

# The number of Latvian residents estimate via logistic regression

Jelena Valkovska

Central Statistical Bureau of Latvia, e-mail: Jelena.Valkovska@csb.gov.lv

## Introduction

The main aim of this research is to investigate the logistic regression as the potential instrument for estimating the number of residents and to get such estimates of beta coefficients, that can be used to estimate the number of Latvian residents and number of emigrants now and in the future.

## Logistic regression

The aim of an analysis using logistic regression is the same as that of any model-building technique used in statistic: to find the best fitting and most parsimonious model to describe the relationship between dependent variable and a set of independent variables. What distinguishes logistic regression from the linear regression model is that the outcome variable in logistic regression is dichotomous.

$$\pi(x) = E(Y|x) = \frac{e^{\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_t x_t}}{1 + e^{\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_t x_t}} \quad (1)$$

The method of estimation used is the maximum likelihood. The main principle of this method is that is used as the estimate of  $\beta$  the value which maximizes the expression:

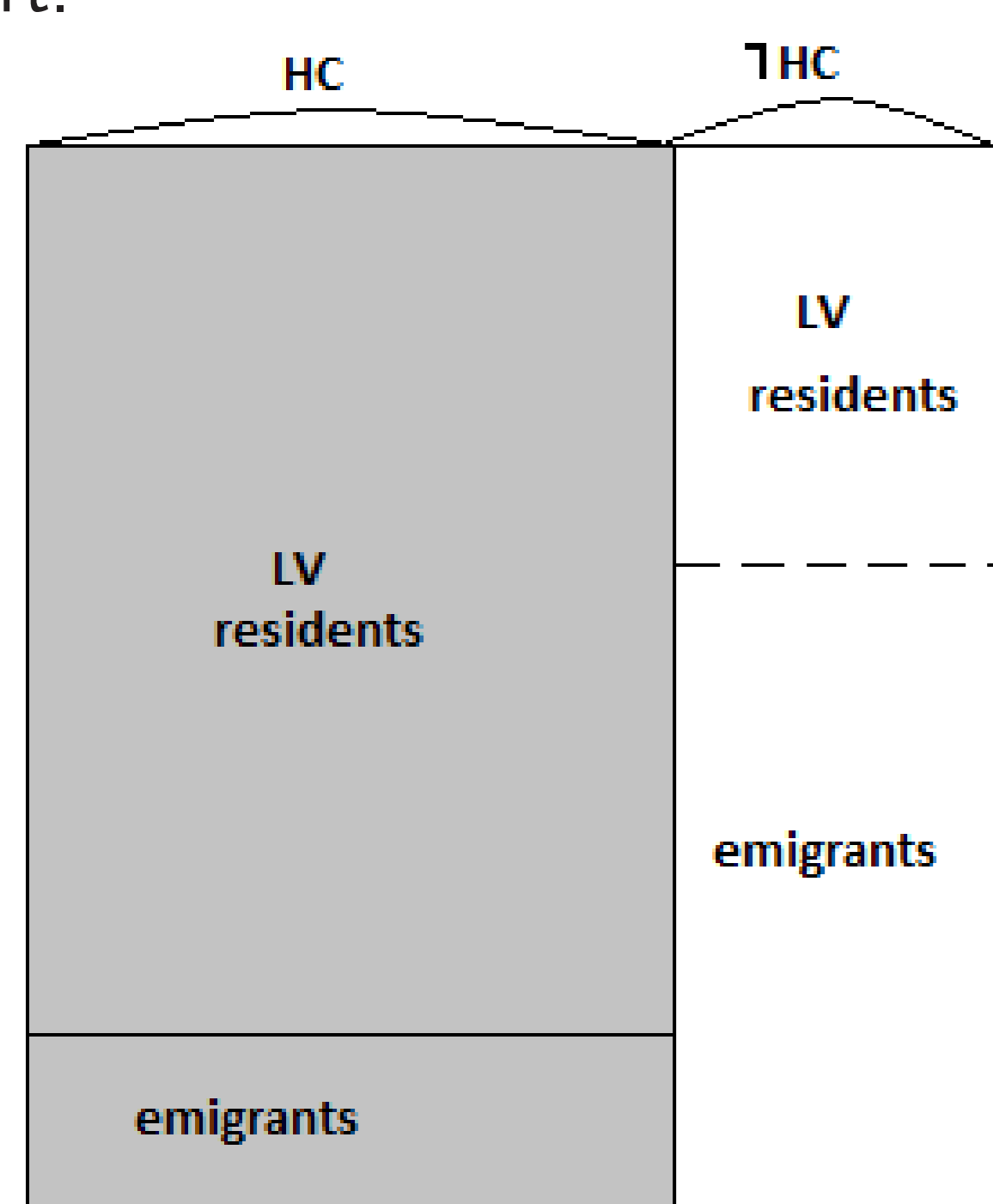
$$l(\beta) = \prod_{i=1}^n \pi(x_i)^{y_i} [1 - \pi(x_i)^{(1-y_i)}] \quad (2)$$

## Population and Housing Census 2011

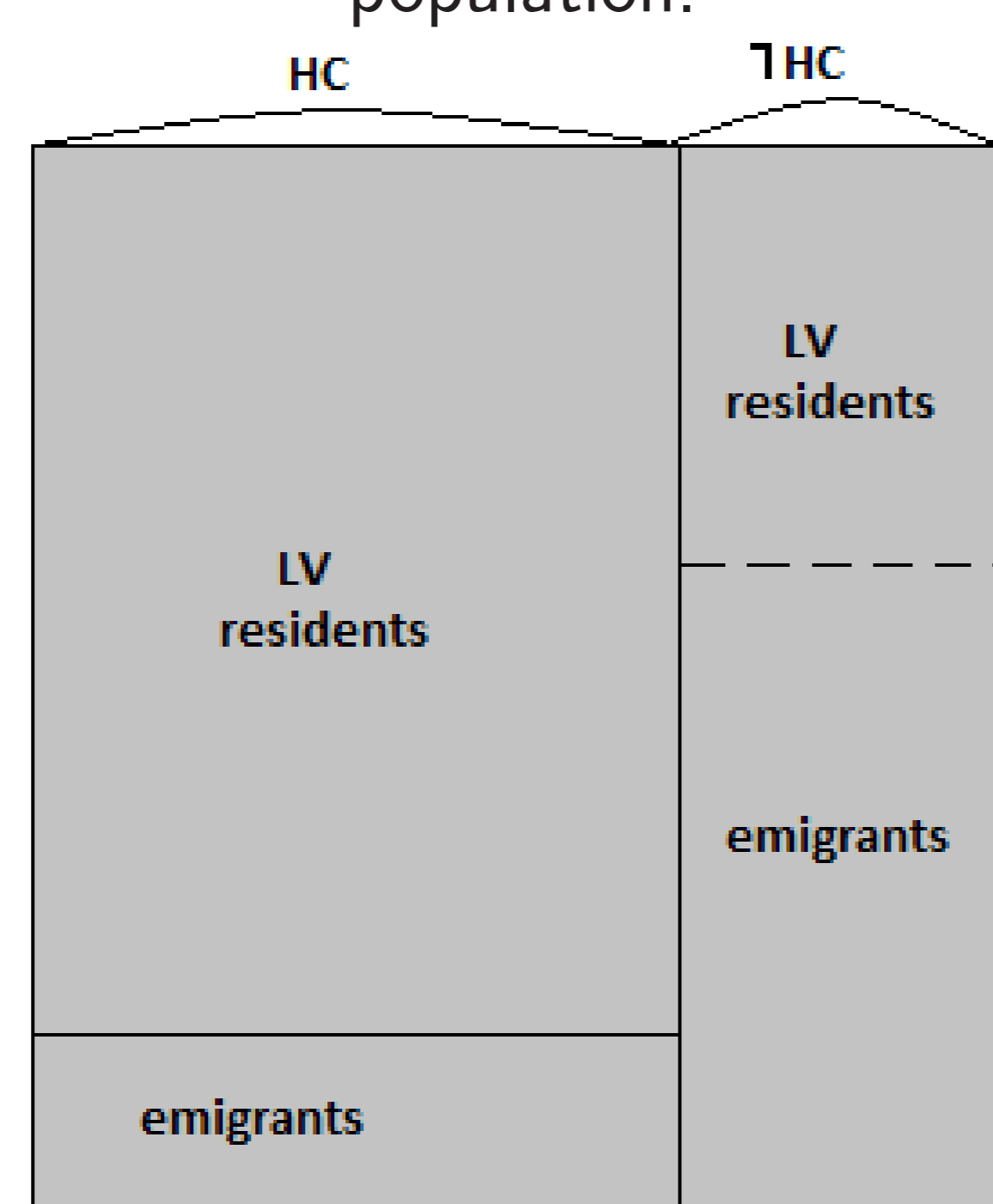
The Population and Housing Census was organised in March 2011. The main aim of it was to obtain detailed enough view on structure and characteristics of the population. There was obtained information about 1,94 mill. persons and counted up about 1,88 mill. residents in the Census. The information (is the person the resident or emigrant) about approximately 300 thousands persons is unknown. The number was estimated and according to the results of Census there were 2 070 371 residents of Latvia on 1 March 2011.

## Models

First of all, the logistic regression model was used for known population part.



Beta coefficients were estimated, using logistic regression for the whole population.



Tab. 1: Latvian residents

	01.01.2012.
Real value (CSP)	2 041 763
Logit model	2 140 073

The obtained value is much greater than the real value.

Tab. 2: Latvian residents

	01.01.2012.
Real value (CSP)	2 041 763
Logit model	2 055 350

The obtained result is close to the real value.

## Conclusion

The main problem may be that the number of emigrants in the investigated population part is too negligible (only 3% of the population), and in the unknown population part this relation is different in accordance to estimated value. Using the logistic regression for the known population we have obtained the conditional probability  $g_i = P(i \in LV | i \in HC)$ . To obtain the probability of unknown population part, we have to estimate the probability  $p_i = P(i \in LV | i \notin HC)$ .

## Model summaries

Tab. 5: Model summaries

Model	-2 Log likelihood	Cox & Snell R Square	Nagelkerke R Square
1.	1063065,512	,050	,121
2.	868832,663	,129	,314
3.	866943,933	,130	,316

## Auxiliary vectors

Tab. 4: Auxiliary information

	Auxiliary info	Transcription	beta		
			1st model	2nd model	3rd model
Sex	X 1	Sex	-,154	-,104	-,107
Marital status	X2 1	Single	-,662	-,488	-,484
	X2 2	Married		-,168	-,171
	X2 3	Divorced	-,659	-,631	-,631
Country of birth	X3	Latvia	,671	,732	,733
Nationality	X4 1	Latvians	1,039	,916	,912
	X4 2	Lithuanians	,768	,765	,767
	X4 3	Estonians	,437	,414	,410
	X4 4	Germans	-,597	-,707	-,705
	X4 5	Belarussians	1,015	,908	,908
	X4 6	Russians	,666	,658	,659
	X4 7	Ukrainians	,675	,638	,639
	X4 8	Poles	,661	,557	,556
	X4 9	Jews	-1,286	-1,447	-1,447
	X4 10	Roma	-,385	,112	,119
Citizenship	X4 19	Unknown	,626	,588	,587
	X5 1	LV citizens	,460	,279	,277
	X5 2	LV non-citizens	,356	,309	,307
Living region	X5 19	no information	,454	,373	,367
	X8 1	Riga	-,143	-,142	-,145
	X8 2	Latgale	,058	,050	,056
	X8 3	Zemgale	,021		
	X8 4	Vidzeme	,058		
Age groups	X8 5	Kurzeme	-,046	-,053	-,051
	X41 00	0- 4		-,154	-,156
	X41 01	4-9	-,348	-,512	-,514
	X41 02	10-14	-,229	-,403	-,406
	X41 2	15-19	-,610	-,839	-,845
	X41 3	20-24	-1,508	-2,245	-2,264
	X41 4	25-29	-1,786	-2,851	-2,873
	X41 5	30-34	-1,659	-2,722	-2,741
	X41 6	35-39	-1,302	-2,369	-2,385
	X41 7	40-44	-1,050	-2,062	-2,075
Employment	X41 8	45-49	-,663	-1,529	-1,538
	X41 9	50-54	,315	,091	,087
	X41 10	55-59	,843	,823	,822
Employer's report code	sum	income		,000	,000
	sum cik	employment		1,173	1,086
	men pnz ien	self-employment		3,841	3,842
Constant	ZK 1	code			1,453
	zaud st	21.-25.			-,702
	apdr atv	11.			-,717
		50., 51.			-1,313
			2,168	2,292	2,295

## Conclusions

- Using the logistic regression model for a known part of the population, we have obtained high probabilities that a person is Latvian resident. This result could be explained by the fact that the number of emigrants in the population is negligible.
- $\beta$  coefficients were calculated, using the logistic regression model for the whole population. The estimated value is close to the true value.