

---

The Truncated Poisson Models with covariates  
to estimate the size of drug users' population at  
risk of cautioning for personal possession

Flavia Mascioli<sup>1</sup>   Carla Rossi<sup>2</sup>   Daria Scacciatelli<sup>2</sup>

<sup>1</sup>Department of Mathematics  
University of Rome "La Sapienza"

<sup>2</sup>Centre for Biostatistics and Bioinformatics  
University of Rome Tor Vergata

Illicit Drug Market Workshop, 2009

# Outline

Introduction

Background

Results

Conclusion

## Introduction

### ▶ **Study Objective**

Estimate the number of drug users' at risk of cautioning for personal possession in Italy, in the year 2007.

Show the robustness of our results by a bootstrap simulation study.

### ▶ **Data Source**

Real data, provided by the Italian Ministry of Interiors in which individuals are divided by sex, age class, geographical area and number of cautionings

Simulated data, where the parameters were chosen according to the real data.

### ▶ **Methods**

Truncated Poisson model with covariates, Horvitz-Thompson, Zelterman and Chao estimators and the Bootstrap method, are considered.

## Formulation of the problem

- ▶ Let  $n_1, n_2, \dots, n_m$  be the frequencies of individuals identified 1, 2, ...,  $m$  times and let  $p_1, p_2, \dots, p_m$  be the associated probabilities.
- ▶ The frequency  $n_0$ , the unobserved number of individuals who are identified zero times, is missing and constitutes the target of the inference
- ▶ A population has  $N$  units  $n$  of which are identified by some mechanism. If the probability to identify an unit is  $(1 - p_0)$ :

$$N = Np_0 + (1 - p_0)N = \text{unobserved} + \text{observed cases} = Np_0 + n$$

## The zero-truncated Poisson Model

If  $p_0$  is unknown, as happen in most of the real cases, a conventional approach assumes that the frequencies arise from a Zero-Truncated Poisson Model:

$$P_j = P_j(\lambda) = \exp(-\lambda) \frac{\lambda^j}{j!} \quad j = 1, \dots,$$

where the parameter  $\lambda$  can be estimated by different methods (i.e. Maximum likelihood approach).

## Underlying Assumptions

- ▶ **Closed population** (i.e. no in-migration or out-migration in the time period studied)
- ▶ **Homogeneous population** (i.e. no subgroups with markedly different probabilities to be observed and re-observed )
- ▶ **Constant probability of being observed** (i.e. the probability of being re-observed should not be influenced by the experience of a previous visit )

### Solutions:

- ▶ Considering a short time study period
- ▶ Introducing observable covariates
- ▶ Considering a short time study period and stratifying by geographical areas.

## Population size Estimators

1. **Horvitz-Thompson** estimator (1952)

$$\hat{N}_{HT} = \frac{n}{(1 - \hat{p}_0)}$$

this model is not flexible enough in capturing population heterogeneity, and will generally underestimate the population size

2. lower bound estimator by **Chao** (1987, 1989)

$$\hat{N}_C = n + \frac{n_1^2}{2n_2}$$

based on a mixture of Poisson distributions

3. robust estimator of **Zelterman** (1988)

$$\hat{N}_Z = \frac{n}{(1 - e^{-2\frac{n_2}{n_1}})}$$

based only on the lower frequencies  $n_1$  and  $n_2$

## Confidence intervals for $\hat{N}$

- ▶ A natural approach to construct the confidence intervals for  $\hat{N}$  assumes that  $\hat{N}$  is asymptotically normal:

$$\hat{N} \pm 1.96 \sqrt{\text{var}(\hat{N})}$$

- ▶ Chao (1989) proposed to reparametrize the model using:  $\log(\hat{N} - n)$ , where  $n$  is the number of different individuals identified.
- ▶ For Zelterman and Chao estimators, we consider the new variances formulation proposed by Böhning (2008).



## Real Data

- ▶ Datasets provided by the Italian Ministry of Interiors: individuals are identified for the first time in 2007 and are divided by sex, age class, geographical area and number (one or more than once) of cautionings.
- ▶ In 2008, Mascioli and Rossi compare the three estimators with covariates **sex** and **age class**.

Both estimators produce almost the same estimates of  $N$ , for each age class and for the total population. As expected, when  $\lambda$  is small, Zelterman estimates are slightly greater than Chao estimates (Böhning and Brittain, 2007).

<i>Males</i>	<i>Estimates</i>			<i>Estimates</i>	
	Age Classes	$\hat{\lambda}$	$\hat{N}_C$	95%CI	$\hat{N}_Z$
< 15	0.076	5706	3599-9213	5701	3565-9290
15-17	0.066	31179	24543-39767	31188	24455-39939
18-19	0.069	59742	50573-70714	59755	50440-70937
20-24	0.067	111588	98334-126768	111629	98164-127086
25-29	0.043	96439	77994-119480	96505	77873-119834
30-34	0.045	53480	40800-70339	53508	40700-70590
35-39	0.052	31345	23201- 42556	31350	23116-42733
> 39	0.051	25026	17700-35611	25036	17631-35783
Total (No covariate)	0.059	398857	369676-430492	399012	369407-431145

<i>Females</i>	<i>Estimates</i>			<i>Estimates</i>	
	Age Classes	$\hat{\lambda}$	$\hat{N}_C$	95%CI	$\hat{N}_Z$
< 15	0.0500	801	182-4104	800	180-4165
15-17	0.0430	3268	1201-9359	3267	1190-9451
18-19	0.0370	7401	3294-17088	7399	3271-17208
20-24	0.0300	21719	12007 -39743	21716	11954-39915
25-29	0.0310	12484	5898-26941	12482	5865-27089
30-34	0.0290	7282	2632-20850	7281	2616-20987
35-39	0.0430	3362	1235-9628	3361	1224-9722
> 39	0.0190	5833	1209-29861	5832	1203-30024
Total (No covariate)	0.033	61725	43761-87409	61745	43656-87682

For the highlighted areas the capture probability is smaller than in the other areas and the corresponding confidence intervals are larger.

Areas	Estimates			Estimates	
	$\hat{\lambda}$	$\hat{N}_C$	95% CI	$\hat{N}_Z$	95% CI
<i>Basilicata</i>	0.083	3087	1722-5714	3090	1704-5789
<i>Calabria</i>	0.057	12736	8364-19606	12731	8313-19719
<i>Campania</i>	0.117	11696	9390-14660	11693	9327-14754
<i>EmiliaRomagna</i>	0.059	35063	27437-44980	35040	27319-45120
<i>FriuliVeneziaGiulia</i>	0.053	4499	2139-9778	4500	2121-9873
<i>Lazio</i>	<b>0.036</b>	<b>54901</b>	<b>39210-77189</b>	<b>54938</b>	<b>39123-77470</b>
<i>Liguria</i>	0.051	23760	16789-33845	23753	16712-33989
<i>Lombardia</i>	0.050	66949	54142-82982	66934	53986-83190
<i>Marche</i>	0.048	15113	9551-24175	15112	9500-24309
<i>Abruzzo + Molise</i>	0.061	10605	6883-16541	10603	6839-16650
<i>Piemonte + Vald' Aosta</i>	0.051	50010	39156-64076	50009	39035-64276
<i>Puglia</i>	0.063	26622	20274-35129	26629	20195-35291
<i>Sardegna</i>	<b>0.035</b>	<b>25293</b>	<b>15459-41755</b>	<b>25297</b>	<b>15397-41938</b>
<i>Sicilia</i>	0.079	40013	33420-48034	40037	33324-48235
<i>Trentino - AltoAdige</i>	0.051	6466	3359-12749	6467	3334-12855
<i>Toscana</i>	<b>0.044</b>	<b>45111</b>	<b>33576-60853</b>	<b>45101</b>	<b>33461-61040</b>
<i>Umbria</i>	<b>0.020</b>	<b>14202</b>	<b>5166-40008</b>	<b>14196</b>	<b>5140-40183</b>
<i>Veneto</i>	0.050	37247	28095-49590	37238	27989-49759
<b>Total (No covariate)</b>	<b>0.057</b>	<b>434490</b>	<b>403930-467510</b>	<b>434230</b>	<b>403260-467730</b>

Both Chao and Zelterman estimates are greater than Horvitz-Thompson estimates.

Areas Age Classes	Data			Estimates				
	$n$	$n_1$	$n_{>1}$	$n_2$	$\hat{\lambda}$	$\hat{N}_C$	$\hat{N}_Z$	$\hat{N}_{HT}$
<i>Basilicata</i>	243	233	10	10	0.083	3087	3090	3035
<i>Calabria</i>	712	692	20	20	0.057	12736	12731	12913
<i>Campania</i>	1294	1220	74	72	0.117	11696	11693	11753
<i>EmiliaRomagna</i>	2054	1994	60	60	0.059	35063	35040	35849
<i>FriuliVeneziaGiulia</i>	230	224	6	6	0.053	4499	4500	4486
<i>Lazio</i>	1885	1851	34	32	0.036	54901	54938	52885
<i>Liguria</i>	1190	1160	30	30	0.051	23760	23753	24002
<i>Lombardia</i>	3301	3219	82	81	0.050	66949	66934	67552
<i>Marche</i>	712	695	17	17	0.048	15113	15112	15149
<i>Abruzzo + Molise</i>	630	611	19	19	0.061	10605	10603	10657
<i>Piemonte + Vald' Aosta</i>	2471	2409	62	61	0.051	50010	50009	50071
<i>Puglia</i>	1608	1558	50	49	0.063	26622	26629	26398
<i>Sardegna</i>	863	848	15	15	0.035	25293	25297	25115
<i>Sicilia</i>	2985	2869	116	111	0.079	40013	40037	39414
<i>Trentino – AltoAdige</i>	318	310	8	8	0.051	6466	6467	6427
<i>Toscana</i>	1943	1901	42	42	0.044	45111	45101	45596
<i>Umbria</i>	297	294	3	3	0.020	14202	14196	14801
<i>Veneto</i>	1845	1799	46	46	0.050	37247	37238	37620
<b>Total (No covariate)</b>	<b>24581</b>	<b>23887</b>	<b>694</b>	<b>696</b>	<b>0.057</b>	<b>434490</b>	<b>434230</b>	<b>443590</b>

## Simulated Data

covariates: sex, age class, geographical area

- ▶ For every covariate we simulate, by a Binomial distribution, a population of 10000 units and we make a Bootstrap with 100 replications, to estimate the variance of  $\hat{N}$
- ▶ The observed population size  $n$  and capture-recapture probability were chosen according to the real data.

Comparison between the results obtained both on real data and on simulated data. For males population, the results are superimposable.

<i>Males</i>	<i>Estimates</i>		<i>Estimates</i>	
Age Classes	$N_C$	95%CI	$N_Z$	95%CI
< 15	5706	3599-9213	5701	3565-9290
15-17	31179	24543-39767	31188	24455-39939
18-19	59742	50573-70714	59755	50440-70937
20-24	111588	98334-126768	111629	98164-127086
25-29	96439	77994-119480	96505	77873-119834
30-34	53480	40800-70339	53508	40700-70590
35-39	31345	23201- 42556	31350	23116-42733
> 39	25026	17700-35611	25036	17631-35783
Total (No covariate)	398857	369676-430492	399012	369407-431145

<i>Males</i>	<i>Bootstrap Estimates</i>		<i>Bootstrap Estimates</i>	
Age Classes	$\hat{N}_C^B(STD)$	95%CI	$\hat{N}_Z^B(STD)$	95%CI
< 15	6305(1922)	3574-11405	6307(1922)	3576-11408
15-17	31924(3964)	25106-40752	31927(3960)	25115-40744
18-198	61178(4738)	52619-71242	61188(4736)	52631-71124
20-24	110790(6809)	98281-125016	110793(6816)	98270-125036
25-29	94258(10612)	75738-117558	94274(10616)	75748-117584
30-34	54068(8901)	39348-74637	54075(8915)	39336-74682
35-39	32227(4689)	24326-42873	32232(4691)	24327-42883
> 39	26158(6213)	16637-41561	26157(6219)	16629-41578
Total (No covariate)	394016(16177)	363638-427100	393975(16227)	363507-427165

For females population, the results are less significant because the recapture probability is very low.

<i>Females</i>	<i>Data</i>		
Age Class	<i>n</i>	<i>n<sub>1</sub></i>	<i>n<sub>&gt;1</sub></i>
< 15	40	39	1
15-17	140	137	3
18-19	272	267	5
20-24	659	649	10
25-29	387	381	6
30-34	209	206	3
35-39	142	139	3
> 39	108	107	1
Total (No covariate)	1957	1925	32

The results obtained with the bootstrap simulation study are comparable with the ones obtained on the real data;

<i>Areas</i>	<i>Estimates</i>		<i>Bootstrap Estimates</i>	
Chao	$\hat{N}_C$	95% CI	$\hat{N}_C^B$ (STD)	95% CI
<i>Basilicata</i>	3087	1722-5714	3197 (1151)	1656-6415
<i>Calabria</i>	12736	8364-19606	14064(3774)	8466-23703
<i>Campania</i>	11696	9390-14660	12052(1490)	9505-15389
<i>EmiliaRomagna</i>	35063	27437-44980	35975(5006)	27494-47283
<i>FriuliVeneziaGiulia</i>	4499	2139-9778	5240 (2211)	2422-11682
<i>Lazio</i>	54901	39210-77189	56565(14468)	34725-92930
<i>Liguria</i>	23760	16789-33845	25166 (5444)	16640-38399
<i>Lombardia</i>	66949	54142-82982	69463 (7092)	56958-84882
<i>Marche</i>	15113	9551-24175	17018 (4953)	9821-29901
<i>Abruzzo + Molise</i>	10605	6883-16541	11418 (2697)	7288-18109
<i>Piemonte + Vald' Aosta</i>	50010	39156-64076	50750 (7375)	38320-67490
<i>Puglia</i>	26622	20274-35129	27314 (4047)	20524-36541
<i>Sardegna</i>	25293	15459-41755	30442 (13678)	13344-70960
<i>Sicilia</i>	40013	33420-48034	40249 (3579)	33868-47948
<i>Trentino – AltoAdige</i>	6466	3359-12749	7175 (3587)	2933-18299
<i>Toscana</i>	45111	33576-60853	48038 (7930)	34925-66364
<i>Umbria</i>	14202	5166-40008	39253 (6225)	28902-53563
<i>Veneto</i>	37247	28095-49590	38707 (5441)	29490-50997
<i>Total (No covariate)</i>	434490	403930-467510	443230 (18125)	409190-480290



<i>Areas</i>	<i>Estimates</i>		<i>Bootstrap Estimates</i>	
Zelterman	$\tilde{N}_Z$	95%CI	$\tilde{N}_Z^B$ (STD)	95%CI
<i>Basilicata</i>	3090	1704-5789	3198(1152)	1657-6419
<i>Calabria</i>	12731	8313-19719	14066 (3778)	8463-23718
<i>Campania</i>	11693	9327-14754	12055 (1494)	9503-15401
<i>EmiliaRomagna</i>	35040	27319-45120	35973 (5007)	27490-47283
<i>FriuliVeneziaGiulia</i>	4500	2121-9873	5243 (2212)	2423-11687
<i>Lazio</i>	54938	39123-77470	56582 (14483)	34722-92992
<i>Liguria</i>	23753	16712-33989	25167 (5449)	16634-38414
<i>Lombardia</i>	66934	53986-83190	69466 (7104)	56942-84913
<i>Marche</i>	15112	9500-24309	17024 (4960)	9819-29930
<i>Abruzzo + Molise</i>	10603	6839-16650	11419 (2698)	7288-18115
<i>Piemonte + Vald' Aosta</i>	50009	39035-64276	50748 (7382)	38308-67508
<i>Puglia</i>	26629	20195-35291	27315 (4048)	20524-36545
<i>Sardegna</i>	25297	15397-41938	30450 (13689)	13343-71010
<i>Sicilia</i>	40037	33324-48235	40270 (3576)	33895-47961
<i>Trentino – AltoAdige</i>	6467	3334-12855	7177 (3588)	2934-18305
<i>Toscana</i>	45101	33461-61040	48051 (7941)	34922-66406
<i>Umbria</i>	14196	5140-40183	39255 (6234)	28892-53589
<i>Veneto</i>	37238	27989-49759	38708 (5452)	29476-51026
<b>Total (No covariate)</b>	<b>434230</b>	<b>403260-467730</b>	<b>443190(18213)</b>	<b>408990-480440</b>

Comparison between the three estimators considering the model with or without covariates. Considering covariates allow to obtain better results

Real Data	$\hat{N}_C$	$\hat{N}_{HT}$	$\hat{N}_Z$
Males (all age covariate)	414504	409692	414672
Males (no covariate)	398857	394898	399012
Females (all age covariate)	63110	61391	63136
Females (no covariate)	66298	66394	60487

Simulate Data	$\hat{N}_C$	$\hat{N}_{HT}$	$\hat{N}_Z$
Males (all age covariate)	416907	415048	416952
Males (no covariate)	394016	395404	393975
Females (all age covariate)	75102	72272	75127
Females (no covariate)	64215	62101	64248

Real Data	$\hat{N}_C$	$\hat{N}_{HT}$	$\hat{N}_Z$
Areas (all area covariate)	483374	483722	483367
Areas (no covariate)	434490	443590	434230

Simulate Data	$\hat{N}_C$	$\hat{N}_{HT}$	$\hat{N}_Z$
Areas (all area covariate)	532086	527670	532167
Areas (no covariate)	443230	444680	443190

Comparison between the three estimators considering the model with or without covariates. Considering covariates allow to obtain better results

Real Data	$\hat{N}_C$	$\hat{N}_{HT}$	$\hat{N}_Z$
Males (all age covariate)	414504	409692	414672
Males (no covariate)	398857	394898	399012
Females (all age covariate)	63110	61391	63136
Females (no covariate)	66298	66394	60487

Simulate Data	$\hat{N}_C$	$\hat{N}_{HT}$	$\hat{N}_Z$
Males (all age covariate)	416907	415048	416952
Males (no covariate)	394016	395404	393975
Females (all age covariate)	75102	72272	75127
Females (no covariate)	64215	62101	64248

Real Data	$\hat{N}_C$	$\hat{N}_{HT}$	$\hat{N}_Z$
Areas (all area covariate)	483374	483722	483367
Areas (no covariate)	434490	443590	434230

Simulate Data	$\hat{N}_C$	$\hat{N}_{HT}$	$\hat{N}_Z$
Areas (all area covariate)	532086	527670	532167
Areas (no covariate)	443230	444680	443190

## Conclusion and Future Work

### Conclusion

- ▶ Considering covariates in the model produces more accurate estimates
- ▶ The bootstrap simulation study confirm the coherence of the results obtained with real data (with the exception of females population)

### Future Work

- ▶ Analysis of differences between large and small areas ( several results have been already obtained)
- ▶ Inclusion of more covariates in the model, socioeconomic background and behavior
- ▶ Use more sophisticated simulation schemes

## References

- ▶ Böhning, D. (2008). A Simple Variance Formula for Population Size Estimators by Conditioning. *Statistical Methodology*, 5, 410-423
- ▶ Böhning, D. and Del Rio Vilas, V. (2008). Estimating the hidden number of scrapie affected holdings in Great Britain using a simple, truncated count model allowing for heterogeneity. *Journal of Agricultural, Biological, and Environmental Statistics*, 13, 1-22.
- ▶ Chao, A. (1987). Estimating the population size for capture-recapture data with unequal catchability. *Biometrics*, 43, 783-791.
- ▶ Chao, A. (1989). Estimating population size for sparse data in capture-recapture experiments. *Biometrics*, 45, 427-438.
- ▶ Horvitz, D. G. and Thompson, D. J. (1952). A generalization of sampling without replacement from a finite universe. *Journal of the American Statistical Association*, 47, 663-685.
- ▶ Mascioli F. and Rossi C.(2009). Capture-recapture methods to estimate prevalence indicators used in the evaluation of drug policies. 3rd Annual Conference of the International Society for the Study of Drug Policy, (Vienna, 2-3 March 2009).
- ▶ Rossi C. and Ricci R. (2009). Modelling and estimating illicit drug market as a tool to evaluate drug policy: the case of Italy. 3rd Annual Conference of the International Society for the Study of Drug Policy, (Vienna, 2-3 March 2009).
- ▶ Scalia Tomba GP., Rossi C., Taylor C., Klempova D., Wiessing L. (2008). Guidelines for Estimating the Incidence of Problem Drug Use. EMCDDA, Lisbon.
- ▶ Zelterman, D. (1988). Robust estimation in truncated discrete distributions with application to capture-recapture experiments. *Journal of Statistical Planning and Inference*, 18, 225-237.