

# Weighting data in national samples

Marcin W. Zieliński

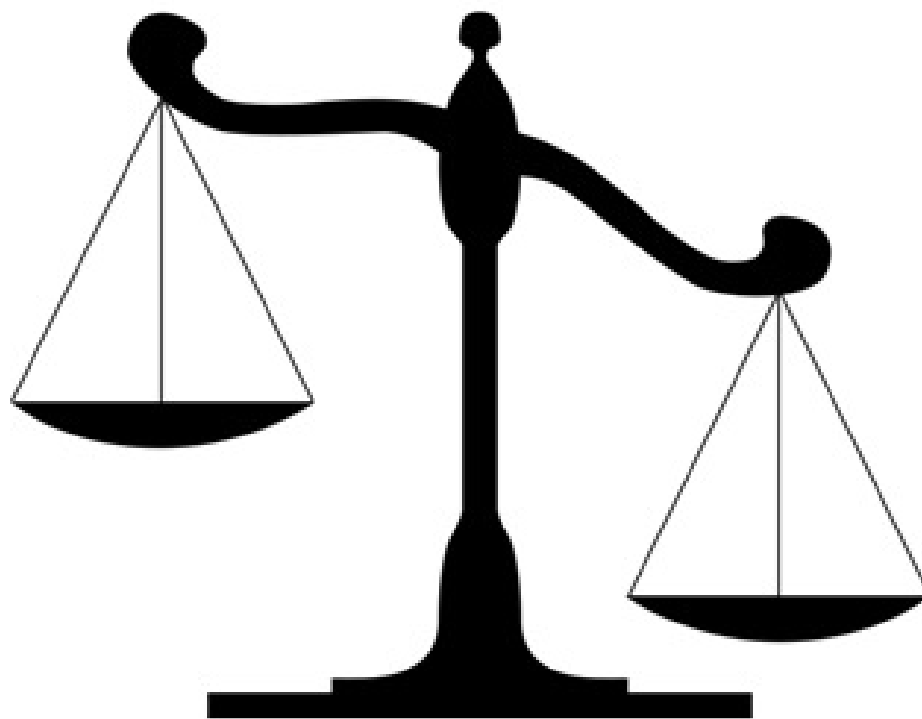
Polish Academy of Science

University of Warsaw

# Traditional vs survey weighting



# Traditional vs survey weighting



# survey weighting

**sample**



**population**



# Survey weight

- Value assigned to each case to make the sample more representative of the population
- Main reasons for non-representativeness of the sample:
  - sample design -> design weights
  - non-responses -> poststratification weights

# Design weights

- Needs information how the sample was constructed

Corrects for:

- Oversampling
- Household size
- Probabilities of selection on different stages in multistage samples

# Design weights

1. Calculate inverse probability of inclusion to the sample

$$WF_{des} = 1/\pi$$

where:

$WF_{des}$  – weighting factor

$\pi$  – probability of being included to the sample

2. Scaling:

$$sWF_{des} = WF_{des}/mean(WF_{des})$$

# An example

- Household sample

Kish grid

HH size	1	2	3	4	5	6	7	8	9
Person interviewed	1	2	1	3	5	4	7	3	2



# An example

- Household sample

Kish grid

HH size	1	2	3	4	5	6	7	8	9
Person interviewed	1	2	1	3	5	4	7	3	2
$\pi$	1	1/2	1/3	1/4	1/5	1/6	1/7	1/8	1/9

# An example

Resp. no	HH size	$\pi$	$WF_{des}$	$sWF_{des}$
1	4	1 / 4	4	4/3,2=1,2500
2	2	1 / 2	2	2/3,2=0,6250
3	2	1 / 2	2	2/3,2=0,6250
4	6	1 / 6	6	6/3,2=1,8750
5	8	1 / 8	8	8/3,2=2,5000
6	4	1 / 4	4	4/3,2=1,2500
7	1	1	1	1/3,2=0,3125
8	1	1	1	1/3,2=0,3125
9	1	1	1	1/3,2=0,3125
10	3	1 / 3	3	3/3,2=0,9375
			<b>mean=3,2</b>	<b>mean=1</b>

# Poststratification weights

- Needs external sources of information about the population, mostly about its demographic structure, like
  - - age
  - - gender
  - - place of residence (urban/rural, size of the city)
  - - region of residence
  - - education

# Poststratification weights

$$W_{post} = \frac{P_{prop}}{S_{prop}}$$

where:

$P_{prop}$  – proportion of the group in population

$S_{prop}$  – proportion of the group in the sample

# Poststratification weights

	Sample (N)	Sample (prop)	Popul. (prop)	WEIGHT
MALE	638	0,638	0,45	$0,45 / 0,638 = 0,705329$
FEMALE	362	0,362	0,55	$0,55 / 0,362 = 1,519337$
SUM	1000	1	1	



<b>SAMPLE</b>	18-20	21-29	30-39	40-49	50-59	60-69	70+	
MALE	0,0053	0,0233	0,0586	0,1651	0,1172	0,0533	0,0313	<b>0,4541</b>
FEMALE	0,0133	0,0533	0,1065	0,0999	0,1332	0,0932	0,0466	<b>0,546</b>
	<b>0,0186</b>	<b>0,0766</b>	<b>0,1651</b>	<b>0,265</b>	<b>0,2504</b>	<b>0,1465</b>	<b>0,0779</b>	<b>1</b>

<b>POPUL</b>	18-20	21-29	30-39	40-49	50-59	60-69	70+	
MALE	0,0126	0,0506	0,0826	0,1252	0,1278	0,0739	0,0379	<b>0,5106</b>
FEMALE	0,0166	0,0679	0,1218	0,0866	0,0826	0,0752	0,0386	<b>0,4893</b>
	<b>0,0292</b>	<b>0,1185</b>	<b>0,2044</b>	<b>0,2118</b>	<b>0,2104</b>	<b>0,1491</b>	<b>0,0765</b>	<b>1</b>

<b>WEIGHT</b>	18-20	21-29	30-39	40-49	50-59	60-69	70+
MALE	0,0126/ 0,0053= 2,3774	2,1717	1,4096	0,7583	1,0904	1,3865	1,2109
FEMALE	1,2481	1,2739	1,1437	0,8669	0,6201	0,8069	0,8283

# Problems:

## 1. Empty cells

<b>SAMPLE</b>	18-20	21-29	30-39	40-49	50-59	60-69	70+	
MALE	0,0053	0,0233	0,0586	0,1651	0,2504	0,0533	0,0313	<b>0,587</b>
FEMALE	0,0133	0,0533	0,1065	0,0999		0,0932	0,0466	<b>0,413</b>
	<b>0,0186</b>	<b>0,0766</b>	<b>0,1651</b>	<b>0,265</b>	<b>0,2504</b>	<b>0,1465</b>	<b>0,0779</b>	<b>1</b>

- Joining cells
- Assigning a very small proportion

## 2. No external cross-classified data

- IPF/raking procedure



# Iterative Proportional Fitting (IPF)/raking



# Iterative Proportional Fitting (IPF)/raking



1. Calculate multiplier of the  $V_1$  that equates sampling marginal distribution to the population marginal distribution using the formula:

$$WFV_1 = \frac{P_{prop}V_1}{S_{prop}V_1}$$

2. Weight  $V_2$  by  $WFV_1$

3. Calculate multiplier of the  $V_2$  that equates weighted (by  $WFV_1$ ) sampling distribution to the population marginal distribution using the formula

$$WFV_2 = \frac{P_{prop}V_2}{S_{prop}V_2}$$

4. Iterate until you're satisfied (or until it makes sense)

# In short:

1. Calculate multiplier of V1 to make V1 distribution equal to V1 population distribution
  2. Weight by multiplier from step 1
  3. Calculate weighted multiplier of V2
  4. Weight by multiplier from step 3
  5. Calculate weighted multiplier of V1
  6. Weight by multiplier from step 5
  7. Calculate weighted multiplier of V2
- [....]

$$FinalWeight = \frac{Weighted\ in\ last\ step_{prop}}{Sample_{prop}}$$

# What changes weighting the data?

- Equalizes proportions in the sample to the population in weight components
- Weight is a multiplier, so it affects all the data, statistics and analysis
- If the weight component is related to the outcome it has bigger effect on the outcome; if not, the effect of weighting is marginal
- Some claim that if you want to extend your findings to the population you should use weights

# Is weighting a golden mean?

- The general aim of weighting is to weight the whole sample. It corrects for unit-nonresponses. What about item-nonresponses?
- Weights almost always increase standard errors of estimates (effect on precision)
- Problem of very large/small weights (may lead to bias results)
- Weighting works only on MAR assumption but not MCAR or NMAR. How to check if data are NMAR? If data are NMAR do we still can expand our findings to the population?
- What if we selected e.g. a group that combines characteristics used in the process of weight calculation? What is the effect of weighting?

# When weight really matters?

**Dependent: life satisfaction**

People from bigger households tend to be more satisfied and there was HH sample used

What kind of weight corrects for this issue?

# When weight really matters?

**Dependent: life satisfaction**

Men are more satisfied than women and women tend to take part in a survey more frequently

What kind of weight should we use? Does it corrects anything for this problem?



# When weight really matters?

## **Dependent: life satisfaction**

People from big cities tend to be less satisfied than people from villages and those from villages tend to take part in a survey more frequently

What kind of weight should we use? Does it corrects anything for this problem?

# When weight really matters?

**Dependent: life satisfaction**

People who are less satisfied tend not to take part in a survey and there was HH sample used

What kind of weight should we use? Does it corrects anything for this problem?

# When weight really matters?

**Dependent: life satisfaction**

Younger people are just as satisfied as older people but younger tend to take part in a survey less frequently

What kind of weight should we use? Does it corrects anything for this problem?

# How to interpret weights?

- Weight = 0
- Weight = 1
- Weight > 1
- Weight > 0 and < 1