

# Estimating Population Size from Register Data: a Capture-Recapture Approach

**Lucy Brown**<sup>1</sup>: lyb3@kent.ac.uk

Eleni Matechou<sup>1</sup>

Bruno Santos<sup>2</sup>

Eleonora Mussino<sup>3</sup>

*1 University of Kent, UK*

*2 University of Lisbon, Portugal*

*3 Stockholm University, Sweden*



# Motivating Case Study

## Overcoverage:

- Due to imperfect emigration and/or death registration
- Leads to serious bias in population estimates
- Negatively influences policy-making and research

## Overcoverage Estimation:

- Existing approaches<sup>[1][2]</sup> rely on multiple systems estimation (MSE) and only consider annual snapshots of the register data
- Instead, we have employed a longitudinal approach, following individuals, and hence registers, over different years

[1] Andrea Monti, Sven Drefahl, Eleonora Mussino, and Juho Härkönen. Over-coverage in population registers leads to bias in demographic estimates. *Population Studies*, 74(3):451–469, 2020

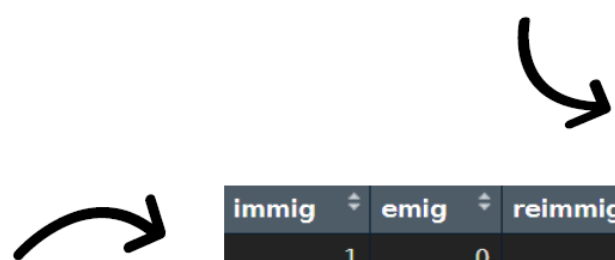
[2] Eleonora Mussino, Bruno Santos, Andrea Monti, Eleni Matechou, and Sven Drefahl. Multiple systems estimation for studying over-coverage and its heterogeneity in population registers. *Quality & Quantity*, pages 1–24, 2023

# Swedish Register Data

- Provided by the Swedish National Institute of Statistics (Statistics Sweden - SCB)
- All foreign-born residents who first entered Sweden between the years 2003 and 2016 as adults

(This is not real data)

**Nine Registers**



**Emigration, (Re-)Immigration and Death records**

	immig	emig	reimmig	marr	div	swecit
	1	0	0	0	0	0
	0	1	0	0	0	0
	1	0	1	0	0	0
	0	0	0	1	0	0
	0	1	0	0	0	0
	0	0	0	0	0	0
	0	0	0	0	0	0

	emp	stud	intmove	faminc	amf	child
	0	0	0	0	1	0
	1	0	0	1	0	0
	1	0	0	1	0	0
	1	0	0	1	0	0
	1	0	0	1	0	0

	id	year	age	cob	sex	death
1	XXX1	2003	27	France	1	0
2	XXX1	2004	28	France	1	0
3	XXX1	2007	31	France	1	0
4	XXX1	2008	32	France	1	0
5	XXX1	2009	33	France	1	0
6	XXX2	2014	68	India	0	0
7	XXX2	2015	69	India	0	1

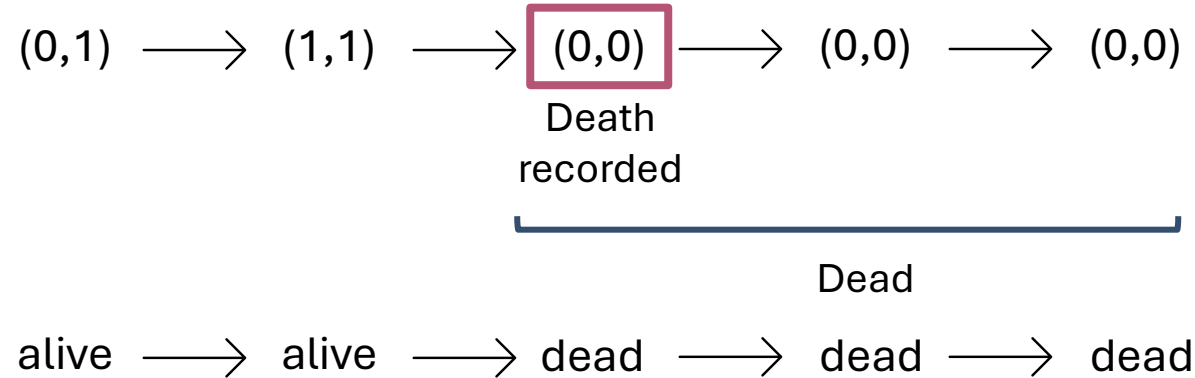
**Covariates (treated as categorical):**

- Sex
- Country of birth
- Age
- Time Since first entering Sweden

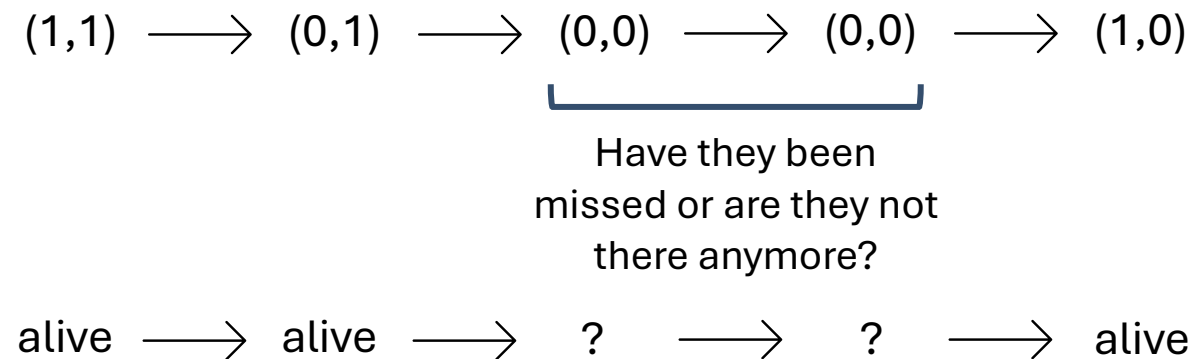
# Observations

An individual is observed on some **combination of R registers**, on which they can be observed (1) or not observed (0), i.e. 2 possible outcomes for each register  $\Rightarrow 2^R$  register combinations

Example Person #1



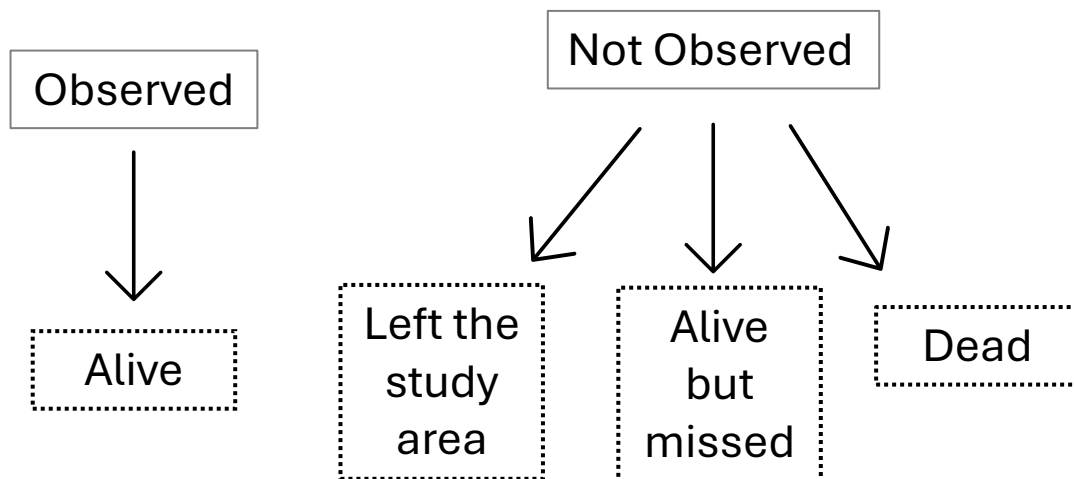
Example Person #2



# Capture-Recapture Models

## Capture Recapture (CR) Models:

1. An initial capture is made
2. Unmarked individuals are marked in some unique way (tags, rings etc)
3. Individuals are released back into the population
4. Subsequent captures are made, allowing individuals to be tracked over time



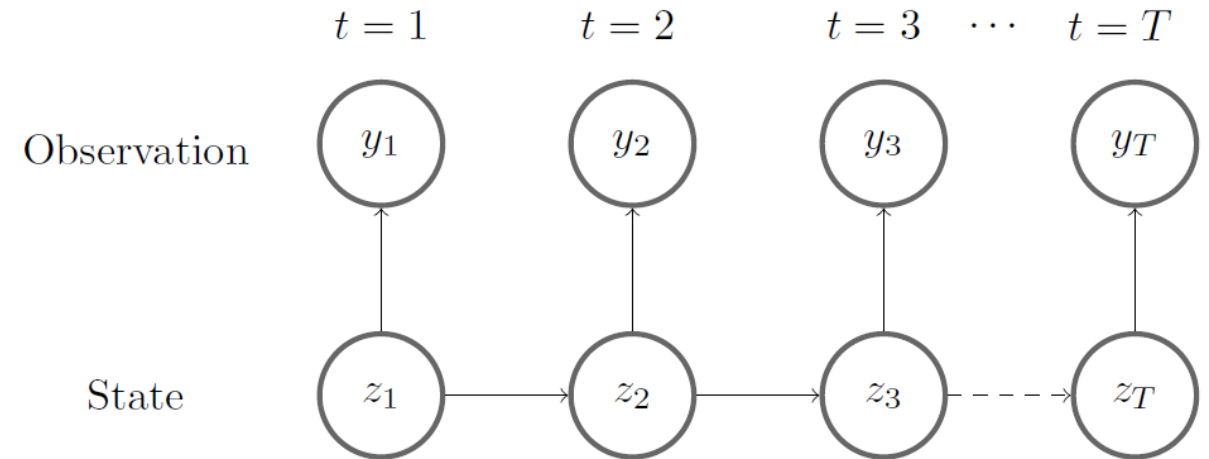
- CR models are commonly used in ecology to estimate population size for wild animals
- An individual's true state  $z_t$  is unobservable (latent)

Kenneth H. Pollock. Capture-recapture models. *Journal of the American Statistical Association*, 95(449):293–296, 2000. ISSN 01621459, 1537274X.  
URL <http://www.jstor.org/stable/2669550>

# HMM Formulation

J. L. Laake. Capture-recapture analysis with hidden markov models. 2013.

- Hidden Markov model (HMM) formulation marginalises over the latent state to calculate the marginal likelihood
  - The marginal likelihood is the probability of the observed data and the latent states, given the model parameters
- Computationally expensive so use the Forward Algorithm
- This iterative method is very efficient and allows latent states to be inferred



**Probability of history at the previous time point** (green text) points to  $\alpha_{t-1}$  in the equation.

**Probability of transitioning to each new state given the state at  $t - 1$**  (blue text) points to  $\Gamma_{t-1}$  in the equation.

**Observation probabilities dependent on the current state** (red text) points to  $P(y_t)$  in the equation.

$$\alpha_1 = \delta P(y_1)$$

$$\alpha_t = \alpha_{t-1} \Gamma_{t-1} P(y_t)$$

$$\Rightarrow L_T = \alpha_T \mathbf{1}'$$

# Simple CR Model

$$\Gamma = \begin{bmatrix} \phi & 1 - \phi \\ 0 & 1 \end{bmatrix} \begin{matrix} z_t = \text{alive} \\ z_t = \text{dead} \end{matrix}$$

- States  $z_t$ : alive (1) or dead (0)
- Observations  $y_t$ : observed (1) or not observed (0)
- Probability of survival:  $\phi$
- Probability of observation:  $p$
- All individuals are alive and observed when initially captured, i.e. initial state  $\delta = [1,0]$

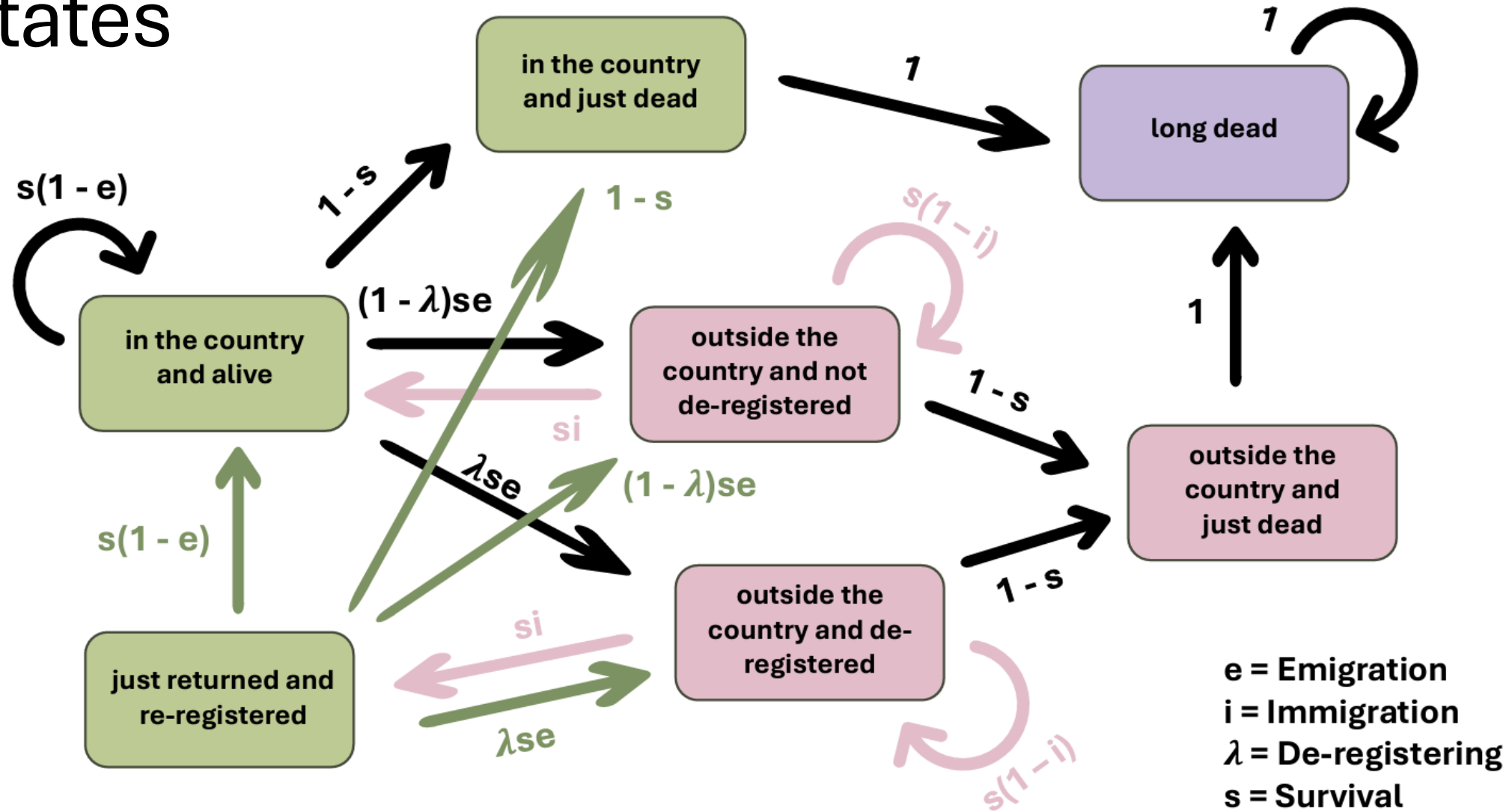
$$P(0) = \begin{bmatrix} 1 - p & 0 \\ 0 & 1 \end{bmatrix} \quad P(1) = \begin{bmatrix} p & 0 \\ 0 & 0 \end{bmatrix}$$

If an individual has capture history  $y_t = [1,0,1]$ , they have a likelihood:

$$L_T = \delta P(y_1) \Gamma P(y_2) \Gamma P(y_3) 1'$$

$$L_T = [1,0] \begin{bmatrix} p & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} \phi & 1 - \phi \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 - p & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} \phi & 1 - \phi \\ 0 & 1 \end{bmatrix} \begin{bmatrix} p & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

# States





# Transition Matrix

$$\begin{bmatrix} s(1-e) & 1-s & \lambda se & (1-\lambda)se & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & s(1-i) & 0 & 1-s & si & 0 \\ si & 0 & 0 & s(1-i) & 1-s & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 \\ s(1-e) & 1-s & \lambda se & (1-\lambda)se & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

1. In the country and alive
2. In the country and just dead
3. Outside the country and de-registered
4. Outside the country but didn't de-register

5. Outside the country and just dead
6. Just returned to the country having de-registered
7. Long dead

# Model Coefficients Specification

- Parameters in the transition matrix are specified using logistic regression:

$$\text{logit}(\theta) = \beta_0 + \underbrace{\beta_1 C_1 + \beta_2 C_2 + \dots}_{\text{Individual- and time-dependent covariates } C} + \underbrace{\epsilon}_{\text{Random effects}} \quad \text{where } \epsilon \sim N(0, \sigma^2)$$

- Parameters in the observation matrix are specified using multinomial regression<sup>[4]</sup>:

$$X = \begin{bmatrix} 1 & \dots & 1 & 1 & \dots & 1 & 1 & \dots & 1 \\ 1 & \dots & 1 & 1 & \dots & 0 & 1 & \dots & 0 \\ \vdots & & \vdots & & \vdots & & \vdots & & \vdots \\ 0 & \dots & 0 & 0 & \dots & 0 & 0 & \dots & 0 \end{bmatrix} \quad p_j = \frac{\exp(\gamma X_{j*} + \epsilon_j)}{1 + \sum_h \exp(\gamma X_{h*} + \epsilon_h)} \quad \text{where } \epsilon \sim N(0, \sigma^2)$$

$j = 1, \dots, J; \quad h = 1, \dots, J - 1$

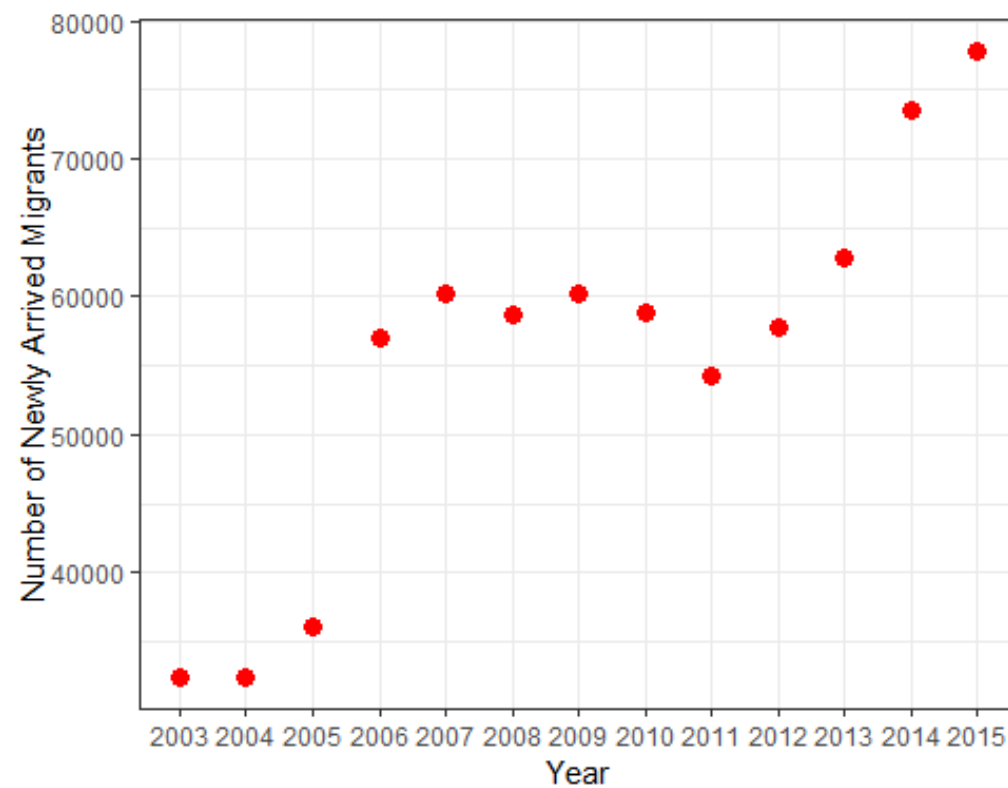
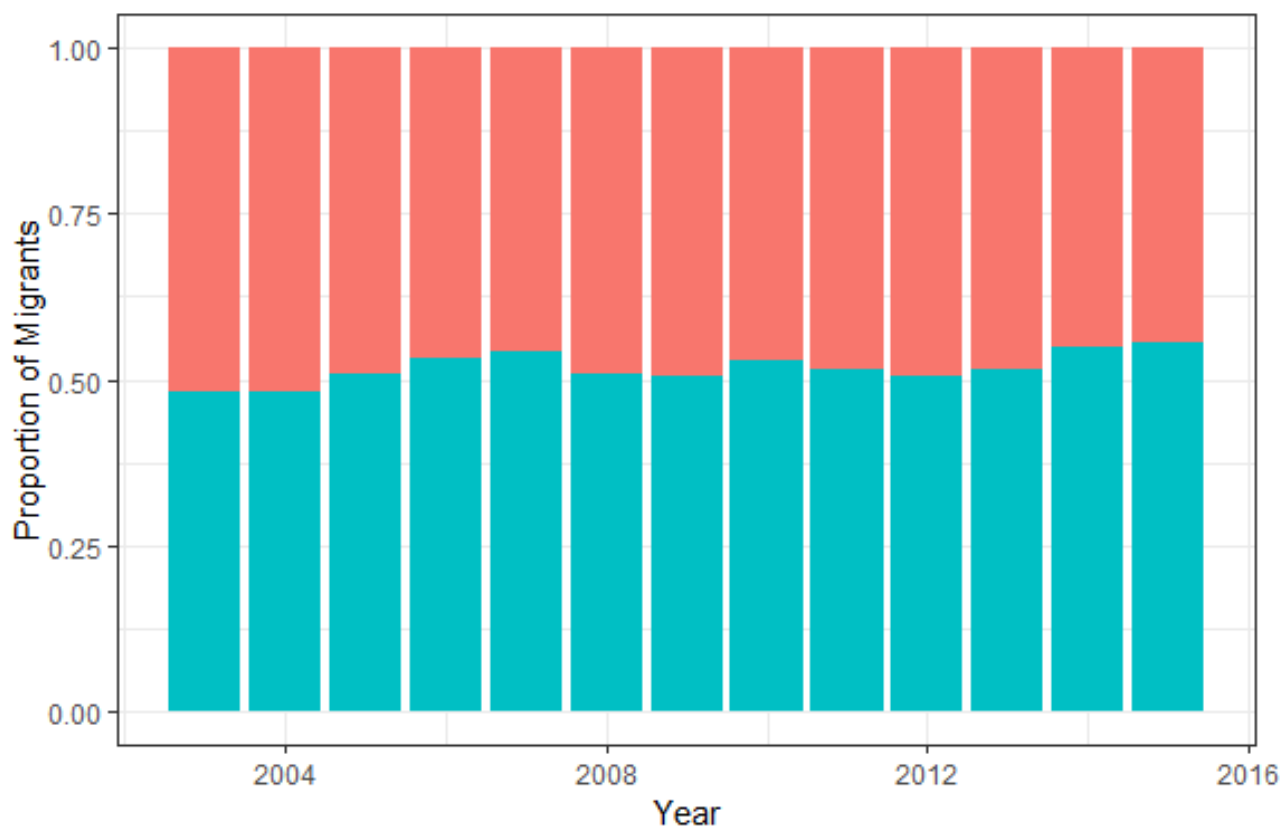
$\underbrace{\hspace{10em}}_{\text{Registers}}$ 
 $\underbrace{\hspace{10em}}_{\text{Covariates}}$ 
 $\underbrace{\hspace{10em}}_{\text{Interactions}}$

where  $p_j$  are the observation probabilities of each register combination

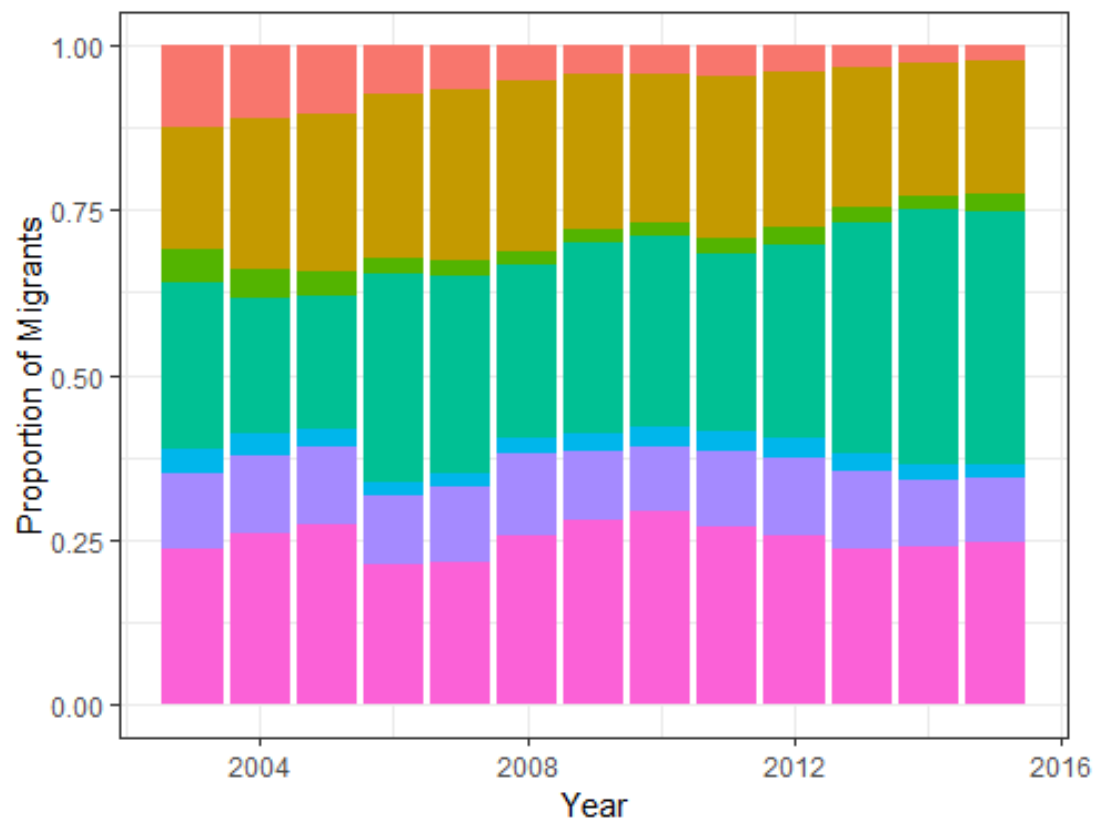
[4] A. Agresti. An introduction to categorical data analysis. John Wiley & Sons, 2nd edition, 2007.

# Results

Distribution of sex for newly arrived migrants



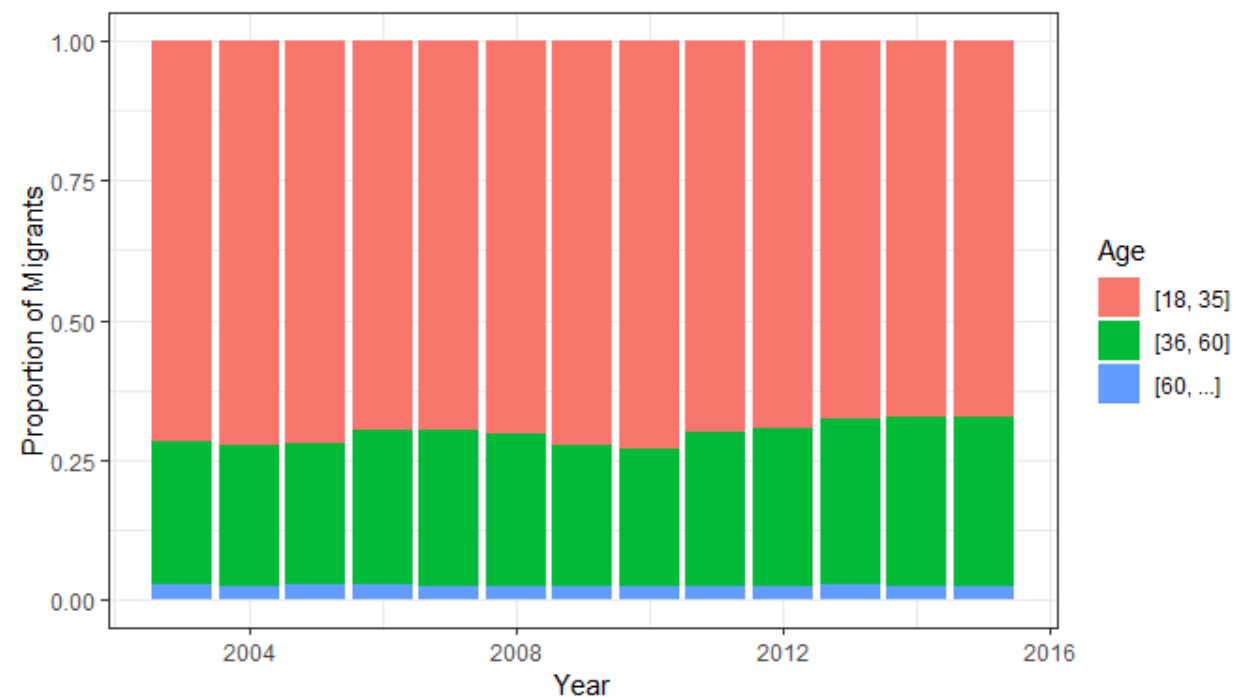
Number of newly arrived migrants



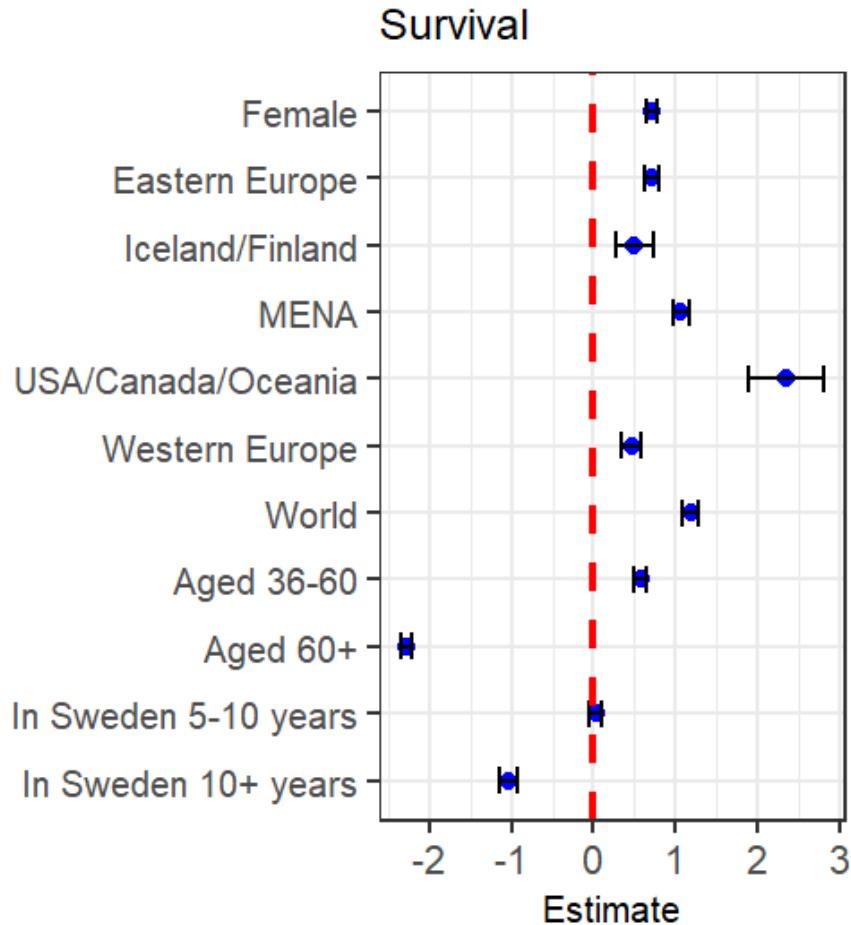
Distribution of country of birth  
for newly arrived migrants



Distribution of age for newly  
arrived migrants



# Survival Probability



**Baseline:** men aged 18-35, born in Denmark/Norway who first entered Sweden less than 5 years ago

$$\text{logit}(\lambda) = 5.617 \Rightarrow \lambda = 0.996$$

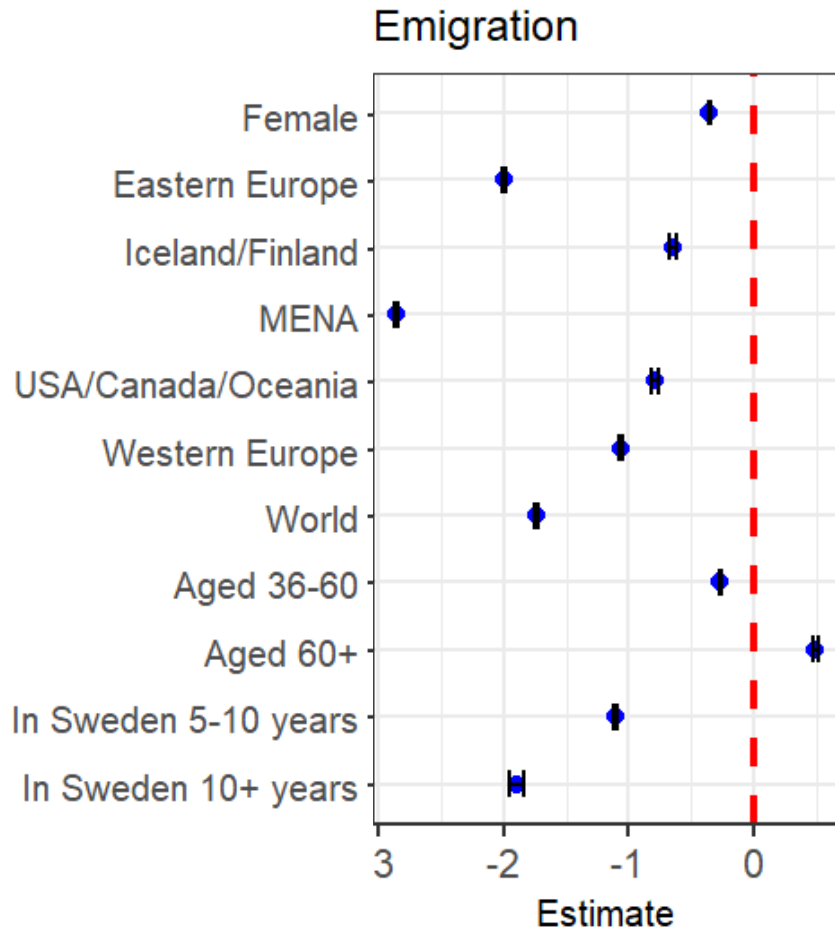
**Women** aged **36-60**, born in **USA/Canada/Oceania** who first entered Sweden **5-10 years ago**

$$\text{logit}(\lambda) = 5.617 + 0.710 + 2.340 + 0.572 + 0.023 \Rightarrow \lambda = 0.9999$$

Men **aged 60+**, born in Denmark/Norway who first entered Sweden **10+ years ago**

$$\text{logit}(\lambda) = 5.617 - 2.272 - 1.028 \Rightarrow \lambda = 0.910$$

# Emigration Probability



**Baseline:** men aged 18-35, born in Denmark/Norway who first entered Sweden less than 5 years ago

$$\text{logit}(\lambda) = -0.427 \Rightarrow \lambda = 0.395$$

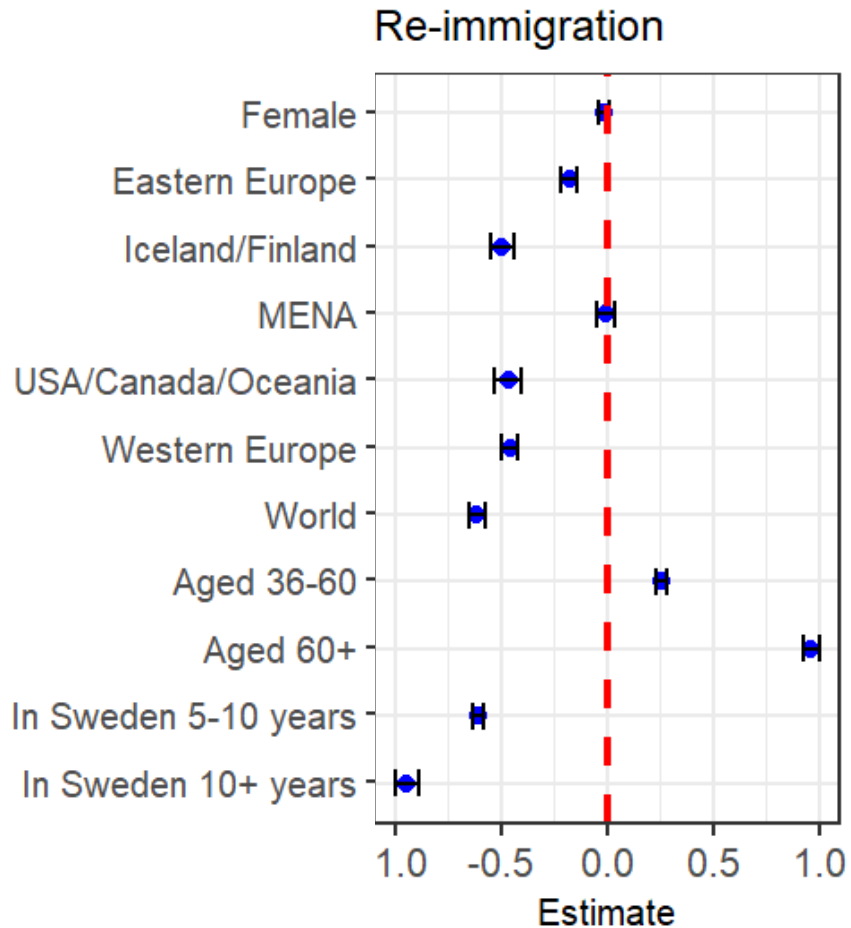
Men aged **60+**, born in Denmark/Norway who first entered Sweden less than 5 years ago

$$\text{logit}(\lambda) = -0.427 + \mathbf{0.485} \Rightarrow \lambda = 0.515$$

**Women** aged **36-60**, born in **MENA** who first entered Sweden **10+ years ago**

$$\text{logit}(\lambda) = -0.427 - \mathbf{0.356} - \mathbf{2.855} - \mathbf{0.273} - \mathbf{1.899} \Rightarrow \lambda = 0.00299$$

# Re-Immigration Probability



**Baseline:** men aged 18-35, born in Denmark/Norway who first entered Sweden less than 5 years ago

$$\text{logit}(\lambda) = -2.978 \Rightarrow \lambda = 0.0532$$

Men aged **60+**, born in Denmark/Norway who first entered Sweden less than 5 years ago

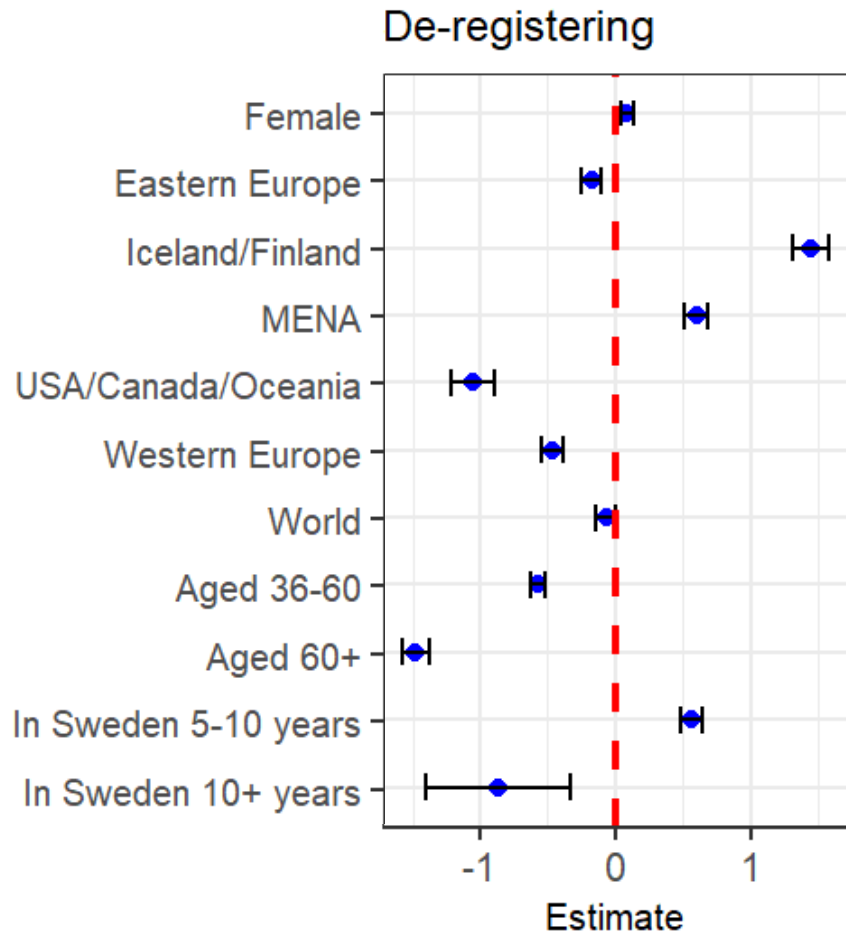
$$\text{logit}(\lambda) = -2.978 + \mathbf{0.963} \Rightarrow \lambda = 0.128$$

**Women** aged 18-35, born in **the rest of the World** who first entered Sweden **10+ years ago**

$$\text{logit}(\lambda) = -2.978 - 0.0159 - 0.613 - 0.946 \Rightarrow \lambda = 0.0115$$



# De-registering Probability



**Baseline:** men aged 18-35, born in Denmark/Norway who first entered Sweden less than 5 years ago

$$\text{logit}(\lambda) = -0.332 \Rightarrow \lambda = 0.418$$

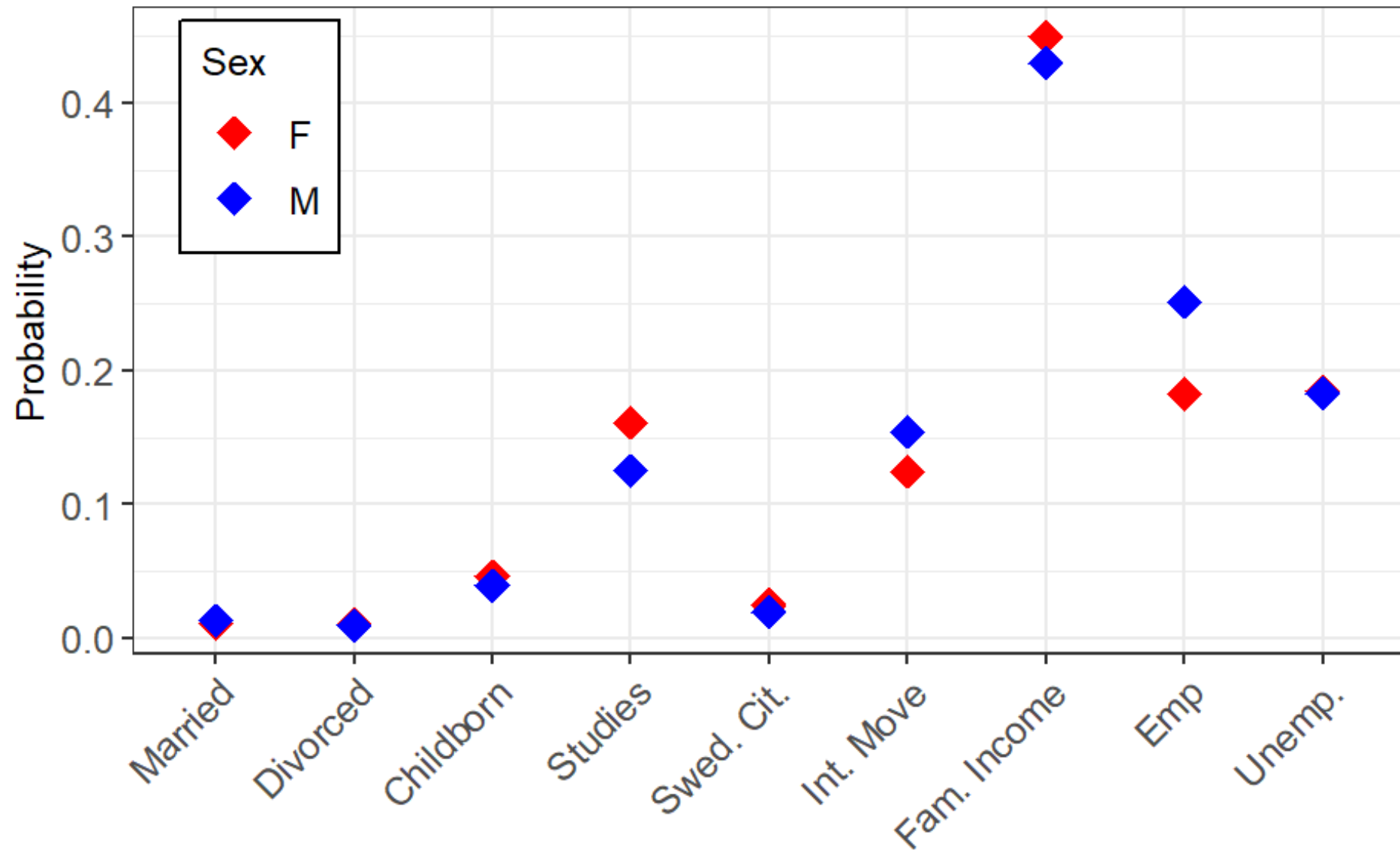
**Women** aged 18-35, born in **Iceland/Finland** who first entered Sweden **5-10 years ago**

$$\text{logit}(\lambda) = -0.332 + 0.088 + 1.452 + 0.568 \Rightarrow \lambda = 0.855$$

Men **aged 60+**, born in **USA/Canada/Oceania** who first entered Sweden **10+ years ago**

$$\text{logit}(\lambda) = -0.332 - 1.053 - 1.474 - 0.865 \Rightarrow \lambda = 0.024$$

## Register Observation Probability by Sex



# Register Combination Probabilities

Highest Probability Register Combinations	Male	Female
Family Income and Employment	0.0965	0.0762
Family Income	0.0368	0.0611
Family Income and Active Unemployment	0.0305	0.0388

Lowest Probability Register Combinations	Male	Female
All except Internal Move and Family Income	$3.6107 \times 10^{-12}$	$2.3061 \times 10^{-12}$
All except Family Income	$6.2981 \times 10^{-12}$	$2.3061 \times 10^{-12}$
All except Family Income and Active Unemployment	$1.0539 \times 10^{-11}$	$6.1599 \times 10^{-12}$

Probability of being unobserved

- Male: 0.0251
- Female: 0.0220



Probability of being observed in at least one register

- Male: 0.9749
- Female: 0.9780

# Future Work

- ★ Calculate population size estimates using the Viterbi Algorithm
- ★ Considering further extensions such as dependence between family units
- ★ Application to equivalent data from Norway, provided by Statistics Norway
  - Additional complexity that checks are regularly done and individuals manually removed if not present in the country

# Thank you for your attention!

Email: [lyb3@kent.ac.uk](mailto:lyb3@kent.ac.uk)

