

Imputacja brakujących danych binarnych w modelu autologistycznym ¹

Marta Zalewska

Warszawski Uniwersytet Medyczny

Statystyka Matematyczna
Wiśła, grudzień 2009

¹Współautorzy: Wojciech Niemirow, UMK Toruń i UW, Bolesław Samoliński, WUM

Plan

- 1 Imputacja Bayesowska i Próbnik Gibbsa
- 2 Model autologistyczny
- 3 Estymacja
 - Maksimum pseudowiarogodności
 - Maksimum wiarogodności i estymacja Bayesowska
- 4 Wyniki symulacyjne
 - Metodologia
 - Wyniki

Plan

- 1 Imputacja Bayesowska i Próbnik Gibbsa
- 2 Model autologistyczny
- 3 Estymacja
 - Maksimum pseudowiarogodności
 - Maksimum wiarogodności i estymacja Bayesowska
- 4 Wyniki symulacyjne
 - Metodologia
 - Wyniki

Plan

- 1 Imputacja Bayesowska i Próbnik Gibbsa
- 2 Model autologistyczny
- 3 Estymacja
 - Maksimum pseudowiarogodności
 - Maksimum wiarogodności i estymacja Bayesowska
- 4 Wyniki symulacyjne
 - Metodologia
 - Wyniki

Plan

- 1 Imputacja Bayesowska i Próbnik Gibbsa
- 2 Model autologistyczny
- 3 Estymacja
 - Maksimum pseudowiarogodności
 - Maksimum wiarogodności i estymacja Bayesowska
- 4 Wyniki symulacyjne
 - Metodologia
 - Wyniki

- 1 Imputacja Bayesowska i Próbnik Gibbsa
- 2 Model autologistyczny
- 3 Estymacja
 - Maksimum pseudowiarygodności
 - Maksimum wiarygodności i estymacja Bayesowska
- 4 Wyniki symulacyjne
 - Metodologia
 - Wyniki

Imputacja Bayesowska i Próbnik Gibbsa

Dane:

$$\mathbf{X} = (\mathbf{X}_{\text{obs}}, \mathbf{X}_{\text{mis}}) \sim \mathbf{p}_{\beta}.$$

- \mathbf{X}_{obs} – obserwowane dane
- \mathbf{X}_{mis} – brakujące dane
- $\beta \sim \pi(\cdot)$ – parametr o rozkładzie *a priori*

Próbnik Gibbsa. Powtarzamy dwa kroki:

- generujemy β z rozkładu $\pi(\beta | \mathbf{X}_{\text{obs}}, \mathbf{X}_{\text{mis}})$,
- generujemy \mathbf{X}_{mis} z rozkładu $\pi(\mathbf{X}_{\text{mis}} | \mathbf{X}_{\text{obs}}, \beta)$.

Metoda MCMC: łańcuch Markowa zbieżny do $\pi(\mathbf{X}_{\text{mis}}, \beta | \mathbf{X}_{\text{obs}})$.

Imputacja Bayesowska i Próbnik Gibbsa

Dane:

$$\mathbf{X} = (\mathbf{X}_{\text{obs}}, \mathbf{X}_{\text{mis}}) \sim \mathbf{p}_{\beta}.$$

- \mathbf{X}_{obs} – obserwowane dane
- \mathbf{X}_{mis} – brakujące dane
- $\beta \sim \pi(\cdot)$ – parametr o rozkładzie *a priori*

Próbnik Gibbsa. Powtarzamy dwa kroki:

- generujemy β z rozkładu $\pi(\beta | \mathbf{X}_{\text{obs}}, \mathbf{X}_{\text{mis}})$,
- generujemy \mathbf{X}_{mis} z rozkładu $\pi(\mathbf{X}_{\text{mis}} | \mathbf{X}_{\text{obs}}, \beta)$.

Metoda MCMC: łańcuch Markowa zbieżny do $\pi(\mathbf{X}_{\text{mis}}, \beta | \mathbf{X}_{\text{obs}})$.

- 1 Imputacja Bayesowska i Próbnik Gibbsa
- 2 Model autologistyczny**
- 3 Estymacja
 - Maksimum pseudowiarogodności
 - Maksimum wiarogodności i estymacja Bayesowska
- 4 Wyniki symulacyjne
 - Metodologia
 - Wyniki

Model

Rozkład autologistyczny na $\mathcal{X} = \{0, 1\}^d$:

$$p_{\beta}(\mathbf{x}) := \frac{1}{Z(\beta)} \exp \left(\sum_{i,j=1}^d \beta_{ij} x_i x_j \right),$$

$\mathbf{x} \sim \text{AL}(\beta)$. Parametry: macierz symetryczna $\beta = (\beta_{ij})$.

Pełne rozkłady warunkowe:

$$p_{\beta}(x_i = 1 \mid \mathbf{x}_{-i}) = \frac{\exp \left(\beta_{ii} + \sum_{j \neq i} x_j \beta_{ij} \right)}{1 + \exp \left(\beta_{ii} + \sum_{j \neq i} x_j \beta_{ij} \right)},$$

gdzie $\mathbf{x}_{-i} = (x_j, j \neq i)$.

Rozkład taki sam jak w standardowym modelu regresji logistycznej.

Model

Rozkład autologistyczny na $\mathcal{X} = \{0, 1\}^d$:

$$p_{\beta}(x) := \frac{1}{Z(\beta)} \exp \left(\sum_{i,j=1}^d \beta_{ij} x_i x_j \right),$$

$x \sim \text{AL}(\beta)$. Parametry: macierz symetryczna $\beta = (\beta_{ij})$.

Pełne rozkłady warunkowe:

$$p_{\beta}(x_i = 1 \mid x_{-i}) = \frac{\exp \left(\beta_{ii} + \sum_{j \neq i} x_j \beta_{ij} \right)}{1 + \exp \left(\beta_{ii} + \sum_{j \neq i} x_j \beta_{ij} \right)},$$

gdzie $x_{-i} = (x_j, j \neq i)$.

Rozkład taki sam jak w standardowym modelu regresji logistycznej.

Symulacja i estymacja

Wniosek:

- symulowanie x jest łatwe przy pomocy próbnika Gibbsa,
- parametry β można estymować standardowymi metodami GLM.

Próbka: $x(1), \dots, x(n)$ i.i.d. $\sim AL(\beta)$.

Symulacja i estymacja

Wniosek:

- symulowanie x jest łatwe przy pomocy próbnika Gibbsa,
- parametry β można estymować standardowymi metodami GLM.

Próbka: $x(1), \dots, x(n)$ i.i.d. $\sim AL(\beta)$.

- 1 Imputacja Bayesowska i Próbnik Gibbsa
- 2 Model autologistyczny
- 3 Estymacja**
 - Maksimum pseudowiarogodności
 - Maksimum wiarogodności i estymacja Bayesowska
- 4 Wyniki symulacyjne
 - Metodologia
 - Wyniki

Maksimum pseudowiarogodności

Częstkowe wiarogodności:

$$L_i(\beta|\mathbf{x}) = L_i(\beta_i|\mathbf{x}) := \log p_{\beta}(x_i | \mathbf{x}_{-i}).$$

Pseudowiarogodność:

$$L_{\text{ps}}(\beta|\mathbf{X}) = \sum_{k=1}^n \sum_{i=1}^d L_i(\beta|\mathbf{x}(k)).$$

$\mathbf{X} = (\mathbf{x}(1), \dots, \mathbf{x}(n))$ – próbka.

Maksimum wiarogodności

Estymator największej wiarogodności via MCMC: Geyer i Thompson (1992, JRSS).

Rodzina wykładnicza:

$$p_{\beta}(x) = \frac{1}{Z(\beta)} e^{\beta^T T(x)},$$

gdzie $T(x)$ – wektor statystyk dostatecznych.

$$\begin{aligned} Z(\beta) &= \sum_{x^*} e^{\beta^T T(x^*)} = \sum_{x^*} e^{(\beta - \beta^*)^T T(x^*)} p_{\beta^*}(x^*) Z(\beta^*) \\ &= \mathbb{E} e^{(\beta - \beta^*)^T T(x^*)} \times Z(\beta^*). \end{aligned}$$

gdzie $x^* \sim p_{\beta^*}$.

Aproksymacja wiarogodności: $x^*(1), \dots, x^*(n^*) \sim p_{\beta^*}$,

$$L_{\text{MCMC}}(\beta|x) = \beta^T T(x) - \log \sum_{k=1}^{n^*} e^{(\beta - \beta^*)^T T(x^*(k))} + \text{const.}$$

Maksimum wiarogodności

Estymator największej wiarogodności via MCMC: Geyer i Thompson (1992, JRSS).

Rodzina wykładnicza:

$$p_{\beta}(x) = \frac{1}{Z(\beta)} e^{\beta^T T(x)},$$

gdzie $T(x)$ – wektor statystyk dostatecznych.

$$\begin{aligned} Z(\beta) &= \sum_{x^*} e^{\beta^T T(x^*)} = \sum_{x^*} e^{(\beta - \beta^*)^T T(x^*)} p_{\beta^*}(x^*) Z(\beta^*) \\ &= \mathbb{E} e^{(\beta - \beta^*)^T T(x^*)} \times Z(\beta^*). \end{aligned}$$

gdzie $x^* \sim p_{\beta^*}$.

Aproksymacja wiarogodności: $x^*(1), \dots, x^*(n^*) \sim p_{\beta^*}$,

$$L_{\text{MCMC}}(\beta|x) = \beta^T T(x) - \log \sum_{k=1}^{n^*} e^{(\beta - \beta^*)^T T(x^*(k))} + \text{const.}$$

Próbkowanie z rozkładu *a posteriori*

Wybieramy $\beta^* = \hat{\beta}_{\text{ps}}$, $n^* = n$.

Estymator „jednokrokowy” Newtona-Raphsona

$$\hat{\beta} = \hat{\beta}_{\text{ps}} - \nabla^2 L_{\text{MCMC}}(\hat{\beta}_{\text{ps}})^{-1} \nabla L_{\text{MCMC}}(\hat{\beta}_{\text{ps}}),$$

gdzie

- $\hat{\beta}_{\text{ps}}$ – estymator największej pseudowiarogodności,
- L_{MCMC} – aproksymacja wiarogodności na podstawie „sztucznej próbki” $X^* = (x^*(1), \dots, x^*(n))$ wygenerowanej z rozkładu $p_{\hat{\beta}_{\text{ps}}}$.

$\hat{\beta}$ ma w przybliżeniu rozkład $\pi(\beta|X)$.

Próbkowanie z rozkładu *a posteriori*

Wybieramy $\beta^* = \hat{\beta}_{\text{ps}}$, $n^* = n$.

Estymator „jednokrokowy” Newtona-Raphsona

$$\hat{\beta} = \hat{\beta}_{\text{ps}} - \nabla^2 L_{\text{MCMC}}(\hat{\beta}_{\text{ps}})^{-1} \nabla L_{\text{MCMC}}(\hat{\beta}_{\text{ps}}),$$

gdzie

- $\hat{\beta}_{\text{ps}}$ – estymator największej pseudowiarogodności,
- L_{MCMC} – aproksymacja wiarogodności na podstawie „sztucznej próbki” $X^* = (x^*(1), \dots, x^*(n))$ wygenerowanej z rozkładu $p_{\hat{\beta}_{\text{ps}}}$.

$\hat{\beta}$ ma w przybliżeniu rozkład $\pi(\beta|X)$.

- 1 Imputacja Bayesowska i Próbnik Gibbsa
- 2 Model autologistyczny
- 3 Estymacja
 - Maksimum pseudowiarogodności
 - Maksimum wiarogodności i estymacja Bayesowska
- 4 **Wyniki symulacyjne**
 - Metodologia
 - Wyniki

Metodologia

Dane: ECAP: (Choroby alergiczne w Polsce, 2006-2008).
18617 jednostek (przypadków, respondentów) i 1225
zmiennych (większość binarnych, ale również liczbowe).

Do naszych eksperymentów wybraliśmy małą podmacierz:
 $n = 2962$ jednostek, $d = 6$ zmiennych, **bez brakujących danych**.

- generujemy sztucznie „braki”,
- imputujemy i estymujemy,
- porównujemy i sprawdzamy.

Metodologia

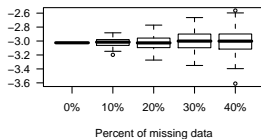
Dane: ECAP: (Choroby alergiczne w Polsce, 2006-2008).
18617 jednostek (przypadków, respondentów) i 1225
zmiennych (większość binarnych, ale również liczbowe).

Do naszych eksperymentów wybraliśmy małą podmacierz:
 $n = 2962$ jednostek, $d = 6$ zmiennych, **bez brakujących danych**.

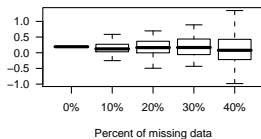
- generujemy sztucznie „braki”,
- imputujemy i estymujemy,
- porównujemy i sprawdzamy.

Estymatory

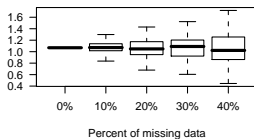
Estimates of β_{11}



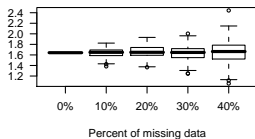
Estimates of β_{12}



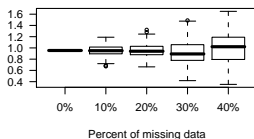
Estimates of β_{13}



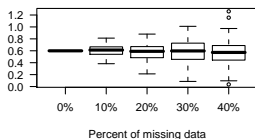
Estimates of β_{14}



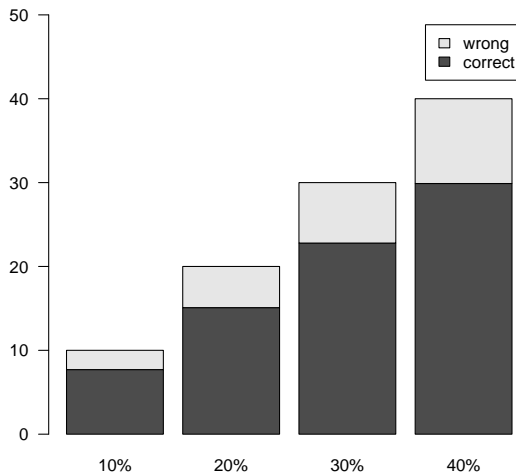
Estimates of β_{15}



Estimates of β_{16}



Rezultaty imputacji



Percent of missings correctly and wrongly imputed