# Model based clustering using multivariate mixed type panel data

## Arnošt Komárek[1], Jan Vávra[1], Vladislav Bína[2]

[1]Faculty of Mathematics and Physics, Charles University, Prague, Czech Republic

[2]Faculty of Management, University of Economics, Prague, Czech Republic

komarek@karlin.mff.cuni.cz, vavraj@karlin.mff.cuni.cz, vladislav.bina@vse.cz
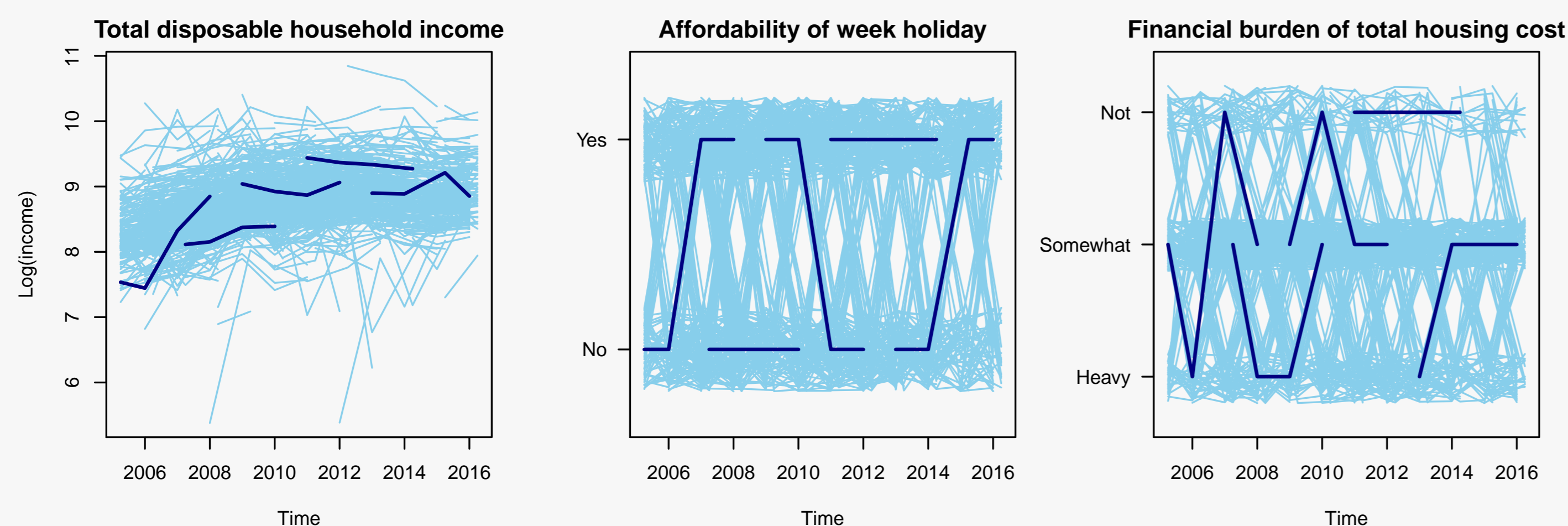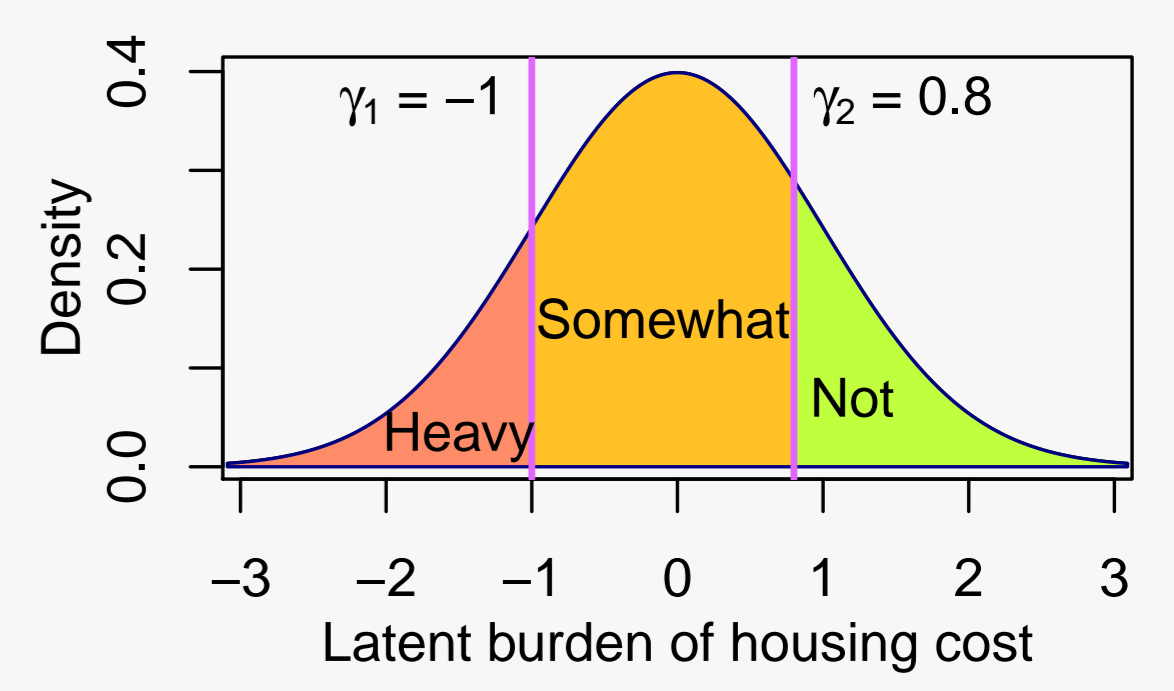
## EU-SILC dataset - mixed type data

▷ EU-SILC = European Union - Statistics on Income and Living Conditions

▷ Longitudinal multidimensional data on income, poverty, social exclusion and living conditions measured on private households

▷ Annually gathered data via questionnaires targeted on both households and individuals living there

▷ Available data: $n = 29\,292$ households from the Czech Republic (years $2005 - 2016$)

▷ Outcomes
- ▹ Numeric - Total disposable household income, ...
- ▹ Binary - Affordability of week holiday away from home, ... } mixed type data
- ▹ Ordinal - Financial burden of total housing cost, ...

▷ Explanatory variables:
- ▹ year, region, level of urbanization, dwelling type, weighted family size, ...



## Research goals

▷ To discover unobserved heterogeneity in various socio-economic characteristics.

▷ To identify hidden groups of similar longitudinal evolution of these characteristics.

▷ To partition households into these groups to determine the level of social-economic status.

▷ To construct a set of general rules for classification of households.

▷ To uncover poverty and social exclusion temporal patterns.



## Notation

▷ household $i \in \{1, \ldots, n\}$,   visit number $j \in \{1, \ldots, n_i, \}$,   outcome $r \in \{1, \ldots, R\}$

▷ $Y_{i,j}^r$ - measured value of an outcome $\rightsquigarrow Y_i^r, Y_i$

▷ $t_{i,j}$ - time of the measurement $\rightsquigarrow t_i$

▷ $C_i$ - all explanatory information known to household $i$

## Model based clustering (Banfield and Raftery, 1993)

▷ $K$: number of unobserved groups (initially assumed to be known)

▷ $U_i \in \{1, \ldots, K\}$: latent indicators of membership to one of the $K$ groups for each household $i$

▷ $0 < w_k = P[U_i = k]$: unknown probability of group $k \in \{1, \ldots, K\}$

▷ $f_k(y_i; C_i, \psi, \psi^{(k)})$: PDF of the probabilistic model for $Y_i$ when household $i$ belongs to group $k$
- ▹ $\psi$: parameters of the probabilistic model that are common to all groups
- ▹ $\psi^{(k)}$: group-specific parameters of model for group $k$

▷ $\theta = (w, \psi, \psi^{(1)}, \ldots, \psi^{(K)})$: unknown parameters of interest

▷ Mixture likelihood:
$$L(\theta) = \prod_{i=1}^{n} \sum_{k=1}^{K} w_k f_k(Y_i; C_i, \psi, \psi^{(k)})$$

▷ By Bayes rule:
$$p_{i,k}(\theta) = P[U_i = k | Y_i = y_i; C_i, \theta] = \frac{w_k f_k(y_i; C_i, \psi, \psi^{(k)})}{\sum_{\ell=1}^{K} w_\ell f_\ell(y_i; C_i, \psi, \psi^{(\ell)})}$$

▷ Classification rule:
$$\widehat{U}_i := k \iff k = \arg\max_{\ell \in \{1, \ldots, K\}} \widehat{p_{i,\ell}}(\theta)$$

▷ Estimation: Bayesian approach and MCMC methods (Gibbs sampling)
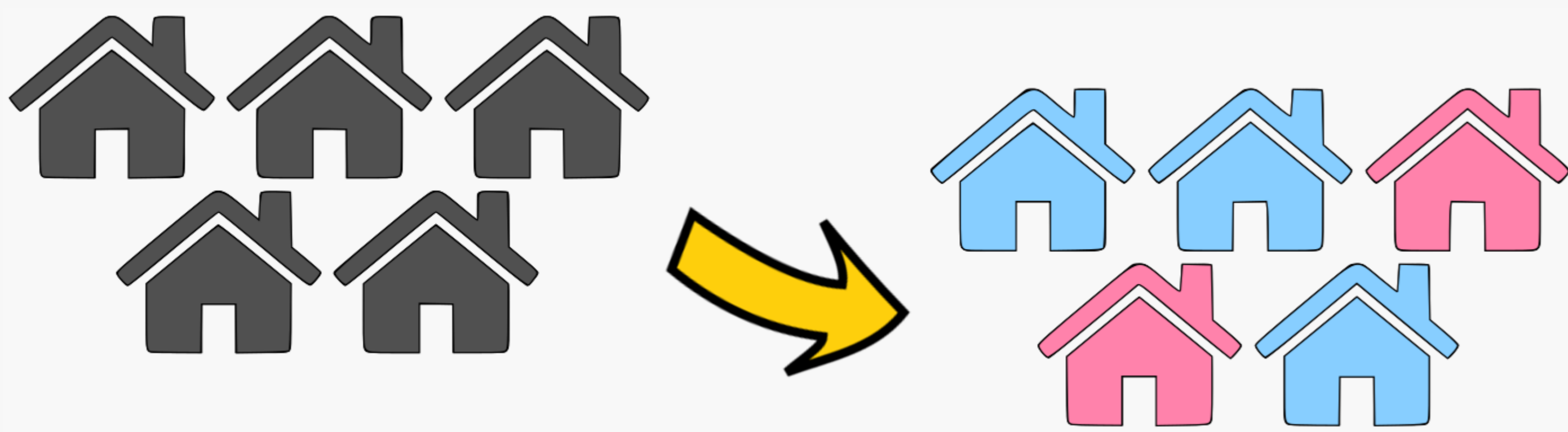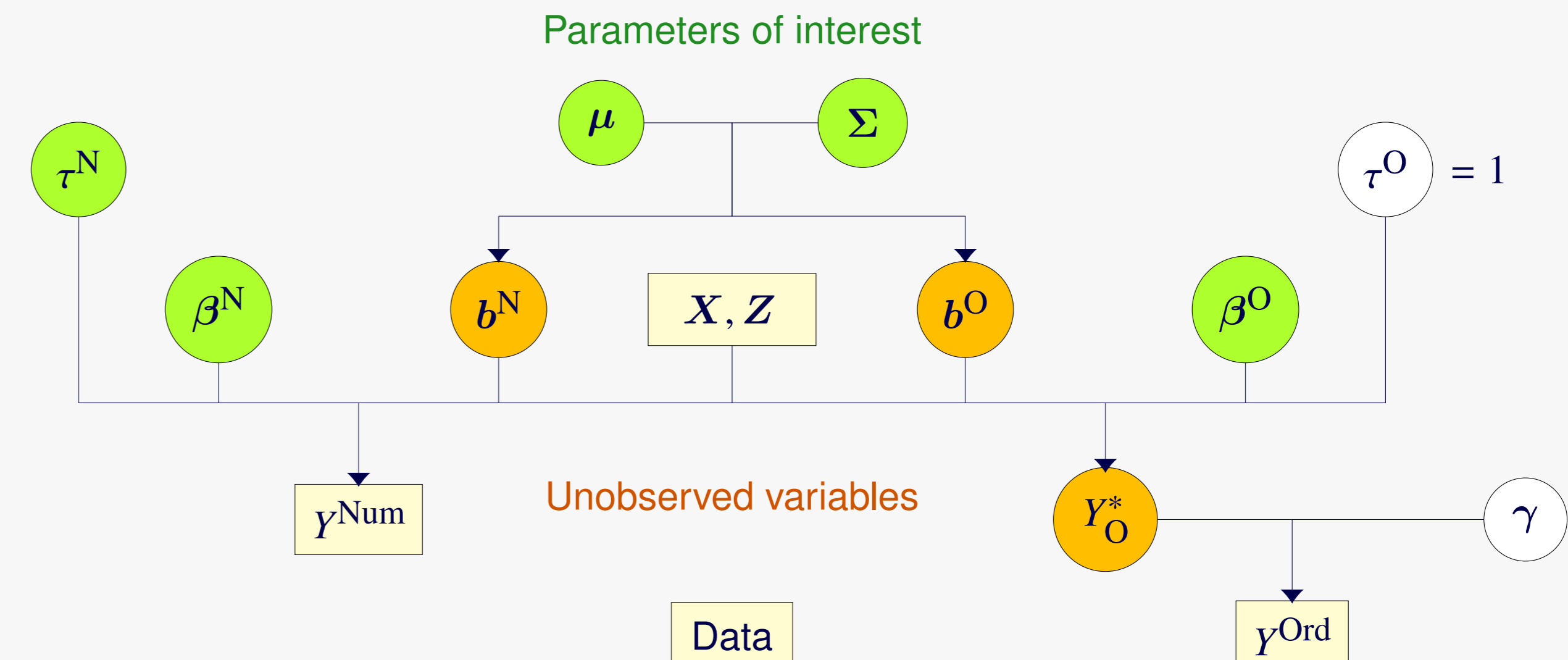
▷ Software: implemented in ®R using programming language C

## Models for single outcomes (in certain group)

▷ Numeric - Linear Mixed-effects Model (LMM)

- ▹ Model formula:
$$Y_{i,j}^N \sim N\left(\left(X_{i,j}^N\right)^\top \beta^N + \left(Z_{i,j}^N\right)^\top b_i^N, \ \left(\tau^N\right)^{-1}\right)$$

- ▹ Random effects:
$$b_i^N \overset{iid}{\sim} N\left(\mu^N, \Sigma^{NN}\right)$$

- ▹ For example: $\log\left(Y_{i,j}^N\right) = b_i^N + \beta^N \cdot t_{i,j} + \varepsilon_{i,j}^N$

▷ Binary + Ordinal - thresholded latent numeric variable following LMM

- ▹ $Y_{i,j}^O = \ell \in \{1, \ldots, L\}$  ($L$ ordered levels)
- ▹ $Y_{i,j}^{O,\star} \sim$ LMM with $\beta^O, b_i^O$ and $\tau^O = 1$
- ▹ Observed $Y_{i,j}^O$ determined by set of thresholds $\gamma$

$$-\infty = \gamma_0 < -1 = \gamma_1 < \gamma_2 < \cdots < \gamma_{L-1} < \gamma_L = \infty$$

$$Y_{i,j}^O = \text{Heavy} \iff Y_{i,j}^{O,\star} \leq \gamma_1$$
$$Y_{i,j}^O = \text{Somewhat} \iff \gamma_1 < Y_{i,j}^{O,\star} \leq \gamma_2$$
$$Y_{i,j}^O = \text{Not} \iff \gamma_2 < Y_{i,j}^{O,\star}$$

- ▹ In general  $Y_{i,j}^O = \ell \iff \gamma_{\ell-1} < Y_{i,j}^{O,\star} \leq \gamma_\ell$



## Joint modelling (in certain group)

▷ Outcomes cannot be considered to be independent of each other.

▷ Individual models are joined through joint distribution of random effects:
$$b_i = \begin{pmatrix} b_i^N \\ b_i^O \end{pmatrix} \overset{iid}{\sim} N\left(\mu = \begin{pmatrix} \mu^N \\ \mu^O \end{pmatrix}, \Sigma = \begin{pmatrix} \Sigma^{NN} & \Sigma^{NO} \\ \Sigma^{ON} & \Sigma^{OO} \end{pmatrix}\right)$$

▷ Leading to a certain dependency structure:
$$\text{var}\left[\begin{pmatrix} Y_{i,j}^N \\ Y_{i,j}^{O,\star} \end{pmatrix} \Big| C_{i,j}\right] = \begin{pmatrix} \left(\tau^N\right)^{-1} + \left(Z_{i,j}^N\right)^\top \Sigma^{NN} Z_{i,j}^N & \left(Z_{i,j}^N\right)^\top \Sigma^{NO} Z_{i,j}^O \\ \left(Z_{i,j}^O\right)^\top \Sigma^{ON} Z_{i,j}^N & 1 + \left(Z_{i,j}^O\right)^\top \Sigma^{OO} Z_{i,j}^O \end{pmatrix}$$

▷ Group-specific parameters:  $\psi^{(k)} = \left(\beta^{(k)}, \mu^{(k)}, \Sigma^{(k)}\right)$

## Hierarchical Bayesian joint model for numeric and ordinal variable



## Classified households

## References

**Banfield, J., D. and Raftery, A., E.** (1993). Model-Based Gaussian and Non-Gaussian Clustering. *Biometrics*, **49**(3), 803–821.

**R Core Team** (2019). *R: A Language and Environment for Statistical Computing.* R Foundation for Statistical Computing, Vienna, Austria. URL: http://www.R-project.org.