



XXVI Congresso

Sociedade Portuguesa de Estatística

TRS method for Census 2021 data at Statistics Portugal

Inês Rodrigues de Sá

Pedro Campos

Pedro.campos@ine.pt



INSTITUTO NACIONAL DE ESTATÍSTICA
STATISTICS PORTUGAL

1



Contents

- Introduction
- Candidate Methods: CKM and TRS
- Risk Measures and Utility Measures
- Challenges, Vantages and Limitations
- Applications to Census 2021
- Results
- Discussion

2

INSTITUTO NACIONAL DE ESTATÍSTICA
STATISTICS PORTUGAL

2



Introduction

- Between September 2016 and August 2017, the Centre of Excellence (CoE) on Statistical Disclosure Control (SDC) maintained a Specific Grant Agreement (SGA) dedicated to the harmonisation of methods for protecting the **confidentiality of aggregated Census data in the European Statistical System**
- Two methods were proposed to protect the confidentiality of these data: **targeted record swapping (TRS)** and **cell key method (CKM)**.

3

3



Candidate Methods

Cell Key Method

- Post-tabular (applied to the table cells)
- Consistently adds unbiased random noise to each table cell

“Countries that do not use a combination of pre and post tabular SDC methods are advised to use the cell key method”

Targeted Record Swapping

- Pre-tabular (applied to the microdata)
- exchange of geography-related variables between linked households.

Recommended by EU-project
“Harmonized Protection of Census Data in the ESS”

4

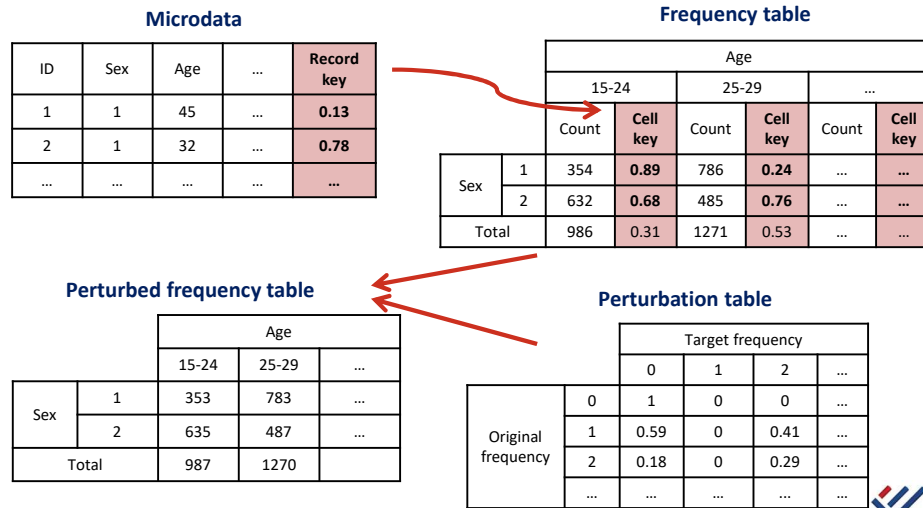
4



Candidate Methods

Cell key method

(Marley & Leaver, 2011; Enderle et al., 2018)



5



Candidate Methods

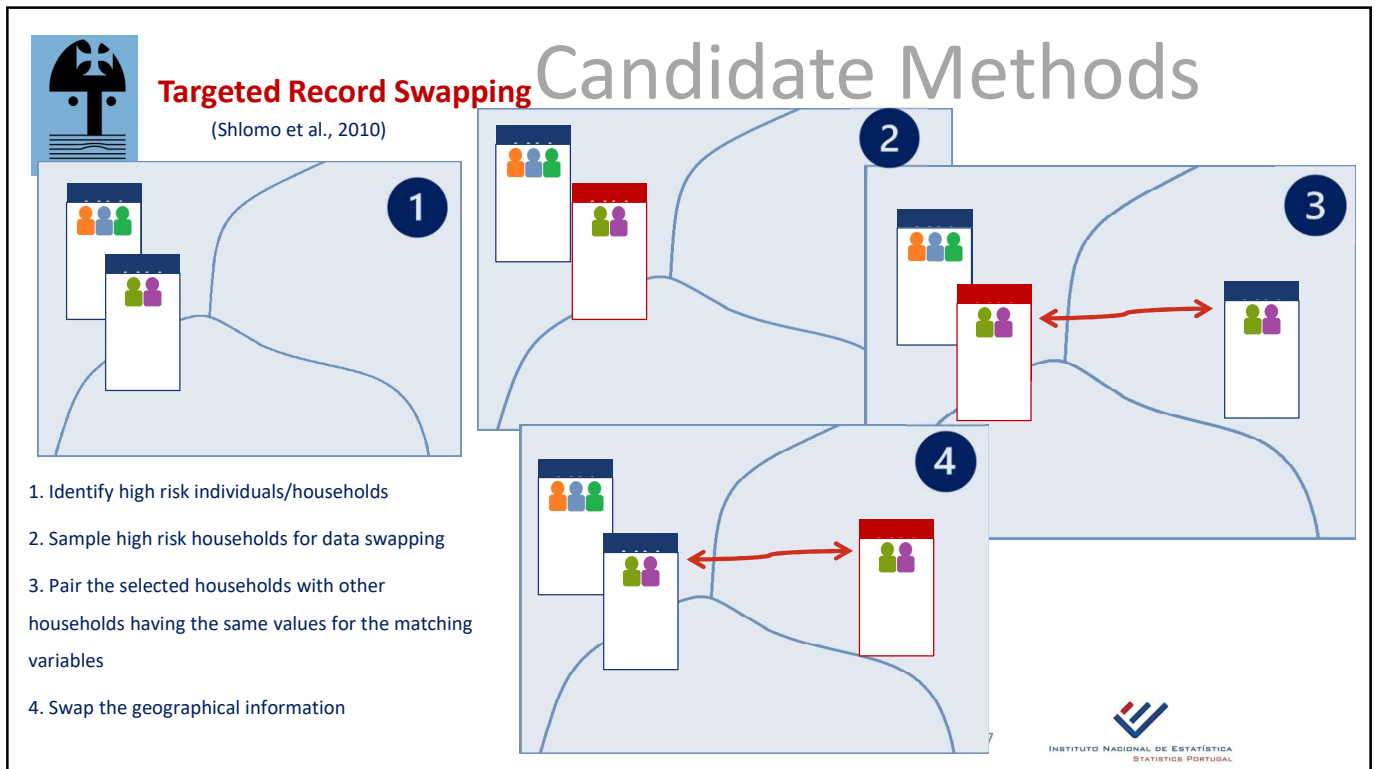
Cell key method

(Marley & Leaver, 2011; Enderle et al., 2018)

- Tested on PT Census 2011 data
- Two groups of EU-hypercubes (Commission Regulation (EU) 2017/712, of 20 April 2017; Commission Implementing Regulation (EU) 2017/543, of 22 March 2017) and some national tables (**no grid data**)
- Risk and utility measures to compare results and support method/parameter choice, **BUT no assessment of disclosure by differencing**

6

6



7

Risk Measures

▪ Let:

- n_c : number of units that fall into cell c in the original table T
- n'_c : number of units that fall into cell c in the protected table T'
- K : total number of cells in table T (or T')

RM 1

Relative change of the number of cells with frequency lower than 3 (change in low frequencies)

$$CLF = \left(\frac{\sum_{c=1}^K I(n'_c < 3)}{\sum_{c=1}^K I(n_c < 3)} - 1 \right) \times 100\%$$

RM 2

Proportion of cells with frequency lower than 3 both in the original and the perturbed table (real low frequencies)

$$RLF = \frac{\sum_{c=1}^K I(n_c < 3 \wedge n'_c < 3)}{K} \times 100\%$$

INSTITUTO NACIONAL DE ESTATÍSTICA
STATISTICS PORTUGAL

8



Utility Measures

UM 1 (Shlomo & Young, 2005; Shlomo, 2007)

Absolute distance (AD) between the original and the perturbed counts

$$AD_c = |n'_c - n_c|$$

UM 2 (Shlomo & Young, 2005; Shlomo, 2007)

Relative absolute distance (RD) between the original and the perturbed counts

$$RD_c = \frac{|n'_c - n_c|}{n_c}$$



Simple descriptive statistics (max, mean, sd, median) across all cells

9

9



Utility Measures

UM 3 (Burton et al., 2017)

Proportion of false zeros

$$FZ = \frac{\sum_{c=1}^K I(n'_c = 0 \wedge n_c \neq 0)}{\sum_{c=1}^K I(n'_c = 0)} \times 100\%$$

UM 4

Proportion of unchanged cells

$$UC = \frac{\sum_{c=1}^K I(n'_c = n_c)}{K} \times 100\%$$

10

10



Challenges

- Disclosure by differencing (e.g. grid cells versus administrative regions) is difficult to measure
- CKM results in loss of table additivity
- CKM can result in false zero frequency cells (but data items on total population shall nevertheless be flagged as 'populated', according to Regulation 1799 of 21 November 2018, Article 6)
- Users need to be aware that perturbative SDC methods were used
- Selected disclosure risk and utility indicators might be published together with data
- The loss of table additivity due to confidentiality protection should be clearly stated

Perturbed frequency table

		Age		
		15-24	25-29	...
Sex	1	353	783	...
	2	635	487	...
Total		987	1270	

Communicating to the users

11

11



CKM vs TRS

The analysis of the limitations associated with CKM - in particular, the loss of additivity of the tables caused by this method and its high operational demands - **leads to the conclusion that TRS is the best option to consider for protecting the confidentiality of Census aggregate data**

12

12



TRS

The Targeted Record Swapping (TRS) - was specified by the UK Office for National Statistics (ONS). It is applied to micro-data (pre-tabular method) and consists of four main stages:

1. identification of records with high risk of disclosure of confidential information;
2. sample selection of households to be exchanged;
3. selection of linking households (pairs);
4. exchange of geography-related variables between linked households.

13

13



TRS

- TRS can be considered a special case of PRAM;
- However, has never been applied in the context of a large-scale census, and TRS has generally been preferred for its ease of implementation

14

14



Advantages of TRS

- Being a pre-tabular method, TRS only needs to be **applied once to the microdata base**;
- then **all tabulations must be obtained** from the perturbed microdata;
- All tables obtained from the microdata base after applying the TRS **are consistent** (the same intersection always results in the same perturbed result) and additive (any (sub-)total of the table, by row or column, is equal to the sum of the cells corresponding to the respective plots);
 - **The application of TRS decreases the risk of disclosure of confidential information by increasing the uncertainty related to the data, namely at lower geographical levels.**

15

15



Limitations of TRS

- All frequencies relating to geographical levels higher than the geographical hierarchy considered in the application of the method (e.g. NUTS II, NUTS I and national total levels) **are not disturbed (and, therefore, are not protected)**;
- **The isolated application of TRS is not recommended**, as it only provides protection for lower geographic levels, as it causes greater loss of information in terms of distances between original and protected cells relative to CKM, and as it is less suitable for protection against microdata base reconstruction attacks
- **Changing the place of residence of a given individual after applying the TRS may lead to inconsistencies between the places of residence/work and the duration of commuting.**

16

16



Is it possible to get around the limitations of TRS?

- The TRS does not modify the frequencies relating to geographical levels higher than those of the geographical hierarchy used, i.e. we can consider that these frequencies remain unprotected;
- **However, the risk of confidential information disclosure associated with such frequencies is, as a rule, very low (due to the lower occurrence of both low frequencies and frequencies concentrated in a single cell per row or column).**

17

17



Application to Portuguese Census

- Data from the **2011 Census were used for Portugal**, covering the entire national territory (Mainland and Autonomous Regions).
- For an initial analysis, the municipality of Olhão (middle size one) was selected.

18

18



Application to Portuguese Census

- Tables :
 - Table 4.12 - Classical households, according to the socio-economic group of the household representative, by type of household;
 - Table 6.05 - Resident population by age group, by nationality and sex;
 - Table 6.07 - Portuguese resident population born abroad by age group, by country of birth and sex;
 - Table 6.21A - Resident population 5 years old and over according to age group and sex
 - Table 6.35 - Employed resident population, by employment status and sex, by branch of economic activity and hours worked in the reference week;
 - Table 6.49 - Resident population aged 15 and over, by response to the question on religion.

19

19



Application to Portuguese Census

The application of the TRS was based on the following parameters:

- Geographical hierarchy (hierarchy):
 - Scenario 1: NUTS3 > Municipality > Parish;
 - Scenario 2: Municipality > Parish;
- Swaprate: 5%;
- Risk_variables: gender, five-year age group, country of birth different from Portugal;
- Threshold for defining high risk (k_anonymity): 3;
- Variables for finding housing-pair (similar):
 - Number of individuals in the dwelling aged < 15 years;
 - Number of males in the dwelling aged 15-64 years;
 - Number of males in the accommodation aged ≥ 65 years;
 - Number of women in the dwelling aged 15-64 years;
 - Number of women in dwelling aged ≥ 65 years;
 - Number of persons employed or students in the dwelling.

20

20



Software

- The two methods were applied using the respective implementations under development by SGA Open source tools for perturbative confidentiality methods .
- In particular, the packages **recordSwapping** (version 0.1.0), **ptable** (version 0.2.0) and **cellKey** (version 0.16.3) of R

21

21



Results for 2011

- Considering the totality of the records, scenario 1 (NUTS3 geographic hierarchy > Municipality > Parish) **results in fewer dwellings/individuals affected by the exchanges** (although the differences are small, especially in relative terms).

22

22



Results

Scenario	Hierarchy	Level	Dwellings swapped		Individuals swapped	
			N.	%	N	%
1	NUTS3 > Municipality > Parish	NUTS3	156500	68,5	463722	69,4
		Municipality	62060	27,2	177574	26,6
		Parish	9870	4,3	27148	4,1
		Total	228430	3.89*	668444	6,3**
2	Municipality > Parish	Municipality	197146	85,3	575914	85,4
		Parish	34096	14,7	98346	14,6
		Total	231242	3.93*	674260	6,5**

23

23



Results for 2021

We started by applying the TRS to all records, using the parameters defined for most cases:

- threshold for minimum frequency = 3
- swaprate = 5%
- variables to determine risk = gender, decennial age groups (last group 80+) and country of birth ≠ PT.

24

24



Results for 2021

All parishes whose % of households affected by swaps was $> 10\%$ in this first TRS application were subject to an adjustment of the TRS parameters:

- threshold for minimum frequency = 2
- swaprate = 3%
- variables to determine risk = sex, major age groups (<15; 15-64; ≥ 65) and country of birth \neq PT.

25

25



Results for 2021

- Additionally, for municipalities with only one parish in this subset, another parish from the same municipality was randomly selected for this adjustment, to allow for exchanges at the parish level.
- We are left with a total of 979 parishes in this subset - the TRS with parameter adjustment was applied to all the records of these 979 parishes; for the records of the remaining parishes, the parameters referred to in the previous paragraph were maintained.

26

26



Results for 2021

However, even after the adjustment of the parameters in this set of parishes, we continue to verify the existence of parishes with % exchanges $> 10\%$



not only in the set of parishes that initially had rates $> 10\%$ (in which 12 parishes with rates $> 10\%$ remained), but also in the set of parishes that after the first application of the TRS had rates $< 10\%$ (in which 269 parishes with rates $> 10\%$ appeared).

27

INSTITUTO NACIONAL DE ESTADÍSTICA
STATISTICS PORTUGAL

27



Results for 2021

- This can be justified by the fact that, in the initial application of the TRS, considering all the records, there may have been exchanges involving dwellings in the parishes that were subject to adjustment and which "made it impossible" to use the respective dwellings for other exchanges, namely with dwellings from these parishes in which the exchange rate increased (because a dwelling can only be selected for exchange once).



By applying the TRS for the two sets of records (with different parameters), some exchanges may have become possible in the set of parishes that were not subject to parameter adjustment.

28

INSTITUTO NACIONAL DE ESTADÍSTICA
STATISTICS PORTUGAL

28



Results for 2021

In order to ensure that no parish was left with a % of dwellings affected by the exchanges > 10%, it was then necessary to reverse some of the exchanges.

- All exchanges involving dwellings in two of the parishes concerned were reversed.
- Besides these, as many exchanges were reverted as necessary to reach the 10% threshold per parish; these exchanges were randomly selected among all the existing exchanges, per parish.
- Thus, all the parishes have this % below 10%

29

29



Results for 2021 -Risk

Risk measures after RRT, scenario 1 (NUTS3>Mun>Freg), municipality of Olhão

Quadro		Freq < 3 Dados orig	Freq < 3 Dados pert	Varição Freq < 3 (%)		Freq < 3 reais	Div atrib Dados orig	Div atrib Dados pert
4.12	Nº	1515	1502	-0,86	Nº	1398	1	0
	%	18,4	18,2		% (Freq < 3)	93,1	0,0	0,0
6.05	Nº	2887	2889	0,07	Nº	2194	0	0
	%	12,2	12,2		% (Freq < 3)	75,9	0,0	0,0
6.07	Nº	2189	2231	1,92	Nº	1743	0	0
	%	10,5	10,7		% (Freq < 3)	78,1	0,0	0,0
6.21A	Nº	926	927	0,11	Nº	843	28	30
	%	11,3	11,3		% (Freq < 3)	90,9	0,3	0,4
6.35	Nº	8577	8571	-0,07	Nº	8010	9	10
	%	9,6	9,6		% (Freq < 3)	93,5	0,0	0,0
6.49	Nº	4	3	-25,0	Nº	3	0	0
	%	6,7	5,0		% (Freq < 3)	45,0	0,0	0,0

30



Results for 2021 - Utility

False zeros and false positives after TRS, scenario 1 (NUTS3>Mun>Freg), municipality of Olhão

Quadro	Cenário	Hierarquia	Falsos zeros		Falsos positivos	
			Nº	%	Nº	%
Q412	1	NUTS3>Mun>Freg	65	1,5	55	1,4
Q605	1	NUTS3>Mun>Freg	437	2,7	530	7,4
Q607	1	NUTS3>Mun>Freg	325	2,0	335	7,7
Q621A	1	NUTS3>Mun>Freg	36	4,7	44	0,6
Q635	1	NUTS3>Mun>Freg	344	0,5	402	2,2
Q649	1	NUTS3>Mun>Freg	-	-	0	0,0

31

31



Discussion

The results showed the effect of the TRS is neither evident nor substantial with respect to cells with low frequencies: depending on the frame and the municipality under analysis, the TRS can lead to a decrease, maintenance or even an increase in the number of cells with low frequencies.

32

32



Discussion

It turns out, however, that as a rule, the number of cells with actual low frequencies (i.e. whose frequency is also low in the original data) is always lower than the number of cells with low frequencies in the perturbed data.

33

33



Discussion

This fact reflects the main effect of the TRS regarding data protection: the increased uncertainty associated with the analysis of the information provides, in itself, some degree of protection.

34

34



Discussion

- With regard to the risk of attribute disclosure, we find that in the generality of the tables this risk is already very low from the outset.
- TRS had no significant effects on this type of disclosure in the analysed tables and municipalities.