# Mixed effects model for SAE

Tomáš Košlab

FNSPE CTU in Prague

June 22, 2018

Tomáš Košlab

Mixed effects model for SAE

FNSPE CTU in Prague

・ロン ・回 と ・ ヨ と ・ ヨ と …

∃ nar

### Outline

Model

Parameter estimation

Prediction of area effects

Simulation experiment

Tomáš Košlab

Mixed effects model for SAE

- N individuals in D domains,  $N_d$  individuals in the d-th domain
- every individual modelled by a Bernoulli-distributed r.v. Y<sub>dj</sub> (above / below poverty line), realization y<sub>dj</sub>
- sample of size n<sub>d</sub> from every domain
- ► task: predict the population mean in every domain,  $\overline{y}_d = \frac{1}{N_d} \sum_{j=1}^{N_d} y_{dj}$
- if n<sub>d</sub> is too small for a direct estimate of sufficient quality small area

イロト イポト イヨト イヨト 一日

### Proposed model

For  $Y_{dj} \sim Be(p_{dj})$  we propose the following logit regression model

(F) 
$$\text{logit}(p_{dj}) = \mathbf{x}_{dj}^T \beta + \mu_d, \quad d = 1, \dots, D_F, \quad j = 1, \dots, n_d,$$
  
(R)  $\text{logit}(p_{dj}) = \mathbf{x}_{dj}^T \beta + u_d, \quad d = D_F + 1, \dots, D, \quad j = 1, \dots, n_d$   
(1)

where

- x<sup>T</sup><sub>dj</sub> = (x<sub>dj,1</sub>,..., x<sub>dj,p</sub>) is the vector of covariates belonging to the *j*-th individual in the *d*-area,
   β = (β<sub>1</sub>,..., β<sub>p</sub>)<sup>T</sup> is the vector of unknown fixed parameters,
   μ = (μ<sub>1</sub>,..., μ<sub>D<sub>e</sub></sub>)<sup>T</sup> is the vector of fixed area effects,
- $\boldsymbol{u} = (u_{D_F+1}, \dots, u_D)^T$  is the vector of random area effects.

Tomáš Košlab

# Model - assumptions & notes

- data in different areas are independent (each area has its own effect)
- data in areas modelled by a fixed effect are independent
- probability p<sub>dj</sub> can be expressed as

(F) 
$$p_{dj} = \frac{\exp(\mathbf{x}_{dj}^{T}\beta + \mu_{d})}{1 + \exp(\mathbf{x}_{dj}^{T}\beta + \mu_{d})}, \quad d = 1, ..., D_{F}, \quad j = 1, ..., n_{d},$$
  
(R)  $p_{dj} = \frac{\exp(\mathbf{x}_{dj}^{T}\beta + u_{d})}{1 + \exp(\mathbf{x}_{dj}^{T}\beta + u_{d})}, \quad d = D_{F} + 1, ..., D, \quad d = 1, ..., n_{d}$ 
(2)

▲ロ → ▲ □ → ▲ □ → ▲ □ → ▲ □ → ▲ □ → ▲ □ → ▲ □ → ▲ □ →

Tomáš Košlab

# PQL method

#### log-likelihood function

$$I(\beta, \mu, \sigma^{2}; \mathbf{y}) = \sum_{d=1}^{D_{F}} \sum_{j=1}^{n_{d}} \left[ y_{dj} \log p_{dj} + (1 - y_{dj}) \log(1 - p_{dj}) \right] \\ + \sum_{d=D_{F}+1}^{D} \log \int_{R} \prod_{j=1}^{n_{d}} p_{dj}^{y_{dj}} (1 - p_{dj})^{1 - y_{dj}} \frac{1}{\sqrt{2\pi\sigma^{2}}} e^{-\frac{u_{d}^{2}}{2\sigma^{2}}} du_{d}$$
(3)

- PQL can be derived from the Laplace approximation of the log-likelihood function, *I<sub>Laplace</sub>*
- omission of the last term in *I<sub>Laplace</sub>* (for computational reasons) leads to *I<sub>PQL</sub>*
- estimates of  $\beta, \mu$  are obtained by N.-R. algorithm, estimate of  $\sigma^2$  by fixed-point

Tomáš Košlab

#### Empirical Best Predictor (EBP)

- closely related to the Best Predictor (BP)
- BP of  $\hat{\overline{y}}_d$  can be expressed as

$$\hat{\overline{y}}_d = \frac{1}{N_d} \left( \sum_{j \in s_d} y_{dj} + \sum_{j \in r_d} \hat{y}_{dj} \right) = \frac{1}{N_d} \left( \sum_{j \in s_d} y_{dj} + \sum_{j \in r_d} \hat{p}_{dj} \right)$$
(4)

where

- predictions are denoted by hats
- s<sub>d</sub> and r<sub>d</sub> denote the indices of observations from the d-th area that are inside and outside the sample respectively

 EBP is obtained by substituting parameter estimates into formulas for BP

Tomáš Košlah

# Plug-in predictor

- the formula for plug-in predictor coincides with (4)
- ► the difference lies in the term p̂<sub>dj</sub> while for EBP (BP) it is calculated by a formula, plug-in is based on the formulas (2)

$$p_{dj} = \frac{\exp(\mathbf{x}_{dj}^{T}\boldsymbol{\beta} + \mu_{d})}{1 + \exp(\mathbf{x}_{dj}^{T}\boldsymbol{\beta} + \mu_{d})}, \quad d = 1, \dots, D_{F}, \quad j = 1, \dots, n_{d},$$

$$p_{dj} = \frac{\exp(\mathbf{x}_{dj}^{T}\boldsymbol{\beta} + u_{d})}{1 + \exp(\mathbf{x}_{dj}^{T}\boldsymbol{\beta} + u_{d})}, \quad d = D_{F} + 1, \dots, D, \quad d = 1, \dots, n_{d}$$
(5)

where the estimates of  $\beta$ ,  $\mu_d$  and the predictions of  $u_d$  are substituted into the equations

requires predictions of u<sub>d</sub>

Tomáš Košlah

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 うの()

# Setup

- D = 30 domains
- $D_F = 5$  areas modelled by a fixed effect
- ▶  $N_d = 1000, d = 1, ..., D$
- design matrix  $\rightarrow$  3 parameters  $\beta_1, \beta_2, \beta_3$ 
  - ► *X*<sub>dj,1</sub> ~ *Be*(0.48)
  - ► *X*<sub>dj,2</sub> ~ *Be*(0.6)
  - if  $x_{dj,2} = 1$ , then  $X_{dj,3} \sim Be(0.5)$ , else  $x_{dj,3} = 1$
- $Y_{dj} \sim Be(p_{dj})$ , values of  $y_{dj}$  are generated using equation (2)
- ▶ different sample sizes for areas with fixed (n<sup>F</sup><sub>d</sub>) and random effects n<sup>R</sup><sub>d</sub>
- ► task: prediction of  $\overline{y}_d$  for every domain for the whole population

### Results



Figure: BIAS of predictions using the respective methods.

<ロ> <同> <同> < 回> < 回>

Tomáš Košlab

Mixed effects model for SAE



Figure: MSE of predictions using the respective methods.

Tomáš Košlab

Mixed effects model for SAE

FNSPE CTU in Prague

# Conclusion & future tasks

- predictions obtained by EBP and plug-in predictor are of comparable quality
- both are superior compared to the direct estimate, especially for small amounts of data
- estimate of prediction error parametric bootstrap
- application on real data and comparison with other models

イロン 不同 とくほう イロン

# Thank you for your attention!

Tomáš Košlab

Mixed effects model for SAE

▶ < 돌 > 돌 つへ( ENSPE CTU in Prague

イロト イポト イヨト イヨト