

Jego wartość oczekiwana wynosi

$$E S_{Imp}^2 = \left(1 + \frac{E|\mathbf{J}|-1}{n(n-1)} - E \frac{1}{|\mathbf{J}|} \right) \sigma^2. \quad (\text{III.2.5})$$

Nieobciążony imputacyjny estymator wariancji estymatora $\hat{\mu}_{Imp}$ ma postać

$$v^2(\hat{\mu}_{Imp}) = \frac{n+1+n(n-1)E \frac{1}{|\mathbf{J}|} - E|\mathbf{J}|}{n(n-1)-1+E|\mathbf{J}|-n(n-1)E \frac{1}{|\mathbf{J}|}} \frac{n-1}{n} S_{Imp}^2. \quad (\text{III.2.6})$$

W przypadku gdy $p_i = p$, $i = 1, \dots, n$

$$E S_{Imp}^2 = \left(1 - \frac{1}{n(n-1)} + \frac{p}{q_n(n-1)} - E \frac{1}{|\mathbf{J}|} \right) \sigma^2$$

oraz $v^2(\hat{\mu}_{Imp}) = \frac{n+1+n(n-1)E \frac{1}{|\mathbf{J}|} - \frac{np}{q_n}}{n(n-1)-1 + \frac{np}{q_n} - n(n-1)E \frac{1}{|\mathbf{J}|}} \frac{n-1}{n} S_{Imp}^2.$

Dowód. Postać estymatora S_{Imp}^2 wynika wprost z postaci podanej w twierdzeniu I.2.2, natomiast ze wzoru (I.2.8) mamy

$$E(S_{Imp}^2 | \mathbf{J}) = \left(1 - \frac{1}{n(n-1)} + \frac{|\mathbf{J}|}{n(n-1)} - \frac{1}{|\mathbf{J}|} \right) \sigma^2.$$

Więc (III.2.5) wynika z równości $ES_{Imp}^2 = EE(S_{Imp}^2 | \mathbf{J})$. Łącząc wzory (III.2.4) i (III.2.5) otrzymujemy postać (III.2.6) nieobciążonego estymatora wariancji estymatora $\hat{\mu}_{Imp}$. ■

III.2.2. Losowanie z rozkładu normalnego

Stosując wykorzystywaną już wcześniej technikę warunkowania przez obserwowaną część próbki X_R bezpośrednio z twierdzenia I.2.3 otrzymujemy pełny opis standardowych własności imputacyjnego estymatora w tym przypadku.

TWIERDZENIE III.2.3. *Dla imputacyjnego estymatora wartości oczekiwanej w przypadku imputacji hot-deck polegającej na losowaniu z rozkładu normalnego wartość oczekiwana oraz wariancja wynoszą odpowiednio*

$$E\hat{\mu}_{Imp} = \mu \quad \text{oraz} \quad \text{Var}\hat{\mu}_{Imp} = \sigma^2 \left(E \frac{1}{|\mathbf{J}|} + \frac{n - E|\mathbf{J}|}{n^2} \right). \quad (\text{III.2.7})$$

Natomiast

$$ES_{Imp}^2 = \left(1 - \frac{n - E|\mathbf{J}|}{n(n-1)} \right) \sigma^2. \quad (\text{III.2.8})$$

Nieobciążony estymator wariancji estymatora $\hat{\mu}_{Imp}$ ma postać

$$v^2(\hat{\mu}_{Imp}) = \frac{n-1}{n} \frac{n^2 E \frac{1}{|\mathbf{J}|} + n - E|\mathbf{J}|}{n(n-1) - n + E|\mathbf{J}|} S_{Imp}^2. \quad (\text{III.2.9})$$

W przypadku gdy $p_i = p$, $i = 1, \dots, n$

$$\text{Var}\hat{\mu}_{Imp} = \sigma^2 \left(E \frac{1}{|\mathbf{J}|} + \frac{q_n - p}{n} \right) \quad \text{oraz} \quad ES_{Imp}^2 = \left(1 - \frac{q_n - p}{(n-1)q_n} \right) \sigma^2.$$

Ponadto nieobciążony estymator wariancji estymatora $\hat{\mu}_{Imp}$ ma postać

$$v^2(\hat{\mu}_{Imp}) = \frac{n-1}{n} \frac{\left(nE \frac{1}{|\mathbf{J}|} + 1 \right) q_n - p}{(n-2)q_n + p} S_{Imp}^2.$$

Dowód. Wzór (III.2.7) jest bezpośrednią konsekwencją wzoru (I.2.10) oraz tego, że $E(\hat{\mu}_{Imp}|\mathbf{J}) = \mu$, natomiast wzór (III.2.8) wynika wprost ze wzoru (I.2.11) przez warunkowanie względem \mathbf{J} . Wzór (III.2.9) na nieobciążony estymator wariancji powstaje z kombinacji wzorów (III.2.7) i (III.2.8). ■

III.3. Imputacja regresyjna

W modelu z podrozdziału I.3 zakładamy, że odpowiedzi i braki odpowiedzi opisane są wektorem \mathbf{J} .

TWIERDZENIE III.3.1. *Imputacyjny estymator wartości oczekiwanej w przypadku imputacji regresyjnej ma postać*

$$\hat{\mu}_{Imp} = \bar{X}_R + \frac{n-|\mathbf{J}|}{n} \lambda (\bar{Y}_{R^c} - v), \tag{III.3.1}$$

gdzie $\bar{Y}_{R^c} = \frac{1}{n-|\mathbf{J}|} \sum_{i \in R^c} Y_i$.

Jego wartość oczekiwana oraz wariancja wynoszą odpowiednio

$$E\hat{\mu}_{Imp} = \mu \text{ oraz } \text{Var}\hat{\mu}_{Imp} = \left(E \frac{1}{|\mathbf{J}|} + \chi^2 \frac{n-E|\mathbf{J}|}{n^2} \right) \sigma^2. \tag{III.3.2}$$

W przypadku $p_i = p, i = 1, \dots, n$

$$\text{Var}\hat{\mu}_{Imp} = \left(E \frac{1}{|\mathbf{J}|} + \chi^2 \frac{q_n - p}{nq_n} \right) \sigma^2.$$

Dowód. Wzór (III.3.1) wynika wprost ze wzoru (I.3.1).

Z pierwszego ze wzorów (I.3.2) mamy $E(\hat{\mu}_{Imp}|\mathbf{J}) = \mu$, więc pierwszy ze wzorów (III.3.2) jest konsekwencją identyczności $E\hat{\mu}_{Imp} = EE(\hat{\mu}_{Imp}|\mathbf{J})$. Ponieważ $\text{Var}E(\hat{\mu}_{Imp}|\mathbf{J}) = 0$, więc

$$\text{Var}\hat{\mu}_{Imp} = E\text{Var}(\hat{\mu}_{Imp}|\mathbf{J}),$$

a z drugiego ze wzorów (I.3.2) mamy

$$\text{Var}(\hat{\mu}_{Imp}|\mathbf{J}) = \sigma^2 \left(\frac{1}{|\mathbf{J}|} + \chi^2 \frac{n-|\mathbf{J}|}{n^2} \right).$$

Stąd wynika drugi ze wzorów (III.3.2). ■

TWIERDZENIE III.3.2. *Imputacyjny estymator wariancji w przypadku imputacji regresyjnej ma postać*

$$S_{Imp}^2 = \frac{(|\mathbf{J}|-1)S_R^2 + \lambda^2 \left((n-|\mathbf{J}|-1)S_{R^c,Y}^2 + \frac{(n-|\mathbf{J}|)|\mathbf{J}|}{n} (\bar{Y}_{R^c} - v)^2 \right)}{n-1},$$

gdzie $S_{R^c,Y}^2 = \frac{1}{n-|\mathbf{J}|-1} \sum_{i \in R^c} (Y_i - \bar{Y}_{R^c})^2$.

Jego wartość oczekiwana wynosi

$$ES_{Imp}^2 = \left(\chi^2 - \frac{1}{n-1} + E|\mathbf{J}| \left(\frac{1}{n-1} - \frac{\chi^2}{n} \right) \right) \sigma^2. \quad (\text{III.3.3})$$

Nieobciążony imputacyjny estymator wariancji estymatora $\hat{\mu}_{Imp}$ ma postać

$$v^2(\hat{\mu}_{Imp}) = \frac{E \frac{1}{|\mathbf{J}|} + \frac{n-E|\mathbf{J}|}{n^2} \chi^2}{\chi^2 - \frac{1}{n-1} + E|\mathbf{J}| \left(\frac{1}{n-1} - \frac{\chi^2}{n} \right)} S_{Imp}^2. \quad (\text{III.3.4})$$

W przypadku $p_i = p, i = 1, \dots, n$

$$E S_{Imp}^2 = \left(\chi^2 - \frac{1}{n-1} + \frac{np}{q_n} \left(\frac{1}{n-1} - \frac{\chi^2}{n} \right) \right) \sigma^2$$

oraz

$$v^2(\hat{\mu}_{Imp}) = \frac{E \frac{1}{|\mathbf{J}|} + \frac{q_n - p}{nq_n} \chi^2}{\chi^2 - \frac{1}{n-1} + \frac{np}{q_n} \left(\frac{1}{n-1} - \frac{\chi^2}{n} \right)} S_{Imp}^2.$$

Dowód. Zgodnie ze wzorem (I.3.4) mamy

$$S_{Imp}^2 = \frac{(|\mathbf{J}|-1)S_R^2 + \lambda^2 \left((n-|\mathbf{J}|-1)S_{R^c, Y}^2 + \frac{(n-|\mathbf{J}|)|\mathbf{J}|}{n} (\bar{Y}_{R^c} - v)^2 \right)}{n-1}.$$

Wzór (I.3.5) implikuje

$$E(S_{Imp}^2 | \mathbf{J}) = \sigma^2 \left(\frac{|\mathbf{J}|-1}{n-1} + \left(1 - \frac{|\mathbf{J}|}{n} \right) \chi^2 \right),$$

co prowadzi natychmiast do (III.3.3). Wzór (III.3.4) jest natychmiastową konsekwencją wzorów (III.3.2) i (III.3.3). ■

IV. LOSOWY ZBIÓR BRAKÓW OBSERWACJI — IMPUTACJA WIELOKROTNA

IV.1. Losowanie spośród respondentów

Zgodnie ze wzorem (II.0.3) estymator imputacji wielokrotnej jest średnią z estymatorów imputacyjnych

$$\hat{\mu}_{MImp} = \frac{1}{m} \sum_{l=1}^m \hat{\mu}_{Imp}^{(l)},$$

czyli zgodnie z twierdzeniem II.1.1

$$\hat{\mu}_{MImp} = \frac{|\mathbf{J}|}{n} \bar{X}_R + \frac{1}{nm} \sum_{l=1}^m \sum_{j \in R^c} X_{K_j^{(l)}}.$$

TWIERDZENIE IV.1.1. *Wartość oczekiwana estymatora wielokrotnej imputacji wynosi*

$$E\hat{\mu}_{MImp} = \mu,$$

a wariancja ma postać

$$\text{Var}\hat{\mu}_{MImp} = \frac{\sigma^2}{mn^2} \left(n+1 - E|\mathbf{J}| - nE\frac{1}{|\mathbf{J}|} \right) + \sigma^2 E\frac{1}{|\mathbf{J}|}. \quad (\text{IV.1.1})$$

W przypadku $p_i = p, i = 1, \dots, n$

$$\text{Var}\hat{\mu}_{MImp} = \frac{\sigma^2}{mn^2} \left(n+1 - \frac{np}{q_n} - nE\frac{1}{|\mathbf{J}|} \right) + \sigma^2 E\frac{1}{|\mathbf{J}|}.$$

Dowód. Ze wzoru na wartość oczekiwaną estymatora imputacyjnego hot-deck wynika, że

$$E(\hat{\mu}_{MImp}|\mathbf{J}) = \mu.$$

Więc wzór na wartość oczekiwaną estymatora $\hat{\mu}_{Imp}$ jest konsekwencją identyczności

$$E\hat{\mu}_{MImp} = EE(\hat{\mu}_{MImp}|\mathbf{J}).$$

Zgodnie ze wzorem (II.1.2) mamy

$$\text{Var}(\hat{\mu}_{MImp}|\mathbf{J}) = \frac{\sigma^2}{mn} \left(\frac{mn-1}{|\mathbf{J}|} + \frac{n+1}{n} - \frac{|\mathbf{J}|}{n} \right).$$

Ponieważ

$$\text{Var}\hat{\mu}_{MImp} = E\text{Var}(\hat{\mu}_{MImp}|\mathbf{J}) + \text{Var}E(\hat{\mu}_{MImp}|\mathbf{J})$$

i drugi składnik po prawej stronie jest równy zero otrzymujemy wzór (IV.1.1). ■

TWIERDZENIE IV.1.2. *Obciążenie estymatora Rubina wariancji \hat{v}_{MImp}^2 wynosi*

$$B\hat{v}_{MImp}^2 = E\hat{v}_{MImp}^2 - \text{Var}\hat{\mu}_{MImp} = -\frac{\sigma^2}{n(n-1)} \left(E|\mathbf{J}| + n(n+1)E\frac{1}{|\mathbf{J}|} - 2n+1 \right). \quad (\text{IV.1.2.})$$

Estymator nieobciążony wariancji estymatora $\hat{\mu}_{Mimp}$ ma postać

$$\hat{v}_*^2 = \frac{E \frac{1}{|\mathbf{J}|}}{1 - E \frac{1}{|\mathbf{J}|}} \left((n-1)\bar{U}_m - B_m \right) + \frac{1}{m} B_m.$$

Dowód. Dla estymatora Rubina wariancji \hat{v}_{Mimp}^2 , określonego wzorami (II.0.7)—(II.0.9), w przypadku wielokrotnej imputacji hot-deck zgodnie ze wzorem z dowodu twierdzenia II.1.2 mamy

$$E(\bar{U}_m | \mathbf{J}) = \frac{\sigma^2 (|\mathbf{J}| - 1)(n(n+1) - |\mathbf{J}|)}{|\mathbf{J}| n^2 (n-1)} \quad \text{oraz} \quad E(B_m | \mathbf{J}) = \frac{\sigma^2 (|\mathbf{J}| - 1)(n - |\mathbf{J}|)}{|\mathbf{J}| n^2}.$$

Zatem

$$E(\hat{v}_{Mimp}^2 | \mathbf{J}) = \frac{\sigma^2 (|\mathbf{J}| - 1)(2n - |\mathbf{J}|)}{|\mathbf{J}| n(n-1)} + \frac{\sigma^2 (\mathbf{J} - 1)(n - \mathbf{J})}{|\mathbf{J}| mn^2}.$$

W konsekwencji

$$E\hat{v}_{Mimp}^2 = \frac{\sigma^2}{n(n-1)} \left(2n - 1 - E|\mathbf{J}| - 2nE \frac{1}{|\mathbf{J}|} \right) + \frac{\sigma^2}{mn^2} \left(n + 1 - E|\mathbf{J}| - nE \frac{1}{|\mathbf{J}|} \right).$$

Ostatecznie, zgodnie ze wzorem (III.2.6) obciążenie estymatora \hat{v}_{Mimp}^2 wynosi

$$E\hat{v}_{Mimp}^2 - \text{Var} \hat{\mu}_{Mimp} = \frac{\sigma^2}{n(n-1)} \left(2n - 1 - E|\mathbf{J}| - 2nE \frac{1}{|\mathbf{J}|} \right) - \sigma^2 E \frac{1}{|\mathbf{J}|},$$

co po uproszczeniach daje wzór (IV.1.2).

Ponieważ $E B_m = \frac{\sigma^2}{n^2} \left(n + 1 - E|\mathbf{J}| - nE \frac{1}{|\mathbf{J}|} \right)$, więc ze wzoru (IV.1.1) wynika, że wystarczy pokazać równość

$$E((n-1)\bar{U}_m - B_m) = \sigma^2 \left(1 - E \frac{1}{|\mathbf{J}|} \right),$$

a to wynika wprost ze wzorów na \bar{U}_m i B_m . ■

Wniosek IV.1.3. *Asymptotycznie względne obciążenie estymatora Rubina \hat{v}_{MImp}^2 wynosi*

$$\lim_{n \rightarrow \infty} \frac{E\hat{v}_{MImp}^2 - \text{Var}\hat{\mu}_{MImp}}{\text{Var}\hat{\mu}_{MImp}} = -\frac{(1-p)^2}{1 + \frac{p(1-p)}{m}}. \quad (\text{IV.1.3})$$

Dowód. Ze wzorów (IV.1.1) i (IV.1.2) wynika, że wzór na względne obciążenie można napisać w postaci

$$K_n = -\frac{E|\mathcal{J}| + n(n+1)E\frac{1}{|\mathcal{J}|} - 2n + 1}{E\frac{n(n-1)}{|\mathcal{J}|} + \frac{n-1}{mn} \left(n+1 - E|\mathcal{J}| - nE\frac{1}{|\mathcal{J}|} \right)}.$$

Ponieważ (Marciniak, Wesołowski, 1999)

$$\lim_{n \rightarrow \infty} nE\frac{1}{|\mathcal{J}|} = \frac{1}{p}, \quad \text{czyli} \quad \lim_{n \rightarrow \infty} E\frac{1}{|\mathcal{J}|} = 0,$$

więc dzieląc licznik i mianownik przez n i wykorzystując fakt, że

$$\lim_{n \rightarrow \infty} \frac{1}{n} E|\mathcal{J}| = p$$

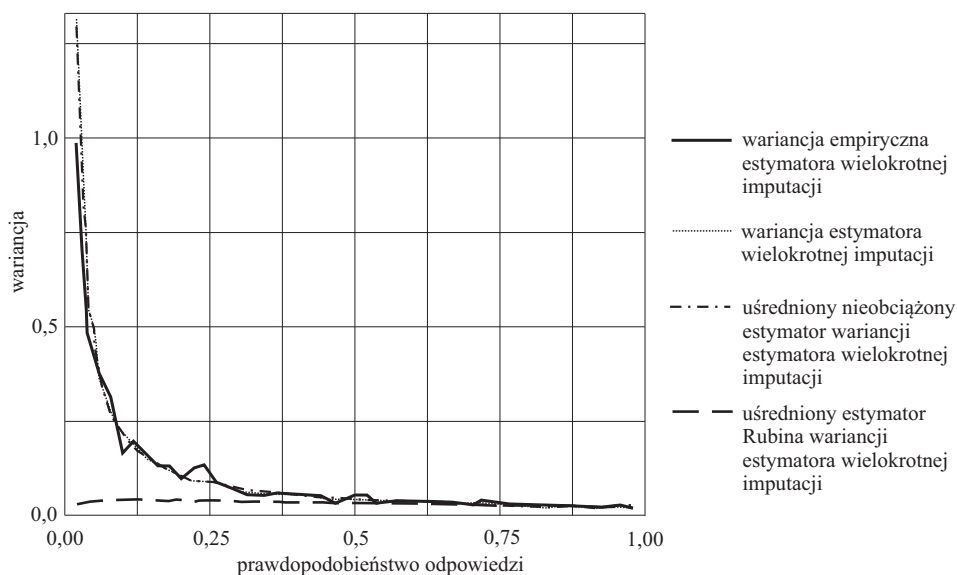
otrzymujemy

$$\begin{aligned} \lim_{n \rightarrow \infty} K_n &= -\frac{\lim_{n \rightarrow \infty} \frac{1}{n} E|\mathcal{J}| + \lim_{n \rightarrow \infty} (n+1)E\frac{1}{|\mathcal{J}|} - 2 + \lim_{n \rightarrow \infty} \frac{1}{n}}{\lim_{n \rightarrow \infty} (n-1)E\frac{1}{|\mathcal{J}|} + \frac{1}{m} \left(\lim_{n \rightarrow \infty} \frac{n-1}{n} \right) \left(1 + \lim_{n \rightarrow \infty} \frac{1}{n} - \lim_{n \rightarrow \infty} \frac{1}{n} E|\mathcal{J}| - \lim_{n \rightarrow \infty} E\frac{1}{|\mathcal{J}|} \right)} = \\ &= \frac{p + \frac{1}{p} - 2}{\frac{1}{p} + \frac{1}{m}(1-p)}, \end{aligned}$$

co po uproszczeniach daje wzór (IV.1.3). ■

Symulacyjna ilustracja wyników zawartych w twierdzeniach IV.1.1 i IV.1.2 przedstawiona jest na wykr. 3. Szczegółowy opis procedury symulacyjnej znajduje się w aneksie.

Wykr. 3. WYNIKI ESTYMACJI WARIANCJI ESTYMATORA ŚREDNIEJ DLA WIELOKROTNEJ IMPUTACJI (HOT-DECK) PRZY LOSOWYM BRAKU OBSERWACJI



Źródło: jak przy wykr. 1.

IV.2. Losowanie z rozkładu normalnego

TWIERDZENIE IV.2.1. *Wartość oczekiwana estymatora imputacji wielokrotnej przy losowaniu z rozkładu normalnego wynosi*

$$E\hat{\mu}_{MImp} = \mu,$$

a jego wariacja ma postać

$$\text{Var}\hat{\mu}_{MImp} = \sigma^2 \left(E \frac{1}{|J|} + \frac{n - E|J|}{mn^2} \right). \tag{IV.2.1}$$

Dowód. Powyższe wzory wynikają wprost ze wzorów z twierdzenia II.1.4. ■

TWIERDZENIE IV.2.2. *Obciążenie estymatora Rubina wariancji \hat{v}_{MImp}^2 wynosi*

$$B\hat{v}_{MImp}^2 = E\hat{v}_{MImp}^2 - \text{Var}\hat{\mu}_{MImp} = -\frac{\sigma^2}{n^2} \left(\frac{n - E|\mathbf{J}|}{n-1} - n^2 E \frac{1}{|\mathbf{J}|} - E|\mathbf{J}| + 2n \right). \quad (\text{IV.2.2})$$

Estymator nieobciążony wariancji estymatora $\hat{\mu}_{MImp}$ ma postać

$$\hat{v}_*^2 = \left(n E \frac{1}{|\mathbf{J}|} \right) \bar{U}_m + \left(\frac{1}{m} + \frac{n}{n-1} E \frac{1}{|\mathbf{J}|} \right) B_m.$$

Dowód. Zgodnie ze wzorami (II.1.1) i (II.1.2) mamy

$$E\bar{U}_m = \frac{\sigma^2}{n^2} \left(n - \frac{n - E|\mathbf{J}|}{n-1} \right) \quad \text{oraz} \quad EB_m = \frac{\sigma^2}{n^2} (n - E|\mathbf{J}|).$$

W konsekwencji

$$E\hat{v}_{MImp}^2 - \text{Var}\hat{\mu}_{MImp} = \frac{\sigma^2}{n^2} \left(n - \frac{n - E|\mathbf{J}|}{n-1} + n - E|\mathbf{J}| \right) - \sigma^2 E \frac{1}{|\mathbf{J}|},$$

co prowadzi do wzoru (IV.2.2).

Ze wzoru (IV.2.1) wynika, że wystarczy znaleźć liczby α i β takie, że

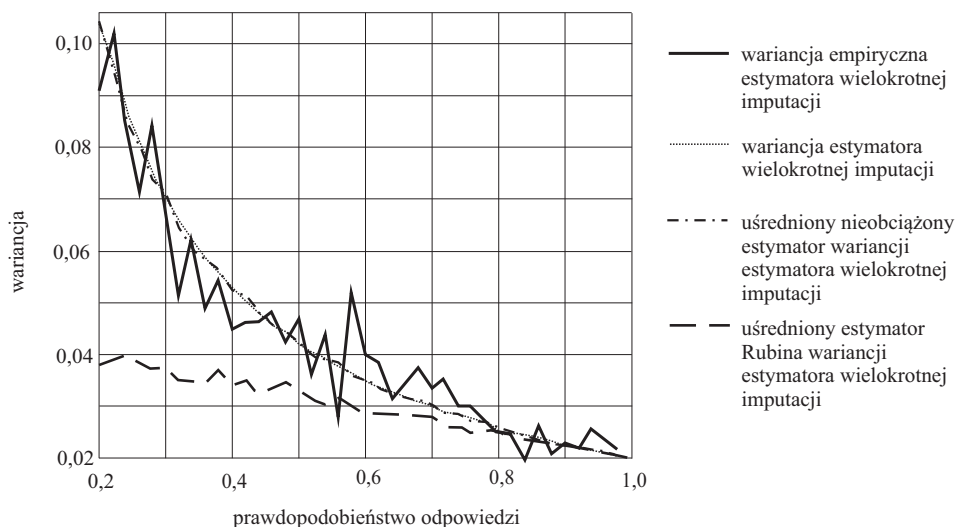
$$\alpha E\bar{U}_m + \beta EB_m = \sigma^2 E \frac{1}{|\mathbf{J}|}.$$

Wtedy $\alpha\bar{U}_m + \left(\beta + \frac{1}{m} \right) B_m$ jest poszukiwanym estymatorem nieobciążonym.

Przyjmując $\alpha = n E \frac{1}{|\mathbf{J}|}$ z powyższej równości otrzymujemy $\beta = \frac{n}{n-1} E \frac{1}{|\mathbf{J}|}$. ■

Symulacyjna ilustracja wyników zawartych w twierdzeniach IV.2.1 i IV.2.2 przedstawiona jest na wykr. 4. Szczegółowy opis procedury symulacyjnej znajduje się w aneksie.

**Wykr. 4. WYNIKI ESTYMACJI WARIANCJI ESTYMATORA ŚREDNIEJ
DLA WIELOKROTNEJ IMPUTACJI (HOT-DECK)
Z ROZKŁADU NORMALNEGO I LOSOWYM BRAKU OBSERWACJI**



Źródło: jak przy wyk. 1.

V. KONKLUZJE

W powyższym opracowaniu przedstawiono podstawowe idee i modele imputacji, w tym imputacji wielokrotnej na podstawie próbki niezależnych obserwacji o jednakowym rozkładzie, w której występują braki. W szczególności rozważono imputację średnią, imputacje typu hot-deck polegające na losowaniu respondentów oraz losowaniu z rozkładu normalnego z parametrami estymowanymi na podstawie obserwowanej części próbki oraz imputację regresyjną. Każda z tych metod jest szczególnym przypadkiem podejścia ogólnego analizowanego w rozdziałach I i III. W przypadku imputacji wielokrotnej wskazano, że popularny estymator wariacji, zwany estymatorem Rubina, nie ma dobrych własności, w szczególności jest obciążony. Zaproponowano wersje nieobciążone estymatorów typu Rubina w przypadku imputacji typu hot-deck metodą losowania respondentów i losowania z rozkładu normalnego. W artykule nie analizowano bayesowskiego modelu obserwacji, który jest ważnym elementem, zaproponowanego przez Rubina, podejścia wykorzystującego imputację wielokrotną. Planujemy przedstawić matematyczne podstawy imputacji w modelu bayesowskim w kolejnym opracowaniu.

LITERATURA

- Andridge, R.R., Little, R.J.A. (2010). A review of hot deck imputation survey non-response. *International Statistical Review*, vol. 78, no. 1, s. 40—64.
- Carpenter, J.R., Kenward, M.G. (2013). *Multiple Imputation and its Application*. Wiley, Chichester.
- de Waal, T., Pannekoek, J., Scholtus, S. (2011). *Handbook of Statistical Data Editing and Imputation*. Wiley, New York.
- Donders, A.R.T., van der Heijden, G.J.M.G., Stijnen, T., Moons, K.G.M. (2006). Review: a gentle introduction to imputation of missing values. *Journal of Clinical Epidemiology*, vol. 59, no. 10, s. 1087—1091.
- Garson, G.D. (2012). Missing Values Analysis & Data Imputation. *Statistical Associates Publishers*, Asheboro.
- Horton, N.J., Lipsitz, S.R. (2001). Multiple imputation in practice: comparison of software packages for regression models with missing data. *American Statistician*, vol. 55, no. 3, s. 244—254.
- Little, R.J.A., Rubin, D.B. (2002). *Statistical Analysis with Missing Data*. Wiley, New York.
- Marciniak, E., Wesołowski, J. (1999). Asymptotic Eulerian expansions for binomial and negative binomial reciprocals. *Proceedings of The American Mathematical Society*, vol. 127, s. 3329—3338.
- Misztal, M. (2012). Imputation of missing data using R package. *Acta Universitatis Lodzianis, Folia Oeconomica*, vol. 269, s. 131—144.
- Norazian Ramli, M.N., Yahaya, A.S., Ramli, N.A., Yusof, N.F.F.M., Abdullah, M.M.A. (2013). Roles of imputation methods for filling the missing values: a review. *Advances in Environmental Biology*, vol. 7, no. 12, s. 3861—3869.
- Rempała, G.A. (2004). Asymptotic factorial powers expansions for binomial and negative binomial reciprocals. *Proceedings of The American Mathematical Society*, vol. 32, no. 1, s. 261—272.
- Rubin, D.B. (1987). *Multiple Imputation for Nonresponse in Surveys*. Wiley, New York.
- van Buuren, S. (2012). *Flexible Imputation of Missing Data*. Chapman&Hall/CRC, London.
- van Buuren, S., Groothuis-Oudshoorn, K. (2011). mice: multivariate imputation by chained equations in R. *Journal of Statistical Software*, vol. 45, no. 3, s. 1—67, <http://www.jstatsoft.org/>.
- Schafer, J.L. (1997). *Analysis of Incomplete Multivariate Data*. Chapman and Hall, Boston.
- Szabłowski, P.J., Wesołowski, J., Wieczorkowski, R. (1996). Estymacja w podpopulacjach. *Wiadomości Statystyczne*, nr 7, s. 1—13.

ANEKS

Badania symulacyjne

Poniżej przedstawiony został opis badań symulacyjnych przeprowadzonych nad rozważanymi w pracy estymatorami wariancji wykorzystującymi imputację wielokrotną metodą typu hot-deck.

Symulacje polegały na $L = 100$ -krotnym powtórzeniu następującego algorytmu:

- 1) generowano dane pochodzące z centralnego rozkładu normalnego z wariancją $\sigma^2 = 4$;
- 2) usuwano (losowo bądź nie, w zależności od rozważanego modelu) obserwacje z wygenerowanych danych;
- 3) generowano $m = 5$ próbek imputacyjnych w sposób określony przez rozważany model imputacyjny;

- 4) wyliczono estymator $\hat{\mu}_{MImp}$ oraz jego estymatory wariancji (Rubina — \hat{v}_{MImp}^2 i nieobciążony — \hat{v}_*^2) na podstawie wygenerowanych próbek imputacyjnych. Na wszystkich rysunkach przedstawiono wykresy następujących wielkości:
 — $\widehat{\text{Var}}(\hat{\mu}_{MImp})$ — wariancji empirycznej estymatora wielokrotnej imputacji $\hat{\mu}_{MImp}$, obliczanej na podstawie L powtórzeń eksperymentu

$$\widehat{\text{Var}}(\hat{\mu}_{MImp}) = \frac{1}{L-1} \sum_{l=1}^L \left(\hat{\mu}_{MImp}^{(l)} - \frac{1}{L} \sum_{j=1}^L \hat{\mu}_{MImp}^{(j)} \right)^2,$$

- gdzie $\hat{\mu}_{MImp}^{(l)}$ oznacza wartość estymatora wielokrotnej imputacji $\hat{\mu}_{MImp}$ w l -tym powtórzeniu eksperymentu, $l=1, \dots, L$;
 — $\widehat{\text{Var}}\hat{\mu}_{MImp}$ — wariancji estymatora wielokrotnej estymacji, obliczanej na podstawie wzorów: II.1.2 (wykr. 1), II.1.7 (wykr. 2), IV.1.1 w wersji $p_i = p$ (wykr. 3) oraz IV.2.1 (wykr. 4);
 — $\overline{\hat{v}_*^2}$ — uśrednionej, na podstawie L powtórzeń eksperymentu, wartości \hat{v}_*^2 nieobciążonego estymatora wariancji estymatora wielokrotnej estymacji:

$$\overline{\hat{v}_*^2} = \frac{1}{L} \sum_{l=1}^L \hat{v}_{*,l}^2,$$

- gdzie $\hat{v}_{*,l}^2$ oznacza wartość estymatora \hat{v}_*^2 w l -tym powtórzeniu eksperymentu, $l = 1, \dots, L$;
 — $\overline{\hat{v}_{MImp}^2}$ — uśrednionej, na podstawie L powtórzeń eksperymentu, wartości \hat{v}_{MImp}^2 estymatora Rubina wariancji estymatora wielokrotnej imputacji:

$$\overline{\hat{v}_{MImp}^2} = \frac{1}{L} \sum_{l=1}^L \hat{v}_{MImp,l}^2,$$

- gdzie $\hat{v}_{MImp,l}^2$ oznacza wartość estymatora Rubina \hat{v}_{MImp}^2 w l -tym powtórzeniu eksperymentu, $l = 1, \dots, L$.

Summary. *The article presents the basics of imputation methodology (including the methodology of multiple imputation), focusing on understanding its mathematical background. We analyze the situation when observations in the original sample are independent random variables with identical distributions,*

and response or its lack is modeled by a random mechanism which is independent of observations. In particular, we point out to problems that arise when the standard Rubin estimate of the multiple imputation variance estimator is used. A possible improvement of this popular estimator is indicated. The starting point of the analysis is when the appearance of response deficiencies is caused by a deterministic mechanism.

Keywords: imputation, multiple imputation, imputation estimator, Rubin estimator, mean imputation, hot-deck imputation, regression imputation.

***Резюме.** В статье представлены основы импутационной методологии (в том числе методологии многократной импутации). Внимание в статье сосредоточено на прояснении математической стороны вопросов. Проанализирована ситуация, когда наблюдения формирующие оригинальную выборку являются независимыми случайными величинами с одинаковыми распределениями, а отсутствие ответов появляется случайно независимо от наблюдения. В частности статья указывает на проблемы, которые возникают когда используется стандартная оценка Рубина дисперсии оценки многократной импутации. В статье указано также на возможное улучшение этой популярной оценки. Отправной точкой анализа является ситуация, когда отсутствие ответов объясняет детерминистический механизм.*

Ключевые слова: импутация, многократная импутация, импутационная оценка, оценка Рубина, импутация средним, импутация типа hot-deck, регрессионная импутация.