# EGM 08: The evolution of

# the traditional model

**Carlos Lozano (AIMC)**
**Julián Sánchez (ODEC/QUINAO)**

**EMRO Conference**
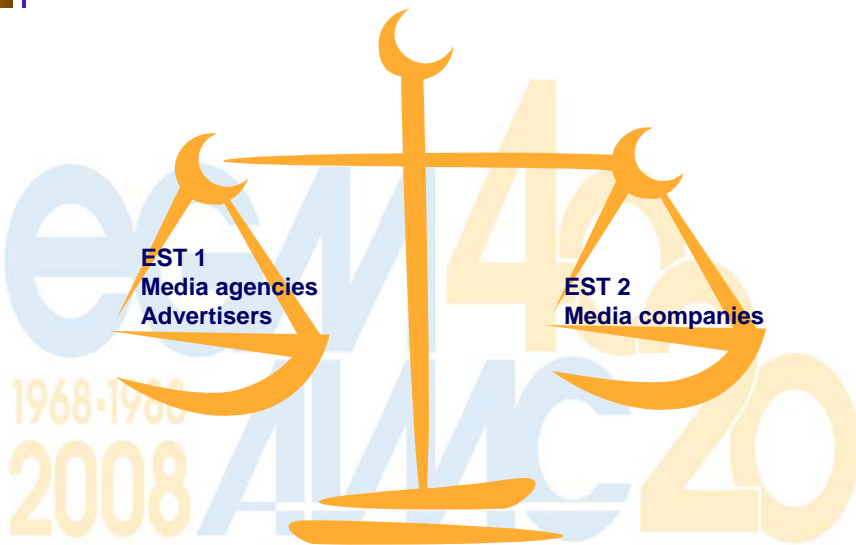**Biarritz 24-28 May 2008**

# THE "NEW EGM"

- **The <u>process</u>: 2006-07**
  **Political and technical agreement.**

- **The <u>conclusion</u>: November 07**
  **Approved in an AIMC´s Assembly (90%)**

- **The <u>results</u>: First wave EGM- 29th April 08**

# THE EQUILIBRATE UNSTABLE

**EST 1**
**Media agencies**
**Advertisers**

**EST 2**
**Media companies**

# THE EQUILIBRATE UNSTABLE

**EST 1**
**Cross Media**
**Consumer centric**

**EST 2**
**Media oriented**
**Single media surveys**

# MODEL EVOLUTION

**2006**
**EGM Prensa**

**2000**
**EGM Radio XXI**

**1968**
**EGM Multimedia**

## 2008: THE "NEW EGM"

Multimedia
Radio
Newspapers
Magazines
20.000 CAPI

## DATA FUSION: The requirements

- **Maintenance the existing "currency"**

- **Maximum automation**

- **Transparency. No "black boxes"**

**ODEC**

**EMRO Conference 2008**

**EGM 08:**
**THE EVOLUTION OF THE TRADITIONAL MODEL**

**OPTIMAL DATA FUSION**

**Biarritz, 26th May 2008**

Julián Sánchez Montenegro
Director Quinao Consultores
Paseo de la Habana 26, 28036 MADRID
jsanchez@quinao.com / +34 609015365

**QUINAO**
consultores

1

---

**ODEC**

**STARTING POINT**

**Coexistence of:**
• **Multimedia Survey (EGM)**
• **Monomedia Surveys**

**Need to integrate different surveys for:**

• **Obtaining a single currency, an unified figure for audiences.**
• **Make the most of the information available.**

**Solution :**

**SURVEY DATA FUSION**

Starting point: Coexistence of different data sources

Interviews:

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| MULTIMEDIA | 10.000 | DEMO GRAPHICS | LIFE STYLES EQUIPMENT GOODS CONSUMPTION | PRESS | RADIO | MAGAZINES | TV | OTHERS Internet Cine Outdoor |
| | + | + | + | + | + | + | + | + |
| MONOMEDIA PRESS | 15.000 | DEMO GRAPHICS | | PRESS | | | | |
| | + | + | | + | + | + | + | + |
| MONOMEDIA RADIO | 12.333 | DEMO GRAPHICS | | | RADIO | | | |
| | + | + | + | + | + | + | + | + |
| MONOMEDIA MAGAZINES | 6.666 | DEMO GRAPHICS | | | | MAGAZINES | | |
| | = | = | = | = | = | = | = | = |
| TOTAL | 43.999 | DEMO GRAPHICS | LIFE STYLES EQUIPMENT GOODS CONSUMPTION | PRESS | RADIO | MAGAZINES | TV | OTHERS Internet Cine Outdoor |

Goal: **Single data file**

**QUINAO**
consultores

2

1

**ODEC**

## INFORMATION AVAILABLE:

| | Interviews: | | LIFE STYLES EQUIPMENT GOODS CONSUMPTION | PRESS | RADIO | MAGAZINES | TV | OTHERS Internet Cine Outdoor |
|---|---|---|---|---|---|---|---|---|
| MULTIMEDIA | 10.000 | DEMO GRAPHICS | LIFE STYLES EQUIPMENT GOODS CONSUMPTION | PRESS | RADIO | MAGAZINES | TV | OTHERS Internet Cine Outdoor |
| MONOMEDIA PRESS | 15.000 | DEMO GRAPHICS | | PRESS | | | | |
| MONOMEDIA RADIO | 12.333 | DEMO GRAPHICS | | | RADIO | | | |
| MONOMEDIA MAGAZINES | 6.666 | DEMO GRAPHICS | | | | MAGAZINES | | |
| TOTAL | 43.999 | DEMO GRAPHICS | LIFE STYLES EQUIPMENT GOODS CONSUMPTION | PRESS | RADIO | MAGAZINES | TV | OTHERS Internet Cine Outdoor |

**COMMON INFORMATION**

**DATA SOURCES AVAILABLE: 1, 2, 3 or 4**

### STARTING POINT

- The problem posed is an imputation problem: "FILLING IN THE GAPS" or "MISSING VALUES".

- Subject to certain restrictions:

"Figures resulting from the fused final file (audiences, behaviors, demographics, …) must equal those resulting from the official source, those measured".

### AIMED GOAL: Single Data File

**FUSION**

| | Interviews | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| MULTIMEDIA | 10.000 | DEMO GRAPHICS | LIFE STILES EQUIPMENT GOODS CONSUMPTION | PRESS | RADIO | MAGAZINES | TV | OTHERS Internet Cine Outdoor |
| MONOMEDIA PRESS | 15.000 | DEMO GRAPHICS | LIFE STILES EQUIPMENT GOODS CONSUMPTION | PRESS | RADIO | MAGAZINES | TV | OTHERS Internet Cine Outdoor |
| MONOMEDIA RADIO | 12.333 | DEMO GRAPHICS | LIFE STILES EQUIPMENT GOODS CONSUMPTION | PRESS | RADIO | MAGAZINES | TV | OTHERS Internet Cine Outdoor |
| MONOMEDIA MAGAZINES | 6.666 | DEMO GRAPHICS | LIFE STILES EQUIPMENT GOODS CONSUMPTION | PRESS | RADIO | MAGAZINES | TV | OTHERS Internet Cine Outdoor |
| TOTAL | 43.999 | DEMO GRAPHICS | LIFE STILES EQUIPMENT GOODS CONSUMPTION | PRESS | RADIO | MAGAZINES | TV | OTHERS Internet Cine Outdoor |

**QUINAO** consultores

---

**ODEC**

Figures resulting from the fused final file (audiences, behaviors, demographics, …) must equal those resulting from the official source, those measured

¿What does this mean?

One goal is to make the most of the information available

Why are monomedia extensions to the multimedia survey made?

- for getting different aggregate audience figures?
- for getting better estimates when we segment or disaggregate the population?

**Reach:**

Let us suppose that the variables used to define a target group have different distribution of frequencies in the different surveys:

Let us suppose a magazine whose target group is people with a university degree. And let us suppose too that the number of these people is different according to the estimates of the different sources.
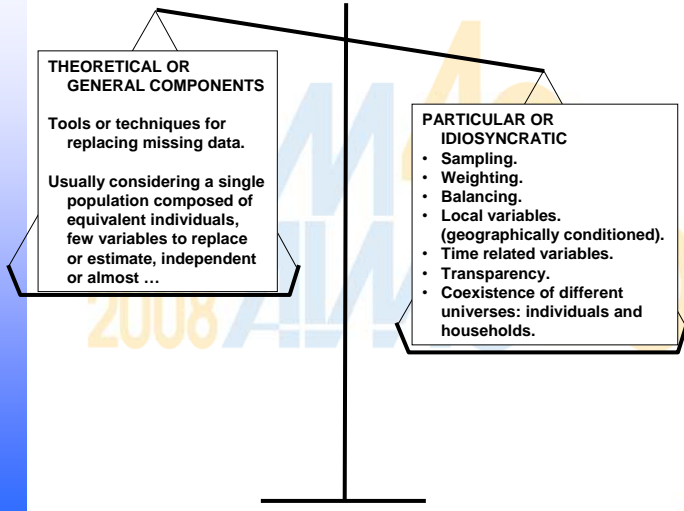
| | MULTIMEDIA (Observed data) | MONOMEDIA RADIO | |
|---|---|---|---|
| | | Estimates based on reach | Estimates based on readers |
| Owners of a university degree | 1.600.000 | 1.200.000 | 1.200.000 |
| Reach (%) | 20% | 20% | 27% |
| Readers | 320.000 | 240.000 | 320.000 |

Is it possible to obtain the same audience figures when files are different? If I make reach in percentage equal in both files then readership in raw figures, number of readers is different.

**QUINAO** consultores

**ODEC**

## PROBLEM COMPONENTS

**THEORETICAL OR GENERAL COMPONENTS**

**Tools or techniques for replacing missing data.**

**Usually considering a single population composed of equivalent individuals, few variables to replace or estimate, independent or almost …**

**PARTICULAR OR IDIOSYNCRATIC**
- **Sampling.**
- **Weighting.**
- **Balancing.**
- **Local variables. (geographically conditioned).**
- **Time related variables.**
- **Transparency.**
- **Coexistence of different universes: individuals and households.**

**QUINAO** consultores

5

---

**ODEC**

## THEORETICAL COMPONENTS

FUSION

TRANSPLANT                    ESTIMATION

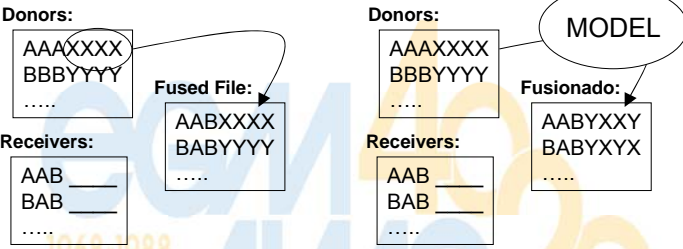**Possible APPROACHES:**
- **Pure fusion**
- **Estimation**

**PURE fusion preserves the structures in the data, correlations among data and is more suited when little information is available to build or fit a model.**

**PURE fusion allows working directly on the raw data, normally using condensed coding of the information (codes of the newspapers read yesterday), while modeling requires create real variables (read or doesn't read for each magazine)**

**Donors:**

AAAXXXX
BBBYYYY
…..

**Fused File:**

AABXXXX
BABYYYY
…..

**Receivers:**

AAB ____
BAB ____
…..

**Donors:**

AAAXXXX
BBBYYYY
…..

MODEL

**Fusionado:**

AABYXXY
BABYXYX
…..

**Receivers:**

AAB ____
BAB ____
…..

Techniques:
- Random hot deck
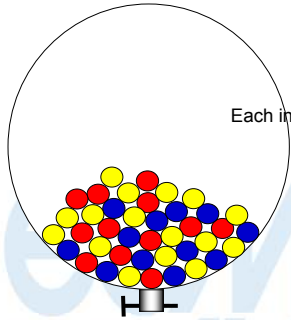- Distance among individuals (KNN)
- Factor reference based fusion
- Segmentation

Techniques:
- Discriminant Analysis
- Logistic Regression
- Artificial Neural Networks
- Bayesian models / nets

**QUINAO** consultores

6

3

**ODEC**

## PROBLEM COMPONENTS AND ADOPTED SOLUTION

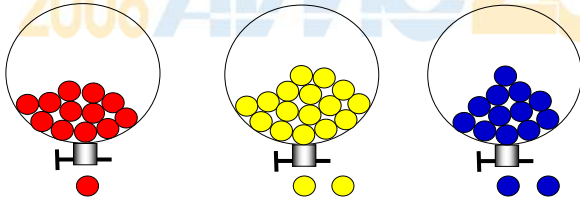1.   **INFERENCE / ESTIMATION : WEIGHTING AND BALANCING**

2.   **COMMON STRATIFICATION**

3.   **INFORMATION TRANSFER: MULTIPLE ASSIGNMENT**

**QUINAO**
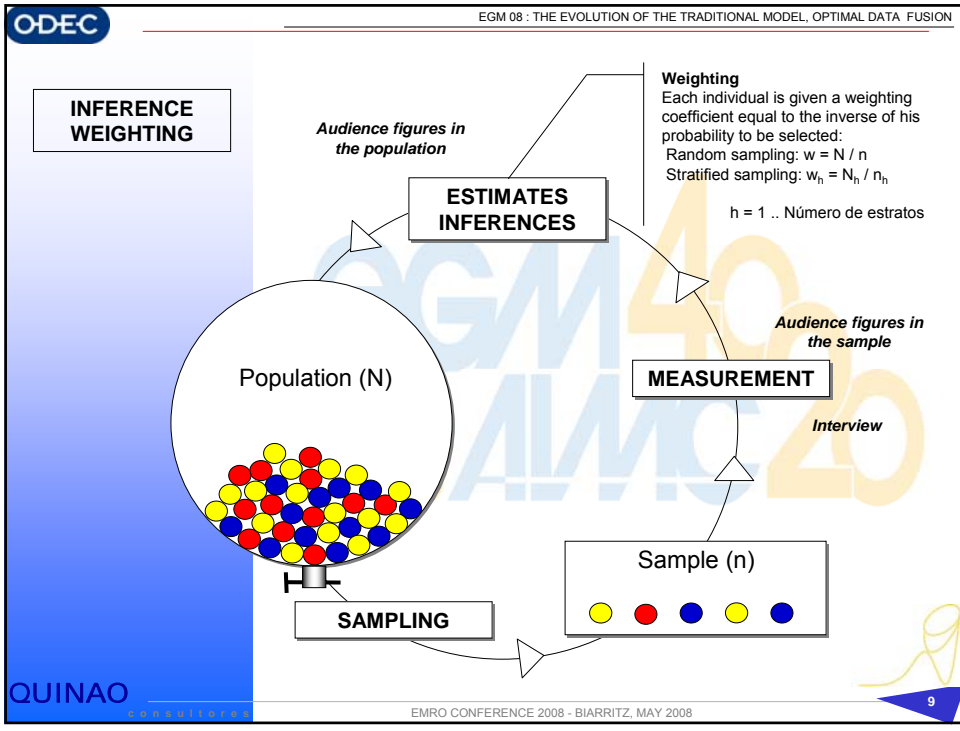consultores

7

---

**ODEC**

## SAMPLING

Absolutely random sampling
Each individual has the same probability to be selected

Stratified sampling
Population is split into different strata, having all the individuals within the same stratus the same probability to be selected.

**QUINAO**
consultores

8

4

**ODEC**

**INFERENCE WEIGHTING**

*Audience figures in the population*

**Weighting**
Each individual is given a weighting coefficient equal to the inverse of his probability to be selected:
Random sampling: $w = N / n$
Stratified sampling: $w_h = N_h / n_h$

$h = 1 ..$ Número de estratos

**ESTIMATES INFERENCES**

*Audience figures in the sample*

**MEASUREMENT**

*Interview*

Population (N)

**SAMPLING**

Sample (n)

**QUINAO**
consultores

9

---

**ODEC**

**SAMPLE BALANCING**

Population

Sample (n)

We want the frequencies distribution of different variables in the sample to fit exactly known distributions

**Sample Balancing**

Global estimators are not always better but better looking.

≠

| | Gender | | Age | | | Town size | | | N |
|---|---|---|---|---|---|---|---|---|---|
| | Male | Female | 14-24 | 25-64 | 65 + | Less than 50. | More than 50. | Capitals | |
| **Population** | 40 | 60 | 35 | 40 | 25 | 10 | 30 | 60 | 100 |
| **Survey** | Male | Female | 14 a 24 | 25 a 64 | 65 y mas | Less than 50. | More than 50. | Capitals | n |
| | 1 | | | | | 1 | | | 1 |
| | 1 | | 1 | | | | 1 | | 2 |
| | 1 | | 1 | | | | | 1 | 10 |
| | 1 | | | 1 | | 1 | | | 3 |
| | 1 | | | 1 | | | 1 | | 4 |
| | 1 | | | 1 | | | | 1 | 6 |
| | 1 | | | | 1 | 1 | | | 1 |
| | 1 | | | | 1 | | 1 | | 8 |
| | 1 | | | | 1 | | | 1 | 1 |
| | | 1 | | | | 1 | | | 9 |
| | | 1 | 1 | | | | 1 | | 6 |
| | | 1 | 1 | | | | | 1 | 2 |
| | | 1 | | 1 | | 1 | | | 8 |
| | | 1 | | 1 | | | 1 | | 8 |
| | | 1 | | 1 | | | | 1 | 8 |
| | | 1 | | | 1 | 1 | | | 6 |
| | | 1 | | | 1 | | 1 | | 8 |
| | | 1 | | | 1 | | | 1 | 9 |
| **Unweighted total** | 36 | 64 | 30 | 37 | 33 | 28 | 36 | 36 | 100 |

Difference cause:

**Random**     **Systematic**

**QUINAO**
consultores

10

5

**ODEC**

## BALANCING

Population

Sample (n)

Balancing is sometimes called:

**POST STRATIFICATION**, because it creates strata in the sample joined by those individuals with the same values in all the balancing variables and because of this with the same weighting or balancing coefficient.

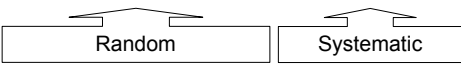**RAKING,** after the algorithm usually applied to turn the sample balanced.

| | Gender | | Age | | | Town size | | | N |
|---|---|---|---|---|---|---|---|---|---|
| | Male | Female | 14-24 | 25-64 | 65 + | Less than 50. | More than 50. | Capitals | |
| Population | 40 | 60 | 35 | 40 | 25 | 10 | 30 | 60 | 100 |
| Survey | Male | Female | 14 a 24 | 25 a 64 | 65 y mas | Less than 50. | More than 50 | Capitals | n |
| | 1 | | | | 1 | | | 1 | 1 |
| | 1 | | | 1 | | | | | 2 |
| | 1 | | | 1 | | | | | 10 |
| | | | | | 1 | 1 | | | 3 |
| | | | | 1 | | | | | 4 |
| | 1 | | | 1 | | | | | 6 |
| | 1 | | | | | | | | 1 |
| | | | | 1 | | | | | 8 |
| | | | | | | | | | 1 |
| | | 1 | 1 | | | | 1 | | 9 |
| | | 1 | 1 | | | | 1 | | 6 |
| | | 1 | 1 | | | | | 1 | 2 |
| | | 1 | | 1 | | 1 | 1 | | 8 |
| | | 1 | | 1 | | 1 | | | 8 |
| | | 1 | | 1 | | | | 1 | 8 |
| | | 1 | | | 1 | 1 | | | 6 |
| | | 1 | | 1 | | | 1 | | 9 |
| | | 1 | | | 1 | | | 1 | 9 |
| Unweighted total | 36 | 64 | 30 | 37 | 33 | 28 | 36 | 36 | 100 |

W. Edwards DEMING

Difference cause:     Random     Systematic

**QUINAO** consultores

11

---

**ODEC**

## JOINT POST STRATIFICATION

Population

Sample (n)

Sample (n)

Sample (n)
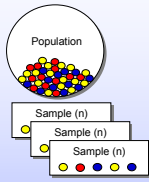
Strata are integrated by individuals absolutely equivalent and being affected equal weighting coefficients and so transfer of the information is done between equivalent individuals.

Stratus:
.
.
.
h
.
.
.

**MULTIMEDIA**

Sample (n)

Women / > 65 years old / capitals
Population 1000
Sample: 20
Weighting coefficient: 50

Readers magazine. X:
Population: 100
Sample: 2

**MONOMEDIA**

Sample (n)

Women / > 65 years old / capitals
Population 1000
Sample: 50
Weighting coefficient: 20

Readers magazine X:
Population: ?
Sample: ?

FUSION

Readers magazine X:
Population: 100

If we could make this transfer, put in the receiver file the information we have in the donor file, this information will be same not only for the equivalent strata in both file but for all their possible combinations

**QUINAO** consultores

12

6

**ODEC**

### INFORMATION TRANSFER

Population

Sample (n)
Sample (n)
Sample (n)

MULTIMEDIA

MONOMEDIA

Stratus:
.
.
.
h
.
.
.

Sample (n)

Sample (n)

Women / > 65 years old / capitals
Population: 180
Sample: **2**
Weighting coefficient: 90

Women / > 65 years old / capitals
Population: 180
Sample: **3**
Weighting coefficient: 60
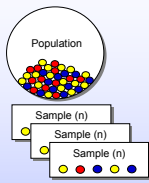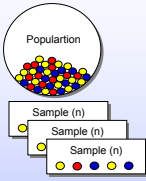
Common | To be fused

| ABCD | XXXXXX |
| ABED | YYYYYY |

| ABED | _ _ _ _ _ _ |
| ABCE | _ _ _ _ _ _ |
| ABCE | _ _ _ _ _ _ |

Control:
Defining the strata
and being equal.

Free:
They don't take part in the strata
definition but make individuals
more or less similar

It is not possible to transfer the
information coming from two
individuals to three and obtaining
the same results

**QUINAO**
consultores

13

---

**ODEC**

### INFORMATION TRANSFER

Population

Sample (n)
Sample (n)
Sample (n)

MULTIMEDIA

MONOMEDIA

Stratus:
.
.
h
.
.
.

Sample (n)

Sample (n)

Women / > 65 years old / capitals
Population: 180
Sample: **6**
Weighting coefficient: 30

Women / > 65 years old / capitals
Population: 180
Sample: **6**
Weighting coefficient: 30

In order to transfer information
exactly we need to make a
MULTIPLE IMPUTATION
(Rubin): replicate the individuals
in each survey stratus.

Common | To be fused

| ABCD | XXXXXX | → | ABED | XXXXXX |
| ABED | YYYYYY | → | ABCE | YYYYYY |
| ABCD | XXXXXX | → | ABCE | XXXXXX |
| ABED | YYYYYY | → | ABED | YYYYYY |
| ABCD | XXXXXX | → | ABCE | XXXXXX |
| ABED | YYYYYY | → | ABCE | YYYYYY |

**QUINAO**
consultores

14

7

**ODEC**

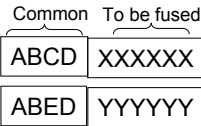**FUSION**

Populartion

Sample (n)
Sample (n)
Sample (n)

Now:

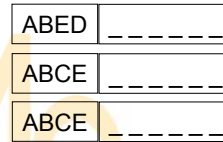Which receiver is going to be given the information from which donor?

**Nearest Neighbor Procedure:** distance among all donors and all receivers is computed and transfer is made between those being the closest.
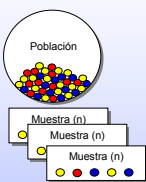
**DONOR**

Common — To be fused

| ABCD | XXXXXX |
| ABED | YYYYYY |

**RECEIVER**

| ABED | _ _ _ _ _ _ |
| ABCE | _ _ _ _ _ _ |
| ABCE | _ _ _ _ _ _ |

Receivers:

| Donors: | ABED | ABCE | ABCE | ABED | ABCE | ABCE |
|---|---|---|---|---|---|---|
| ABCD |  |  |  |  |  |  |
| ABED |  | 8 |  |  |  |  |
| ABCD | 2 | 6 | 3 |  |  |  |
| ABED |  | 1 | 7 |  |  |  |
| ABCD |  |  | 5 |  |  |  |
| ABED |  |  |  |  |  |  |

1) **Distance matrix** is computed and distances are rank ordered from greatest to lowest.
2) Pair with the lowest distance between them is selected, information from donor is transferred to receiver and both, donor and receiver, are deleted from the table.

**QUINAO**
consultores

---

**ODEC**

**EGM BALANCING**

Población

Muestra (n)
Muestra (n)
Muestra (n)

**Individuals balancing matrixes**

| | | | Constrains |
|---|---|---|---|
| Province 50 | | Town size 2 | 100 |
| Region 17 | x | Town size 7 | 119 |
| Region 17 | x | Gender 2 | 34 |
| Region 17 | x | Edad 7 | 119 |
| Age 7 | x | Sexo 2 | 14 |
| Weekday 7 | x | Ama 2 | 14 |
| Housewife 2 | | | 2 |
| Regiony 9 | x | Family size 3 | 27 |
| Regiony 9 | x | Weekday 2 | 18 |
| | | **Total** | **447** |

Different Individuals Possible

| Province x 50 | Town size x 7 | Gender x 2 | Age x 7 | Weekday x 7 | Housewife x 2 | Head x 2 | Town size 3 | = | = 411.600 | Same weighti |
|---|---|---|---|---|---|---|---|---|---|---|

**QUINAO**
consultores

**ODEC**

## JOINT STRATIFICATION INTO SAME SIZE STRATA

Population

Sa,ple (n)

Muestra (n)

Muestra (n)

**TWO APPROACHES:**

1. LOOKING FOR THE THINNEST SPLIT INTO STRATA COMMON TO DONOR AND RECEIVER:

   Segmentation techniques looking to split both donor and surveys in as many COMMON STRATA AS POSSIBLE. These strata are going to have probably a different aggregated weight in each survey.

   Size of the strata unknown to be estimated from the file resulting joining both DONOR and RECEIVER and computing weights.
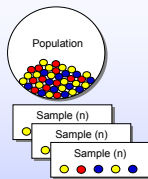
   REWEIGHTING THE SURVEYS for meeting all the restrictions plus the new restriction that strata must weight the same.

2. PREDEFINED STRATA OF KNOWN SIZE IN THE POPULATION.

   PROVINCE (50) x  TOWN SIZE (2) x GENDER (2) = 200

   Adding this new restrictions to our balancing matrixes and recomputing weighting coefficients for all the surveys.

**QUINAO**
consultores

---

**ODEC**

## JOINT STRATIFICATION INTO SAME SIZE STRATA

Population

Sample (n)

Sample (n)

Sample (n)

**SOLUTION:**

1. NATURAL STRATIFICATION (PREDEFINED):

   PROVINCE (50) X TOWN SIZE (2) X GENDER (2) = 200 STRATA

   PROVINCE (50) X WEEKDAY (2) X GENDER (2) = 200 STRATA

   (Depending on the data being fused)

2. DISTANCES AMONG DONORS AND RECEIVERS BASED ON
   AGE
   ROL
   SOCIAL STATUS
   TOWN SIZE
   WEEKDAY
   HOUSEHOLD SIZE
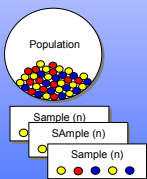   TOWN
   LANGUAGE
   ….

**QUINAO**
consultores

**ODEC**

**TRASFER OF INFORMATION AMONG INDIVIDUALS NOT WEIGHTING THE SAME**

$$\sum w_d = \sum w_r$$

Population

Sample (n)
SAmple (n)
Sample (n)

1) In the real situation individuals in the stratus would have different weighting coefficients, although the sum of all the coefficients should be the same for the same stratus in every source.

2) Distances matrix among donors and receivers is computed using control variables and non control but common variables

3) Distances are rank ordered and the pair of most similar individuals is selected.

**For each stratus h**

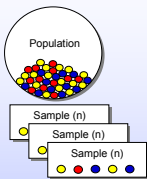| | | | | Receivers | | | | |
|---|---|---|---|---|---|---|---|---|
| | | 1 | 2 | … | … | j | … | q |
| | Weight | wr1 | wr2 | | | wrj | | wrq |
| **Donors** | Weight | | | | | | | |
| 1 | wd1 | | | | | | | |
| 2 | wd2 | | | | Distance Matrix | | | |
| . | … | | | | | | | |
| . | … | | | | | | | |
| i | wdi | | | | | | | |
| . | … | | | | | | | |
| . | … | | | | | | | |
| **p** | wdp | | | | | | | |

4) If donor weight $wd_i$ is greater than receiver weight $wr_j$, information from the donor is transferred, pasted to the receiver, receiver is pasted to the fused file with its original weighting coefficient and donor remains in the table as a potential donor with a weight equal to its original weight minus the weight of the receiver having just got his information.

5) If donor weight is lower than receiver, then information from the donor is transferred to the receiver, and receiver is written to the fused file with a weight equal to that of the donor, donor is deleted from the table, it has already transferred all his information, and receiver remains in the table as a potential receiver with a weight equal to its original weight minus the weight of the donor.

6) Process continues till neither donor nor receivers are available.

**QUINAO**
consultores

---

**ODEC**

**FUSION**

## DONOR and RECEIVERS HAVE DIFFERENT WEIGHTS

Population

Sample (n)
Sample (n)
Sample (n)

1) **Distance matrix is computed and distances are rank ordered from greatest to lowest.**

2) **Pair with the lowest distance between them is selected.**

**Donors and their weights:**

**Receivers and their weights:**

**Total Receivers weight:**

| Donors | 3,10 | 0,60 | 1,40 | 1,80 | 1,10 | 8,00 |
|---|---|---|---|---|---|---|
| 3,00 | | | | | | |
| 1,50 | | 8 | | | | |
| 2,00 | 2 | 6 | 3 | | | |
| 1,00 | | 1 | 7 | | | |
| 0,50 | | | 5 | | | |

**8,00**
**Total weight**

**Distance Matrix**

**QUINAO**
consultores

10

## DONORS and RECEIVERS HAVE DIFFERENT WEIGHTS

**FUSION**

Population

Sample (n)
Sample (n)
Sample (n)

**Receivers and their weights:**

**Total Receivers weight:**

**Donors and their weights:**

| | 3,10 | 0,60 | 1,40 | 1,80 | 1,10 | 8,00 |
|---|---|---|---|---|---|---|
| 3,00 | | | | | | |
| 1,50 | | 8 | | | | |
| 2,00 | 2 | 6 | 3 | | | |
| 1,00 | | 1 | 7 | | | |
| 0,50 | | | 5 | | | |
| 8,00 | | | | | | |

**Total weight**

**Distance Matrix**

**Most similar pair: Donor weight greater than receiver weight**
1) Receiver is pasted donor information
2) Receiver is written to the fused file with its own weight and deleted from the distance table
3) Donor remains in the table with a weight equal to the difference of weights.

---

**FUSION**

Population

Sample (n)
Sample (n)
Sample (n)

**RESULTS:**

Donor file, Receiver file and Fused file contain **exactly the same information** in the imputed variables, and this for all the common strata and for all their possible aggregations.

Internal **relations among fused variables are kept** and are the same for all files and surveys.

For those variable not controlled, distributions are as similar as possible.

**Traceability** is possible, one can know exactly how many times each record is replicated, and how original interviews are the base for each data.

The thinnest the strata the more similar information will result in every file but the possibility of finding individuals similar in very small strata decreases. **A compromise must be made among number and size of strata**.

11

**ODEC**

## SINCRONIZING

Interviews

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| MULTIMEDIA | 10.000 | DEMO GRAPHICS | LIFE STILES EQUIPMENT GOODS CONSUMPTION | PRESS | RADIO | MAGAZINES | TV | OTHERS Internet Cine Outdoor |
| + | + | + | + | + | + | + | + | + |
| MONOMEDIA PRESS | 15.000 | DEMO GRAPHICS | LIFE STILES EQUIPMENT GOODS CONSUMPTION | PRESS | RADIO | MAGAZINES | TV | OTHERS Internet Cine Outdoor |
| + | + | + | + | + | + | + | + | + |
| MONOMEDIA RADIO | 12.333 | DEMO GRAPHICS | LIFE STILES EQUIPMENT GOODS CONSUMPTION | PRESS | RADIO | MAGAZINES | TV | OTHERS Internet Cine Outdoor |
| + | + | + | + | + | + | + | + | + |
| MONOMEDIA MAGAZINES | 6.666 | DEMO GRAPHICS | LIFE STILES EQUIPMENT GOODS CONSUMPTION | PRESS | RADIO | MAGAZINES | TV | OTHERS Internet Cine Outdoor |
| = | = | = | = | = | = | = | = | = |
| TOTAL | 120.000 | DEMO GRAPHICS | LIFE STILES EQUIPMENT GOODS CONSUMPTION | PRESS | RADIO | MAGAZINES | TV | OTHERS Internet Cine Outdoor |

**STEP 1:**

- Information from Multimedia,  TO Multimedia + Press + Radio + magazines

**STEP 1(A):**

- Information for HOUSEWIFE from Multimedia,  TO Multimedia + Press + Radio + magazines

**STEP 2:**

- Information from Multimedia + RADIO TO Multimedia + Press + Radio + magazines

· · · · · ·

Size of the final file is roughly three times the size resulting from simply joining the surveys

**QUINAO** consultores

---

**ODEC**

## AUTOMATIZATION



**QUINAO** consultores

12

**ODEC**

CHARACTERISTICS OF THE PROCESS:

- EFFICIENT

- NO NEED TO MAKE DECISSIONS

- TRANSPARENT

- OPTIMAL

STEPS

- Common stratification.

- Weighting and balancing: same weight for strata

- Transfer of information: distance based with multiple replication

**QUINAO** consultores

---

**ODEC**

SOME RESULTS:

EGM: 1ST WAVE 2008
SPAIN

| | FINAL FILE AFTER FUSION | EGM MULTIMEDIA | RADIO | PRESS | MAGAZINES |
|---|---|---|---|---|---|
| Total POPULATION | 38.261 | 38.261 | 38.261 | 38.261 | 38.261 |
| **LAST PERIOD READERS (000)** | | | | | |
| Daily newspapers | **28.876** | 29.647 | | 28.876 | |
| Sunday suplements | **11.767** | **11.765** | | | |
| Weekly magazines | **21.394** | 18.583 | | | 21.394 |
| Biweekly magazines | **2.802** | 2.094 | | | 2.802 |
| Monthly magazines | **21.660** | 18.869 | | | 21.660 |
| Total Magazines | **28.087** | 25.874 | | | 28.086 |
| **RADIO LISTENERS** | | | | | |
| TOTAL TEMÁTICA | **10.845** | 11.807 | 10.845 | | |
| TEMÁTICA MUSICAL | **9.749** | 10.601 | 9.749 | | |
| TEMÁTICA INFORMATIVA | **1.300** | 1.391 | 1.300 | | |
| OTRAS TEMÁTICAS | **73** | 69 | 73 | | |
| C40 | **3.140** | 3.266 | 3.140 | | |
| Dial | **1.530** | 1.706 | 1.530 | | |
| C100 | **975** | 901 | 975 | | |
| M80 | **533** | 553 | 533 | | |
| **CINEMA** | **17.651** | 17.653 | | | |
| INTERNET | **17.549** | 17.554 | | | |

**QUINAO** consultores