

On Models for Small Area Compositions

Angela Luna

A.Luna-Hernandez@soton.ac.uk

Li-Chun Zhang

L.Zhang@soton.ac.uk

Social Statistics & Demography
University of Southampton

Acknowledgements: This work has been funded by the Office for National Statistics - ONS and the Economic and Social Research Council - ESRC.

Motivation

Compositions

Area	1	...	J	Total
1	Y_{aj}			$Y_{1.}$
2				$Y_{2.}$
3				$Y_{3.}$
...				...
A				$Y_{A.}$
Total	$Y_{.1}$...	$Y_{.J}$	$Y_{..}$

Area	1	...	J	Total
1	Y_{aj}			$Y_{1.}$
2				$Y_{2.}$
3				$Y_{3.}$
...				...
A				$Y_{A.}$
Total	$Y_{.1}$...	$Y_{.J}$	$Y_{..}$

Target: Estimate the within area cell counts Y_{aj} , using proxy information and fixed row/column margins.

Two different (fixed-effects) approaches to this problem are considered:

- ▷ **Structure Preserving Models:** Long tradition in SAE. Assumptions about the relationship between the interactions of two compositions in the log-linear scale. (Proxy information is required).

Two different (fixed-effects) approaches to this problem are considered:

- ▷ **Structure Preserving Models:** Long tradition in SAE. Assumptions about the relationship between the interactions of two compositions in the log-linear scale. (Proxy information is required).
- ▷ **Regression (Generalized Linear) Models:** Multinomial-Logistic: Assumptions about the relationship between the log-odds with respect to a reference category and a set of covariates. (Proxy information can be used as covariate).

In this work, we:

- 1 Introduce a generalization of the Structure Preserving approach, covering the SPREE and GSPREE models and also the logit-multinomial (using proxy information) as particular cases.

In this work, we:

- 1 Introduce a generalization of the Structure Preserving approach, covering the SPREE and GSPREE models and also the logit-multinomial (using proxy information) as particular cases.
- 2 Use data from 2001 and 2011 Population Censuses in England to compare the different models in terms of their Prediction Error.

In this work, we:

- 1 Introduce a generalization of the Structure Preserving approach, covering the SPREE and GSPREE models and also the logit-multinomial (using proxy information) as particular cases.
- 2 Use data from 2001 and 2011 Population Censuses in England to compare the different models in terms of their Prediction Error.
- 3 Show some ongoing work on a model using a mapping matrix between the proxy and desired compositions, which allows to incorporate auxiliary information at the aggregate level.

- 1 Structure Preserving Models
- 2 Model using a Mapping matrix

Structure Preserving Models

Structure Preserving Models

Denote by $\theta_{aj}^X = X_{aj}/X_a$. an auxiliary composition of exactly the same dimension as $\theta_{aj}^Y = Y_{aj}/Y_a$, its log-linear representation given by:

$$\gamma_{aj}^X = \alpha_0^X + \alpha_a^X + \alpha_j^X + \alpha_{aj}^X$$

where $\gamma_{aj}^X = \log \theta_{aj}^X$, $\alpha_0^X = \bar{\gamma}_{..}^X$, $\alpha_a^X = \bar{\gamma}_{a.}^X - \bar{\gamma}_{..}^X$, $\alpha_j^X = \bar{\gamma}_{.j}^X - \bar{\gamma}_{..}^X$ and $\alpha_{aj}^X = \gamma_{aj}^X - \bar{\gamma}_{a.}^X - \bar{\gamma}_{.j}^X + \bar{\gamma}_{..}^X$.

Denote by $\theta_{aj}^X = X_{aj}/X_a$. an auxiliary composition of exactly the same dimension as $\theta_{aj}^Y = Y_{aj}/Y_a$, its log-linear representation given by:

$$\gamma_{aj}^X = \alpha_0^X + \alpha_a^X + \alpha_j^X + \alpha_{aj}^X$$

where $\gamma_{aj}^X = \log \theta_{aj}^X$, $\alpha_0^X = \bar{\gamma}_{..}^X$, $\alpha_a^X = \bar{\gamma}_{a.}^X - \bar{\gamma}_{..}^X$, $\alpha_j^X = \bar{\gamma}_{.j}^X - \bar{\gamma}_{..}^X$ and $\alpha_{aj}^X = \gamma_{aj}^X - \bar{\gamma}_{a.}^X - \bar{\gamma}_{.j}^X + \bar{\gamma}_{..}^X$.

The log-linear representation satisfies the constraints:

$\sum_a \alpha_a^X = 0$, $\sum_j \alpha_j^X = 0$, $\sum_a \alpha_{aj}^X = \sum_j \alpha_{aj}^X = 0$. Analogous for θ_{aj}^Y .

Denote by $\theta_{aj}^X = X_{aj}/X_a$. an auxiliary composition of exactly the same dimension as $\theta_{aj}^Y = Y_{aj}/Y_a$, its log-linear representation given by:

$$\gamma_{aj}^X = \alpha_0^X + \alpha_a^X + \alpha_j^X + \alpha_{aj}^X$$

where $\gamma_{aj}^X = \log \theta_{aj}^X$, $\alpha_0^X = \bar{\gamma}_{..}^X$, $\alpha_a^X = \bar{\gamma}_{a.}^X - \bar{\gamma}_{..}^X$, $\alpha_j^X = \bar{\gamma}_{.j}^X - \bar{\gamma}_{..}^X$ and $\alpha_{aj}^X = \gamma_{aj}^X - \bar{\gamma}_{a.}^X - \bar{\gamma}_{.j}^X + \bar{\gamma}_{..}^X$.

The log-linear representation satisfies the constraints:

$\sum_a \alpha_a^X = 0$, $\sum_j \alpha_j^X = 0$, $\sum_a \alpha_{aj}^X = \sum_j \alpha_{aj}^X = 0$. Analogous for θ_{aj}^Y .

The modelling process is focused on the relationship between α_{aj}^Y and α_{aj}^X . Marginal constraints such as $\sum_a \hat{Y}_{aj} = Y_{.j}$ for $j = 1, \dots, J$ and $\sum_j \hat{Y}_{aj} = Y_a$ for $a = 1, \dots, A$ can be considered using IPF without modifying the parameter estimates. Proxy information (not just covariates) is required.

Structure Preserving Models

In the context of SAE, the following Structure Preserving models have been used:

In the context of SAE, the following Structure Preserving models have been used:

1. Given $\{\theta_{aj}^X\}, \{Y_{1\cdot}, \dots, Y_{A\cdot}\}$:

In the context of SAE, the following Structure Preserving models have been used:

1. Given $\{\theta_{aj}^X\}, \{Y_{1.}, \dots, Y_{A.}\}$:

Synthetic Estimator: Adapted from Gonzalez & Hoza (1978),

$$\hat{Y}_{aj} = \theta_{aj}^X Y_{a.}$$

The underlying model is $\alpha_j^Y = \alpha_j^X, \alpha_{aj}^Y = \alpha_{aj}^X$.

The estimated composition is a rescaled version of the auxiliary composition.

2. Given $\{\theta_{aj}^X\}$, $\{Y_{1.}, \dots, Y_{A.}\}$, $\{Y_{.1}, \dots, Y_{.J}\}$:

2. Given $\{\theta_{aj}^X\}$, $\{Y_{1\cdot}, \dots, Y_{A\cdot}\}$, $\{Y_{\cdot 1}, \dots, Y_{\cdot J}\}$:

SPREE: Purcell & Kish (1980) use IPF to fit the two margins,

$$\hat{Y}_{aj}^{(1)} = \theta_{aj}^X Y_{a\cdot}, \quad \hat{Y}_{aj}^{(2)} = \frac{\hat{Y}_{aj}^{(1)}}{\hat{Y}_{\cdot j}^{(1)}} Y_{\cdot j}, \quad \dots$$

until convergency is achieved. This estimator minimizes the distance between the compositions X and \hat{Y} satisfying the marginal constraints. The underlining model is $\alpha_{aj}^Y = \alpha_{aj}^X$.

3. Given $\{\theta_{aj}^X\}$, $\{Y_{1.}, \dots, Y_{A.}\}$, $\{Y_{.1}, \dots, Y_{.J}\}$ and an estimated $\{\theta_{aj}^Y\}$:

3. Given $\{\theta_{aj}^X\}$, $\{Y_{1.}, \dots, Y_{A.}\}$, $\{Y_{.1}, \dots, Y_{.J}\}$ and an estimated $\{\theta_{aj}^Y\}$:

Generalized Linear Structural Model (GSPREE): Zhang & Chambers (2004) propose the model

$$\alpha_{aj}^Y = \beta \alpha_{aj}^X.$$

3. Given $\{\theta_{aj}^X\}$, $\{Y_{1.}, \dots, Y_{A.}\}$, $\{Y_{.1}, \dots, Y_{.J}\}$ and an estimated $\{\theta_{aj}^Y\}$:

Generalized Linear Structural Model (GSPREE): Zhang & Chambers (2004) propose the model

$$\alpha_{aj}^Y = \beta \alpha_{aj}^X.$$

β can be estimated using ML under the multinomial distribution, when expressing the model as:

$$\mu_{aj}^Y = \lambda_j + \beta \mu_{aj}^X$$

for $\mu_{aj} = \log \theta_{aj} - \frac{1}{J} \sum_I \log \theta_{aj} = \alpha_j + \alpha_{aj}$.

3. Given $\{\theta_{aj}^X\}$, $\{Y_{1..}, \dots, Y_{A..}\}$, $\{Y_{.1}, \dots, Y_{.J}\}$ and an estimated $\{\theta_{aj}^Y\}$:

Generalized Linear Structural Model (GSPREE): Zhang & Chambers (2004) propose the model

$$\alpha_{aj}^Y = \beta \alpha_{aj}^X.$$

β can be estimated using ML under the multinomial distribution, when expressing the model as:

$$\mu_{aj}^Y = \lambda_j + \beta \mu_{aj}^X$$

for $\mu_{aj} = \log \theta_{aj} - \frac{1}{J} \sum_I \log \theta_{aj} = \alpha_j + \alpha_{aj}$.

Given the sum-zero constraint of the α_j , the λ_j are nuisance parameters with no practical interest.

Extension of the Structure Preserving approach

All the previous models can be seen as particular cases of the more general model:

$$\begin{bmatrix} \alpha_{a1}^Y \\ \vdots \\ \alpha_{aJ}^Y \end{bmatrix} = \mathbb{B} \boldsymbol{\beta} \mathbb{B} \begin{bmatrix} \alpha_{a1}^X \\ \vdots \\ \alpha_{aJ}^X \end{bmatrix}$$

Extension of the Structure Preserving approach

All the previous models can be seen as particular cases of the more general model:

$$\begin{bmatrix} \alpha_{a1}^Y \\ \vdots \\ \alpha_{aJ}^Y \end{bmatrix} = \mathbb{B} \boldsymbol{\beta} \mathbb{B} \begin{bmatrix} \alpha_{a1}^X \\ \vdots \\ \alpha_{aJ}^X \end{bmatrix}$$

Where $\mathbb{B}_{J \times J} = \mathbb{I} - J^{-1}11'$ and $\boldsymbol{\beta}_{J \times J} = \{\beta_{jk}\}$ contains all the parameters.

Extension of the Structure Preserving approach

All the previous models can be seen as particular cases of the more general model:

$$\begin{bmatrix} \alpha_{a1}^Y \\ \vdots \\ \alpha_{aJ}^Y \end{bmatrix} = \mathbb{B} \beta \mathbb{B} \begin{bmatrix} \alpha_{a1}^X \\ \vdots \\ \alpha_{aJ}^X \end{bmatrix}$$

Where $\mathbb{B}_{J \times J} = \mathbf{I} - J^{-1} \mathbf{1} \mathbf{1}'$ and $\beta_{J \times J} = \{\beta_{jk}\}$ contains all the parameters.

The multiplication on left and right by \mathbb{B} ensure that the sum zero constraints are satisfied by the predicted α_{aj}^Y , as well as the uniqueness of $\mathbb{G} = \mathbb{B} \beta \mathbb{B}$. Denoting by $\{g_{jk}\}$ the components of \mathbb{G} we can write,

$$\alpha_{aj}^Y = \sum_k g_{jk} \alpha_{ak}^X.$$

As in the GSPREE, the estimation of β can be done using ML under the multinomial distribution writing the model as

$$\eta_a^Y = \lambda + \mathbb{B} \beta \mathbb{B} \eta_a^X.$$

Extension of the Structure Preserving approach

Some particular cases:

a) SPREE: $\beta = \mathbf{I}$

$$\alpha_{aj}^Y = \alpha_{aj}^X - \frac{1}{J} \sum_k \alpha_{ak}^X$$

Extension of the Structure Preserving approach

Some particular cases:

a) SPREE: $\beta = I$

$$\alpha_{aj}^Y = \alpha_{aj}^X - \frac{1}{J} \sum_k \alpha_{ak}^X$$

b) GSPREE: With parameter ϕ , $\beta = \phi I$

$$\alpha_{aj}^Y = \phi \alpha_{aj}^X - \phi \frac{1}{J} \sum_k \alpha_{ak}^X$$

Extension of the Structure Preserving approach

c) **Logit-Multinomial Model:** The model with $J-1$ parameters

$$\eta_{aj}^Y = \gamma_j + \phi_j \eta_{aj}^X$$

for $\eta_{aj} = \log [\theta_{aj}/\theta_{aJ}]$, can be written as a structural model in the form

$$\alpha_a^Y = \mathbb{B}_{(J)} \beta \mathbb{B} \alpha_a^X$$

for $\mathbb{B}_{(J)}$ the $J \times (J-1)$ matrix resulting of dropping the column J from \mathbb{B} and β a $(J-1) \times J$ matrix defined as

$$\beta = \left[\text{Diag} \left\{ \vec{\beta}_{(J)} \right\} \mid -\vec{\beta}_{(J)} \right]$$

for $\vec{\beta}_{(J)}$ a vector of $J-1$ parameters (The category J doesn't have a free parameter).

Extension of the Structure Preserving approach

d) GSPREE with category-specific (J) parameters:

$$\beta = \text{Diag} \left\{ \vec{\beta} \right\}$$

$$\alpha_{aj}^Y = \beta_j \alpha_{aj}^X - \frac{1}{J} \sum_k \beta_k \alpha_{ak}^X$$

The second term on the right hand, included to satisfy the sum-zero constrains without impose restrictions to the β_j , make the predictions of this model not a line anymore.

Extension of the Structure Preserving approach

d) GSPREE with category-specific (J) parameters:

$$\beta = \text{Diag} \left\{ \vec{\beta} \right\}$$

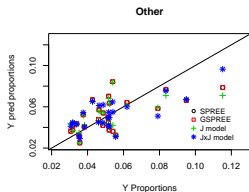
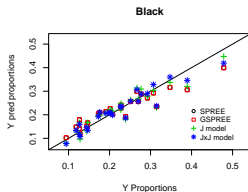
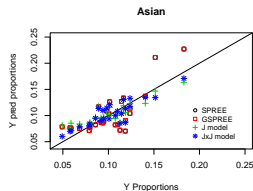
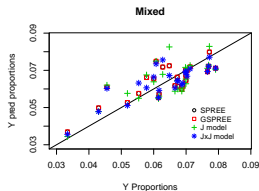
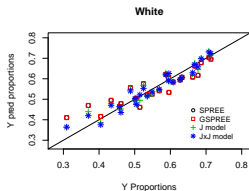
$$\alpha_{aj}^Y = \beta_j \alpha_{aj}^X - \frac{1}{j} \sum_k \beta_k \alpha_{ak}^X$$

The second term on the right hand, included to satisfy the sum-zero constrains without impose restrictions to the β_j , make the predictions of this model not a line anymore.

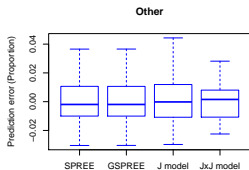
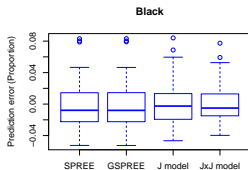
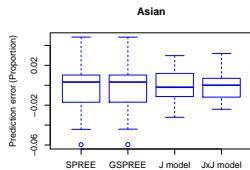
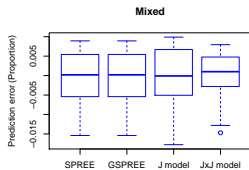
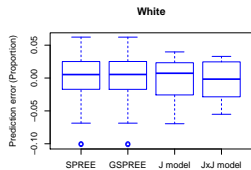
e) GSPREE JxJ model: $\beta = \{\beta_{jk}\}$, $G = \{g_{jk}\} = \mathbb{B}\beta\mathbb{B}$

$$\alpha_{aj}^Y = \sum_k g_{jk} \alpha_{ak}^X$$

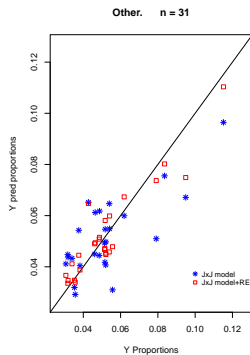
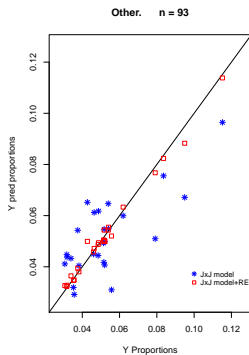
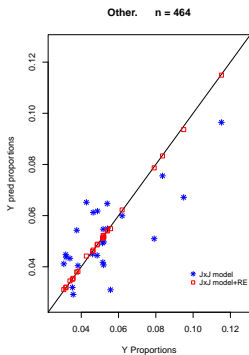
Data from 2001 and 2011 Population Census in England for the Hackney Borough. Ethnicity.



Data from 2001 and 2011 Population Census in England for the Hackney Borough. Ethnicity.



Extending the general model to include random effects



In summary...

- Having an auxiliary composition (register/census data), an estimated (updated) composition and the current margins, we extend the GSPREE model from one to a maximum of $J \times J$ parameters.

- Having an auxiliary composition (register/census data), an estimated (updated) composition and the current margins, we extend the GSPREE model from one to a maximum of $J \times J$ parameters.
- According to our preliminary exercises, the new models show less bias than SPREE and GSPREE models (fixed effects approach).

- Having an auxiliary composition (register/census data), an estimated (updated) composition and the current margins, we extend the GSPREE model from one to a maximum of $J \times J$ parameters.
- According to our preliminary exercises, the new models show less bias than SPREE and GSPREE models (fixed effects approach).
- We are still working on the extension to include cell-specific random effects. As expected, for a big sample size the estimative obtained using a mixed model gets closer to the direct estimate, however, as it is borrowing strength from the auxiliary composition, it would be more stable.

- Having an auxiliary composition (register/census data), an estimated (updated) composition and the current margins, we extend the GSPREE model from one to a maximum of $J \times J$ parameters.
- According to our preliminary exercises, the new models show less bias than SPREE and GSPREE models (fixed effects approach).
- We are still working on the extension to include cell-specific random effects. As expected, for a big sample size the estimative obtained using a mixed model gets closer to the direct estimate, however, as it is borrowing strength from the auxiliary composition, it would be more stable.
- MSE estimation is still need to be addressed.

- 1 Structure Preserving Models
- 2 Model using a Mapping matrix

Denote by $P = \{P_{ij}\}$ the gross flow from the composition X to Y, i.e., assume that for each area:

$$\begin{bmatrix} \theta_{a1}^Y \\ \vdots \\ \theta_{aJ}^Y \end{bmatrix} = \begin{bmatrix} P_{11} & \dots & P_{1J} \\ \vdots & \ddots & \vdots \\ P_{J1} & \dots & P_{JJ} \end{bmatrix} \begin{bmatrix} \theta_{a1}^X \\ \vdots \\ \theta_{aJ}^X \end{bmatrix}$$

The column sum of P is 1.

Denote by $P = \{P_{ij}\}$ the gross flow from the composition X to Y, i.e., assume that for each area:

$$\begin{bmatrix} \theta_{a1}^Y \\ \vdots \\ \theta_{aJ}^Y \end{bmatrix} = \begin{bmatrix} P_{11} & \dots & P_{1J} \\ \vdots & \ddots & \vdots \\ P_{J1} & \dots & P_{JJ} \end{bmatrix} \begin{bmatrix} \theta_{a1}^X \\ \vdots \\ \theta_{aJ}^X \end{bmatrix}$$

The column sum of P is 1.

What is the effect of the mapping matrix P in the log-linear representation of Y ?

$$\theta_a^Y = P\theta_a^X \quad \xrightarrow{?} \quad \log \theta_a^Y \approx M \log \theta_a^X$$

Using a First Order Taylor approximation over

$$\ln(\theta_{aj}^Y) = \ln\left(\sum_l \theta_{al}^X P_{jl}\right)$$

as function of $\ln(\theta_{aj}^X)$, around the distribution of X at an aggregate level, denoted by $\tilde{\theta}^X$, we obtain

$$\ln \theta_{aj}^Y - \ln \tilde{\theta}_j^Y \approx \sum_l (q_{jl} - \tau_j \tilde{\theta}_l^X) [\ln \theta_{al}^X - \ln \tilde{\theta}_l^X]$$

Using a First Order Taylor approximation over

$$\ln(\theta_{aj}^Y) = \ln\left(\sum_l \theta_{al}^X P_{jl}\right)$$

as function of $\ln(\theta_{aj}^X)$, around the distribution of X at an aggregate level, denoted by $\tilde{\theta}^X$, we obtain

$$\ln \theta_{aj}^Y - \ln \tilde{\theta}_j^Y \approx \sum_l (q_{jl} - \tau_j \tilde{\theta}_l^X) \left[\ln \theta_{al}^X - \ln \tilde{\theta}_l^X \right]$$

where $\tilde{\theta}^Y = P\tilde{\theta}^X$, $q_{jl} = \frac{P_{jl}\tilde{\theta}_l^X}{\tilde{\theta}_j^Y}$ is the reverse flow and $\tau_j = \frac{P_{jr_j}}{\tilde{\theta}_j^Y}$ for P_{jr_j} one cell specifically chosen for the j category.

Applying the link function

$$\mu_{aj} = \log \theta_{aj} - \frac{1}{J} \sum_l \log \theta_{aj} = \alpha_j + \alpha_{aj}$$

Applying the link function

$$\mu_{aj} = \log \theta_{aj} - \frac{1}{J} \sum_l \log \theta_{aj} = \alpha_j + \alpha_{aj}$$

we can obtain the relationship

$$\alpha_{aj}^Y \approx \sum_l (q_{jl} - \bar{q}_{.l}) \alpha_{al}^X - \sum_l (\tau_j - \bar{\tau}) \tilde{\theta}_l^X \alpha_{al}^X.$$

According to our empirical studies, the leading term in the approximation is the term associated with the reverse flow. In this sense, a model involving also auxiliary information on the reverse flow at an aggregate level could be of interest.

Thanks!