Swiss TPH
Swiss Tropical and Public Health Institute
Schweizerisches Tropen- und Public Health-Institut
Institut Tropical et de Santé Publique Suisse
Associated Institute of the University of Basel

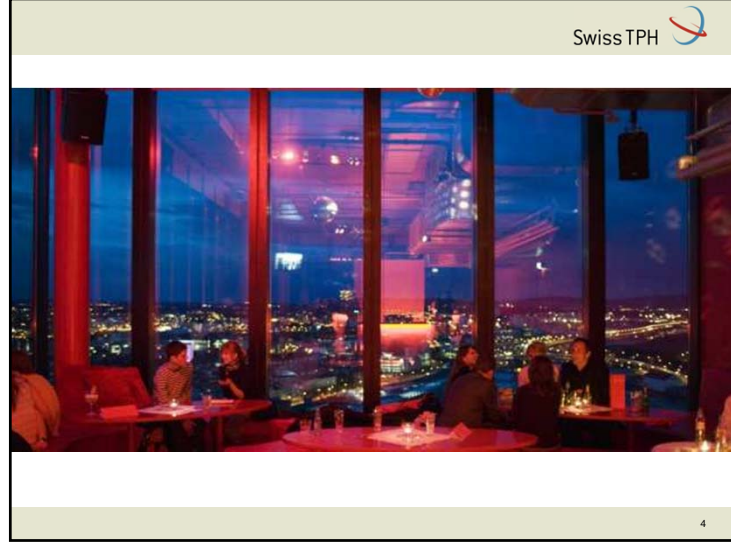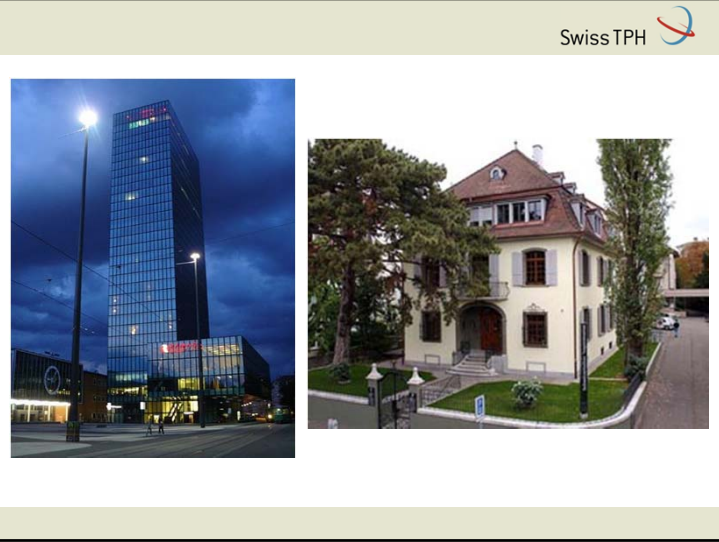**2014 Spanish Stata Users Group meeting**

# Automated harmonisation of variables names and values from several datasets prior to conducting batch statistical analyses

Barcelona, 23rd October 2014

Xavier Bosch-Capblanch
Swiss TPH, Basel (Switzerland)

---



Swiss TPH

---



Swiss TPH

---



Swiss TPH

4

THE NEED AND THE PROBLEM

THE FIRST ATTEMPTS

THE SOLUTION: HARMONISATION

WAY FORWARD

Bosch-Capblanch X. Harmonisation of variables names prior to conducting statistical analyses with multiple datasets: an automated approach. BMC Med Inform Decis Mak. 2011 May 19;11:33. doi: 10.1186/1472-6947-11-33. PubMed PMID: 21595905; PubMed Central PMCID: PMC3123542.

---

THE NEED AND THE PROBLEM

$$\text{Coverage} = \frac{\text{Number of children age A vaccinated with V}}{\text{Number of children age A targetted by V}}$$

---



---

Pakistan



| | | |
|---|---|---|
| × Offical country estimates | ● Administrative data | |
| ■ EPI | ▲ DHS secondary data | |
| ■ DHS primary data | ▲ WHO/UNICEF estimates | |

**Slide 1 (top-left):**

## Global Immunization 1980-2012, DTP3 coverage
### global coverage at 83% in 2012

% coverage

Global | African | American | Eastern Mediterranean | European | South East Asian | Western Pacific

Source: WHO/UNICEF coverage estimates 2012 revision. July 2013
Immunization Vaccines and Biologicals, (IVB), World Health Organization.
194 WHO Member States. Date of slide: 27 August 2013.

unicef | World Health Organization

---

**Slide 2 (top-right):**

511 household surveys (DHS, MICS)
(N from 1,000 Trinidad and Tobago to 51,000 India)   Swiss TPH

Afghanistan 2011: Afghanistan 2011; Albania 2009; Albania 2000; Albania 2000; Angola 2001; Angola 2001; Armenia 2000; Armenia 2000; Armenia 2005; Armenia 2010; Azerbaijan 2006; Azerbaijan 2000; Azerbaijan 2000; Bangladesh 1994; Bangladesh 1994; Bangladesh 1996; Bangladesh 1996; Bangladesh 2000; Bangladesh 2000; Bangladesh 2004; Bangladesh 2007; Bangladesh 2011; Bangladesh 2006; Bangladesh 2006; Belarus 2005; Belize 2006; Belize 2006; Belize 2011; Belize 2011; Benin 1996; Benin 1996; Benin 2001; Benin 2001; Benin 2006; Bhutan 2010; Bhutan 2010; Bolivia 1989; Bolivia 1994; Bolivia 1998; Bolivia 1998; Bolivia 2003; Bolivia 2008; Bolivia 2000; Bosnia and Herzegovina 2000; Bosnia and Herzegovina 2000; Bosnia and Herzegovina 2006; Bosnia and Herzegovina 2011; Bosnia and Herzegovina 2011; Brazil 1986; Brazil 1991; Brazil 1996; Brazil 1996; Burkina Faso 1993; Burkina Faso 1993; Burkina Faso 1999; Burkina Faso 1999; Burkina Faso 2003; Burkina Faso 2010; Burkina Faso 2006; Burkina Faso 2006; Burundi 1987; Burundi 2010; Burundi 2000; Burundi 2000; Burundi 2005; Burundi 2005; Cambodia 2000; Cambodia 2005; Cambodia 2010; Cameroon 1991; Cameroon 1991; Cameroon 1998; Cameroon 1998; Cameroon 2004; Cameroon 2011; Cameroon 2000; Cameroon 2000; Cameroon 2006; Central African Republic 1994; Central African Republic 1994; Central African Republic 2000; Central African Republic 2000; Central African Republic 2006; Chad 1997; Chad 1997; Chad 2004; Chad 2000; Colombia 1986; Colombia 1990; Colombia 1990; Colombia 1995; Colombia 1995; Colombia 2000; Colombia 2000; Colombia 2005; Colombia 2010; Comoros 1996; Comoros 1996; Comoros 2000; Comoros 2000; Congo 2005; Congo 2005; Congo DR 2007; Congo DR 2001; Congo DR 2001; Congo DR 2010; Congo DR 2010; Cuba 2000; Cuba 2006; Cuba 2006; Cuba 2010; Cuba 2010; Côte d'Ivoire 1994; Côte d'Ivoire 1994; Côte d'Ivoire 1999; Côte d'Ivoire 1999; Côte d'Ivoire 2000; Côte d'Ivoire 2000; Côte d'Ivoire 2006; Côte d'Ivoire 2006; Djibouti 2006; Djibouti 2006; Dominican Republic 1986; Dominican Republic 1991; Dominican Republic 1996; Dominican Republic 1996; Dominican Republic 1999; Dominican Republic 1999; Dominican Republic 2002; Dominican Republic 2002; Dominican Republic 2007; Dominican Republic 2000; Dominican Republic 2000; Ecuador 1987; Egypt 1988; Egypt 1992; Egypt 1995; Egypt 1995; Egypt 2000; Egypt 2000; Egypt 2003; Egypt 2003; Egypt 2008; Ethiopia 2000; Ethiopia 2000; Ethiopia 2005; Ethiopia 2011; Gabon 2000; Gabon 2000; Gambia 2000; Gambia 2000; Gambia 2006; Gambia 2006; Georgia 2005; Georgia 2005; Ghana 1988; Ghana 1993; Ghana 1993; Ghana 1998; Ghana 1998; Ghana 2003: Ghana 2008; Ghana 2006; Ghana 2006; Guatemala 1987; Guatemala 1995; Guatemala 1995; Guatemala 1999; Guatemala 1999; Guinea 1999; Guinea 1999; Guinea 2005; Equatorial Guinea 2000; Equatorial Guinea 2000; Guinea-Bissau 2000; Guinea-Bissau 2006; Guinea-Bissau 2006; Guinea-Bissau 2006; Guyana 2009; Guyana 2000; Guyana 2000; Guyana 2006; Guyana 2006; Haiti 1994; Haiti 1994; Haiti 2000; Haiti 2000; Haiti 2006; Honduras 2006; India 1993; India 1993; India 1999; India 1999; India 2006; Indonesia 1987; Indonesia 1991; Indonesia 1994; Indonesia 1997; Indonesia 1997; Indonesia 2002; Indonesia 2007; Indonesia 2000; Indonesia 2000; Iraq 2000; Iraq 2000; Iraq 2006; Iraq 2006; Iraq 2011; Iraq 2011; Jamaica 2000; Jamaica 2000; Jamaica 2005; Jamaica 2005; Jordan 1990; Jordan 1990; Jordan 1997; Jordan 1997; Jordan 2002; Jordan 2002; Jordan 2007; Kazakhstan 1995; Kazakhstan 1995; Kazakhstan 1999; Kazakhstan 1999; Kazakhstan 2006; Kazakhstan 2006; Kazakhstan 2010; Kazakhstan 2010; Kenya 1989; Kenya 1993; Kenya 1993; Kenya 1998; Kenya 1998; Kenya 2003; Kenya 2009; Kenya 2000; Kenya 2000; Kyrgyzstan 1997; Kyrgyzstan 1997; Kyrgyzstan 2005; Kyrgyzstan 2005; Lao People'S Democratic Republic 2000; Lao People'S Democratic Republic 2006; Lao People'S Democratic Republic 2000; Lao People'S Democratic Republic 2006; Lesotho 2004; Lesotho 2009; Lesotho 2000; Lesotho 2000; Liberia 1986; Liberia 2007; Macedonia 2005; Macedonia 2005; Madagascar 1992; Madagascar 1997; Madagascar 1997; Madagascar 2004; Madagascar 2009; Madagascar 2000; Madagascar 2000; Malawi 1992; Malawi 1992; Malawi 2000; Malawi 2000; Malawi 2004; Malawi 2010; Malawi 2006; Maldives 2009; Mali 1987; Mali 1996; Mali 1996; Mali 2001; Mali 2001; Mali 2006; Mauritania 2000; Mauritania 2007; Mauritania 2007; Mexico 1987; Moldova 2005; Moldova 2000; Moldova 2000; Mongolia 2000; Mongolia 2000; Mongolia 2005; Mongolia 2005; Montenegro 2005; Montenegro 2005; Morocco 1987; Morocco 1992; Morocco 1992; Morocco 2003; Mozambique 1997; Mozambique 1997; Mozambique 2003; Mozambique 2011; Mozambique 2008; Myanmar 2000; Myanmar 2000; Namibia 1992; Namibia 1992; Namibia 2000; Namibia 2000; Namibia 2007; Nepal 1996; Nepal 1996; Nepal 2001; Nepal 2001; Nepal 2006; Nepal 2010; Nicaragua 1998; Nicaragua 1998; Nicaragua 2001; Nicaragua 2001; Niger 1992; Niger 1998; Niger 1998; Niger 2006; Niger 2000; Niger 2000; Nigeria 1990; Nigeria 1990; Nigeria 1999; Nigeria 2003; Nigeria 2008; Nigeria 2008; Nigeria 2007; Nigeria 2007; Nigeria 2011; Nigeria 2011; Ondo State 1986; Pakistan 1991; Pakistan 1991; Pakistan 2006; Paraguay 1990; Paraguay 1990; Peru 1986; Peru 1991; Peru 1991; Peru 1996; Peru 1996; Peru 2000; Peru 2000; Peru 2006; Philippines 1993; Philippines 1998; Philippines 1998; Philippines 2003; Philippines 2008; Philippines 1999; Philippines 1999; Rwanda 1992; Rwanda 1992; Rwanda 2005; Rwanda 2005; Rwanda 2010; Rwanda 2000; Rwanda 2000; Sao Tome and Principe 2008; Sao Tome and Principe 2000; Sao Tome and Principe 2000; Senegal 1986; Senegal 1993; Senegal 1997; Senegal 1997; Senegal 2005; Senegal 2011; Serbia 2005; Serbia 2005; Serbia 2010; Serbia 2010; Sierra Leone 2008; Sierra Leone 2000; Sierra Leone 2000; Sierra Leone 2005; Sierra Leone 2005; Sierra Leone 2010; Sierra Leone 2010; Somalia 2006; Somalia 2006; South Africa 1998; Sri Lanka 1987; Sudan 1990; Sudan North 2000; Sudan South 2000; Sudan South 2000; Suriname 2000; Suriname 2000; Suriname 2006; Suriname 2000; Suriname 2010; Swaziland 2006; Swaziland 2000; Swaziland 2000; Swaziland 2010; Syrian Arab Republic 2006; Syrian Arab Republic 2006; Tajikistan 2000; Tajikistan 2000; Tajikistan 2005; Tanzania 1991; Tanzania 1996; Tanzania 1999; Tanzania 1999; Tanzania 2004; Tanzania 2010; Thailand 1987; Thailand 2006; Thailand 2006; Timor-Leste 2009; Togo 1988; Togo 1998; Togo 1998; Togo 2000; Togo 2006; Togo 2006; Togo 2010; Togo 2010; Trinidad and Tobago 1987; Trinidad and Tobago 1987; Trinidad and Tobago 2000; Trinidad and Tobago 2006; Tunisia 1988; Turkey 1993; Turkey 1993; Turkey 1998; Turkey 1998; Turkey 2004; Uganda 1988; Uganda 1995; Uganda 1995; Uganda 2001; Uganda 2006; Ukraine 2007; Ukraine 2005; Ukraine 2005; Uzbekistan 1996; Uzbekistan 1996; Uzbekistan 2000; Uzbekistan 2000; Uzbekistan 2006; Uzbekistan 2006; Vanuatu 2007; Vanuatu 2007; Venezuela 2000; Venezuela 2000; Viet Nam 1997; Viet Nam 1997; Viet Nam 2002; Viet Nam 2000; Viet Nam 2000; Viet Nam 2006; Viet Nam 2006; Viet Nam 2010; Yemen 1991; Yemen 2006; Yemen 2006; Zambia 1992; Zambia 1996; Zambia 1996; Zambia 2002; Zambia 2002; Zambia 2007; Zambia 1999; Zambia 1999; Zimbabwe 1988; Zimbabwe 1994; Zimbabwe 1994; Zimbabwe 1999; Zimbabwe 1999; Zimbabwe 2005; Zimbabwe 2010; Zimbabwe 2009; Zimbabwe 2009:

---

**Slide 3 (bottom-left):**

Bangladesh 2000 (337 variables)   Swiss TPH

hh1; br8co; va1; ca9o; im4by; im8bd; hh16; hc1a; hc11bcg; hc15g; hh2; br8cn; va2; ca9p; im4cd; im8bm; ws1; hc1c; hc11bdb; hc15ha; ln; br8dm; va3; ca9q; im4cm; im8by; ws2; hc2; hc11bdg; hc15hb; uf1; br8df; bf1; ca9r; im4cy; im9; ws3; hc3; hc11beb; hc15hc; uf2; br8do; bf1a; ca9s; im21ad; im10; ws4; hc4; hc11beg; hc15hd; uf4; br8dn; bf1bu; ca9x; im21am; im11; ws5; hc5; hc11bfb; hc15he; uf6; br8em; bf1bn; ca10; im21ay; im12; ws6a; hc6; hc11bfg; hc15hf; uf7; br8ef; bf2; ca11a; im21bd; im13; ws6b; hc7; hc11bgb; hc15hg; uf8d; br8eo; bf3a; ca11p; im21bm; im14; ws6c; hc7a; hc11bgg; hc15hy; uf8m; br8en; bf3b; ca11q; im21by; im15; ws6d; hc8; hc11bxb; hc15ia; uf8y; br8fm; bf3c; ca11r; im21cd; im16; ws6f; hc9a; hc11bxg; hc15ib; uf9; br8ff; bf3d; ca11x; im21cm; im17; ws6g; hc9b; hc11cab; hc15ic; uf10d; br8fo; bf3e; ca11z; im21cy; im19a; ws6x; hc9c; hc11cag; hc15id; uf10m; br8fn; bf3f; ca13; im3ad; im19b; ws6z; hc9d; hc11cbb; hc15ie; uf10y; br9; bf3g; ca14a; im3am; im19c; ws6_1; hc9e; hc11cbg; hc15if; uf11; br10a; bf3h; ca14b; im3ay; hl4; ws6_2a; hc9f; hc11ccb; hc15ig; uf11a; br10b; bf5; ca14c; im3bd; ed3a; ws6_2b; hc9g; hc11ccg; hc15iy; br1; br10c; ca1; ca14d; im3bm; ed3b; ws6_2c; hc9h; hc11d; hc15ja; br2; br10d; ca2a; ca14e; im3by; hh3; ws6_2d; hc9j; hc15; hc15jb; br3; br10e; ca2b; ca14f; im3cd; hh4; ws6_2e; hc9i; hc15a; hc15jy; br4; br10f; ca2c; ca14g; im3cm; hh5d; ws6_2x; hc10a; hc15b; chweight; br4a; br10g; ca3; ca14h; im3cy; hh5m; ws6_2z; hc10b; hc15ca; wlthscor; br6; br10h; ca4; ca14i; im3dd; hh5y; ws6_3a; hc10c; hc15cb; wlthind5; br7; br10x; ca5; ca14j; im3dm; hh6; ws6_3b; hc10d; hc15cc; cmcdoic; br8am; br11; ca6; ca14x; im3dy; hh7; ws6_3c; hc10e; hc15cx; cdob; br8af; br12a; ca7; im1; im3ed; hh7a; ws6_3d; hc10f; hc15cy; cage; br8ao; br12b; ca8; im2d; im3em; hh7b; ws6_3e; hc10g; hc15d; cage_6; br8an; br12c; ca9a; im2m; im3ey; hh9; ws6_3x; hc10h; hc15ea; cage_11; br8bm; br12d; ca9d; im2y; im6d; hh10; ws6_3z; hc11a; hc15eb; melevel2; br8bf; br12e; ca9e; im4ad; im6m; hh11; ws6_4; hc11bab; hc15ec; melevel; br8bo; br12f; ca9h; im4am; im6y; hh12; ws7; hc11bag; hc15ed; br8bn; br12g; ca9i; im4ay; hh13; ws8; hc11bbb; hc15ex; br8cm; br12h; ca9j; im4bd; im8am; hh14; ws9; hc11bbg; hc15ey; br8cf; br12x; ca9k; im4bm; im8ay; hh15; ws9a; hc11bcb; hc15f;

---

**Slide 4 (bottom-right):**

Variable: day vaccination DTP3 (month, year...)   Swiss TPH

NA; **im4cd**; NA; NA; **h7d**; h7d; h7d; h7d; im4cd; NA; h7d; NA; h7d; NA; h7d; NA; h7d; h7d; h7d; im4cd; NA; im4cd; NA; im4fd; NA; im3d3d; h7d; NA; h7d; NA; h7d; NF; h7d; h7d; NA; h7d; NA; h7d; h7d; im4cd; NA; im4cd; NA; im3d3d; h7d; h7d; NA; h7d; NA; h7d; h7d; im4cd; NA; h7d; h7d; im4cd; NA; NA; h7d; h7d; h7d; h7d; NA; h7d; NA; h7d; h7d; im4cd; NA; im4cd; NA; h7d; NA; im4cd; NA; im4cd; h7d; NA; h7d; im4cd; NA; h7d; h7d; NA; h7d; NA; h7d; NA; h7d; h7d; h7d; NA; im4cd; NA; h7d; h7d; NA; h7d; NA; h7d; NA; h7d; h7d; h7d; NA; im4cd; NA; h7d; h7d; im4cd; NA; im3d3d; NF; NA; im3d3d; NA; h7d; NA; im4cd; NA; im4cd; NA; im4cd; h7d; NA; h7d; NA; h7d; NA; h7d; h7d; NA; im4cd; NA; im4cd; NA; im4cd; NA; h7d; h7d; NA; h7d; h7d; NA; h7d; NA; h7d; NA; h7d; NA; h7d; h7d; NA; im4cd; NA; im4cd; NA; im5cd; NA; h7d; NA; h7d; h7d; h7d; NA; NA; h7d; h7d; h7d; h7d; NA; im4cd; NA; im4cd; NA; im3d3d; NA; im4cd; NA; h7d; NA; h7d; NA; NA; h7d; h7d; h7d; NA; h7d; NA; im4cd; NA; im4cd; NA; h7d; h7d; NA; h7d; h7d; NA; h7d; h7d; im4cd; NA; h7d; h7d; im4cd; NA; h7d; h7d; NA; h7d; NA; h7d; h7d; im5cd; NA; h7d; h7d; NA; h7d; NA; h7d; h7d; im4cd; NA; h7d; h7d; NA; im4cd; NA; h7d; NA; h7d; NA; h7d; h7d; NA; im3cd; NA; h7d; h7d; NA; h7d; h7d; NA; h7d; h7d; im4cd; NA; im4cd; h7d; NA; h7d; h7d; NA; h7d; NA; h7d; h7d; h7d; NA; h7d; NA; h7d; h7d; NA; h7d; NA; h7d; im4cd; NA; h7d; NA; h7d; **h7d**; h7d; im4cd; im3d3d; h7d; NA; h7d; h7d; NA; h7d; h7d; NA; h7d; NA; h7d; NA; h7d; h7d; h7d; NA; h7d; NA; h7d; NA; h7d; h7d; im4cd; NA; h7d; h7d; im4cd; NA; h7d; h7d; im3cd; **NA**; NF; h7d; im4cd; NA; im4cd; NA; im3d3d; im4cd; NA; h7d; NA; h7d; NA; im4cd; NA; im4cd; NA; im4cd; NA; im4cd; im3d3d; h7d; im4cd; NA; im3d3d; im4cd; NA; im4cd; NA; h7d; NA; h7d; NA; h7d; h7d; im4cd; NA; h7d; NA; im4cd; NA; im4cd; NA; h7d; h7d; h7d; NA; h7d; h7d; h7d; h7d; NA; im3d3d; h7d; im4cd; NA; im4cd; NA; h7d; NA; h7d; h7d; h7d; NA; h7d; h7d; NA; h7d; NA; h7d; h7d; **im3d3d**; h7d; im4cd; NA; h7d; h7d; NA; NA; h7d; h7d; im4cd; NA; h7d; h7d; NA; h7d; NA; h7d; h7d; **im5cd**;

12

## Slide 1

Value = 3 for variable DTP3 vaccination

Swiss TPH

**not vaccinated**; vacc marked on card; vacc marked on card; not vaccinated; vacc marked on card; not vaccinated; vacc marked on card; vacc marked on card; vacc marked on card; vacc marked on card; vacc marked on card; vacc marked on card; vacc marked on card; not vaccinated; vacc marked on card; vacc marked on card; vacc marked on card; vacc marked on card; vacc marked on card; not vaccinated; vacc marked on card; not vaccinated; pas vaccine; vacc marked on card; vacc marked on card; vacc marked on card; not vaccinated; vacc marked on card; vacc marked on card; non vaccine; vacc marked on card; vacc marked on card; vacc marked on card; vacc marked on card; vacc marked on card; vacc marked on card; vacc marked on card; pas carnet; no vacunado; not vaccinated; vacc marked on card; vacc marked on card; vacc marked on card; vacc marked on card; no vacunado; vacc marked on card; vacc

marked on card; vacc marked on card; vacc marked on card; **vacc marked on card**; not vaccinated; vacc marked on card; not vaccinated; vacc marked on card; vacc marked on card; vacc marked on card; vacc marked on card; vacc marked on card; not vaccinated; not vaccinated; vacc marked on card; vacc marked on card; vacc marked on card; vacc marked on card; vacc marked on card; vacc marked on card; vacc marked on card; not vaccinated; vacc marked on card; vacc marked on card; vacc marked on card; vacc marked on card; vacc marked on card; vacc marked on card; vacc marked on card; vacc marked on card; vacc marked on card; vacc marked on card; vacc marked on card; vacc marked on card; vacc marked on card; not vaccinated; vacc marked on card; vacc marked on card; vacc marked on card; vacc marked on card; vacc marked on card; not vaccinated; vacc marked on card; vacc marked on card; vacc marked on card; vacc marked on card; not vaccinated; vacc marked on card; vacc marked on card; vacc marked on card; vacc marked on card; not vaccinated; vacc marked on card; vacc marked on card; vacc marked on card; vacc marked on card; not vaccinated; vacc marked on card; vacc marked on card; vacc marked on card; vacc marked on card; vacc marked on card; vacc marked on card; vacc marked on card; vacc marked on card; vacc marked on card; vacc marked on card; vacc marked on card; vacc marked on card; pas vaccine; vacc marked on card; vacc marked on card; vacc marked on card; pas vaccinnee; not vaccinated; vacc marked on card; vacc marked on card; not vaccinated; not vaccinated; nao vacinada; not vaccinated; not vaccinated; vacc marked on card; not vaccinated; vacc marked on card; vacc marked on card; pas vaccine; vacc marked on card; vacc marked on card; not vaccinated; vacc marked on card; vacc marked on card; vacc marked on card; vacc marked on card; vacc marked on card; vacc marked on card; vacc marked on card; not vaccinated; not vaccinated; vacc marked on card; vacc marked on card; vacc marked on card; vacc marked on card; not vaccinated; vacc marked on card; vacc marked on card; vaccination marked on card; not vaccinated; vacc marked on card; vacc marked on card; vacc marked on card; vacc marked on card; vacc marked on card; vacc marked on card; vaccination marked on card; vaccination marked on card; vacc marked on card; vaccination marked on card; non vaccine; vaccination marked on card; vacc marked on card; vacc marked on card; vaccination marked on card; vaccination marked on card; vacc marked on card; vaccination marked on card; vaccination marked on card; vacc marked on card; vacc marked on card; not vaccinated; vacc marked on card; vacc marked on card; vacc marked on card; vaccination marked on card; vacc marked on card; vaccination marked on card; vaccination marked on card; vaccination marked on card; vacc marked on card; vaccination marked on card; vacc marked on card; vacc marked on card; vacc marked on card; vaccination marked on card; vaccination marked on card;

## Slide 2

Swiss TPH

**FIRST ATTEMPTS**

$$\text{Coverage} = \frac{\text{Number of children age A vaccinated with V}}{\text{Number of children age A targetted by V}}$$

## Slide 3



## Slide 4

## Slide 1

**INCONVENIENCES**

- Errors

- Code with hundreds of lines:
  - surveys x variables + surveys x variables x values

> If survey = Afghanistan 2000 AND OldVar = h7
>
> rename h7 to DPT3
> recode DTP3 1=2, 2=3, 3=0
> ...

- At survey 68: complains from collaborators

- At survey 174: nonsense, who am I? what is life?

- At survey 362: addictions

## Slide 2

**THE SOLUTION: HARMONISATION**

## Slide 3

```
use Datasetsvars.dta

for each survey S in Datasetsvars.dta
    for each variable of interest V for survey S in Datasetsvars.dta

        rename OldVar to NewVar
        recode OldValues to NewValues

    go to next variable of interest
go to next survey
```

## Slide 4

**Folder structure**

- Setup
  - 0. Start.do: starting `Project' folder
  - 1. Setup.do: Stata environment
  - **2. Main.do**: main loop (z)

- Harmonisation
  - 0.Start.do
  - Code
  - Logs
  - Outputs
    - Graphs
    - Tables
    - Text (.csv)

- Data
  - DatasetsVars.dta
  - Exceptions.dta
  - Thesaurus.dta
  - Surveys
  - Other

6

**Slide 1:**

Initial procedures
• Folders...
• Open logs
• ...

**MAIN LOOP**

Do Include / exclude datasets and variables (DatasetsVars.dta)

• Folder: code → Code C = 1_abc.do, 2_cde.do..., n_xyz.do

Do Code C once

• Folder: datasets → Dataset D = 1_abc.dta, 2_cde.dta..., n_xyz.dta

Do Code C with Dataset D

Final procedures
• Export Outputs / tables -> Outputs / txt (xls)
• Close logs
• ...

---

**Slide 2:**

Harmonisation (overview)
• Any available dataset
• Import (and repair country and year) datasets
• Update DatasetsVars

**1. Variables**

• Search and identify variables
• Decisions on searches outputs
Produce output for manual checks

**2. Value labels**
• Search and identify values
• Decisions on searches outputs
Produce output for manual checks

**3. Rename and recode**

---

**Slide 3:**

**1. Harmonisation of variables**

**STRATEGY B**

Old Labels    vs    New key terms

**Normalisation Thesaurus**

| Old labels | Normalised | Thesaurus | CHK [Day DTP 3] | Outcome |
|---|---|---|---|---|
| day of dpt  III | DAY DTP 3 | DAY DTP 3 | 3 = 3 | Best match |
| dia vaccinação dtp 3 | DIA VACCINACAO DTP 3 | DAY vaccinacao DTP 3 | 4 > 3 | Match |
| dtp 3 | DTP 3 | DTP 3 | 2 < 3 | OUT |

**STRATEGY C**

Old Variables used in other surveys for New Variable, which exist in current dataset: abc, def, ghi, ikl...

---

**Slide 4:**

```
local nT = 0
foreach cTerm of local cVarNewKTerms { // for 'dtp', '3', 'year' in 'year dtp3'.
        local nT = `nT' + 1
        quietly ds
        foreach cVar in `r(varlist)' { // dtp1 dtp2 dtp3...
                local cLabel : variable label `cVar' // dtp 3 day vaccinaçao
                local nFoundInThisVar = 0
                foreach cSyn of local Synos`nT' { // For each synonymous of
                                                  cTerm.
                        local nFoundInThisVar = cond(strpos(" `cLabel' ", "
                          `cSyn' ") > 0 | `nFoundInThisVar' == 1, 1, 0)
                }
                if `nFoundInThisVar' == 0 drop `cVar' // Drops var if none of the
                synonymous of the term is found.
        }
}
ds
```

I drop existing 'old' variables which do not have term 1 (DTP); from the remaining, I drop those which do not have term 2 (DAY)...

Slide 1 (top-left):

| Search | Type | SData | SVarNew | SLabelNew | SKeyTerms | Clear | SVarOld | SLabelOld |
|---|---|---|---|---|---|---|---|---|
| B* | mics | afghanistan-mic | BCGbm | BCG month | bcg month | 1 | im3bm | month bcg immunization |
| B* | mics | centralafricanre | BCGby | BCG year | bcg year | 0 | im2y | annee vaccination bcg |
| B* | mics | cuba-mics_4-ch | BCGby | BCG year | bcg year | 0 | im3by | ano inmunizacion bcg al nacer |
| B* | dhs | am-dhs_61-ch | BCGby | BCG year | bcg year | 0 | h2y | bcg year |
| B- | dhs | gy-dhs_5i-ch | BCGby | BCG year | bcg year | 0 | h2y | bcg year combined vacc cards |
| B- | dhs | gy-dhs_5i-ch | BCGby | BCG year | bcg year | 0 | s2bcgy | bcg year health facility vacc card |
| B- | dhs | gy-dhs_5i-ch | BCGby | BCG year | bcg year | 0 | s1bcgy | bcg year home vacc card |
| B* | mics | afghanistan-mic | BCGby | BCG year | bcg year | 0 | im3by | year bcg immunization |
| C70 | mics | mozambique-m | BCGd | BCG day | bcg day | 1 | im2d | bcg |

NO DUPLICATES: N = 2,898

---

Slide 2 (top-right):

| Search | Type | SData | SVarNew | SLabelNew | SKeyTerms | Clear | SVarOld | SLabelOld | Case | Rationale |
|---|---|---|---|---|---|---|---|---|---|---|
| B- | mics | mozambique-mics_3-ch | BCG | BCG | bcg | 1 | im2d | bcg | A.1 | bcg(); im2d(bcg); im2m(bcg); im2y(bcg); incomplete label, Var name consistent with other surveys |
| B- | mics | mozambique-mics_3-ch | BCG | BCG | bcg | 1 | im2m | bcg | A.1 | bcg(); im2d(bcg); im2m(bcg); im2y(bcg); incomplete label, Var name consistent with other surveys |
| B- | mics | mozambique-mics_3-ch | BCG | BCG | bcg | 1 | im2y | bcg | A.1 | bcg(); im2d(bcg); im2m(bcg); im2y(bcg); incomplete label, Var name consistent with other surveys |

ALL RECORDS: N = 15,153

---

Slide 3 (bottom-left):

---

Slide 4 (bottom-right):

## 2. Harmonisation of values labels

Normalisation   Thesaurus

| | | 0 = abc |
| 0 = old value label for 0 | | 1 = def |
| 1 = old value label for 1 | vs | ... |
| ... | | N = xyz |
| N = old value label for N | | -6 = no label |
| | | -7 = inconsistent |
| | | -8 = don't know |
| | | -9 = missing |

| Old value labels | New value label | Outcome |
|---|---|---|
| 2 = history | 4 = recall | Match: 2 -> 4 |
| - | 1 = yes | No match; empty 1 = yes |
| 44 = in health centre | - | No match: manual review |

8

```
// MATCHING VALUE LABELS. For each Value label (old), check whether there is an equivalent in any of the Value label (new).
foreach cValLblOld_i of local lValLblOld { // Test each value label old... e.g. 'caretaker recall'
        // Each one of the Value label (new) to be tested against the old one above.
        foreach cValLblNew_i of local lValLblNew {
                        if !missing(mValLblNew[rownumb(mValLblNew,"`cValLblNew_i'"),2]) continue // Has already been
                                                                                              found.
                        local cValLblNew_i2 = subinstr("`r(xR)'", "_", " ", .)

                        // E.g. 'mother' against 'caregiver', then 'recall'; and 'recall' against 'caregiver', then 'recall').
                        fNormTxt "`cValLblNew_i'"

                        local cValLblNew_i2 = subinstr("`r(xR)'", "_", " ", .)

                        local nTermsFound = 0
                        foreach cValLblNew_i_t of local cValLblNew_i2 { // All terms of value label new i have to be in an old
                                                                        one; e.g. 'mother'
                                        do "`cDoSynonymous'" "`cValLblNew_i_t'"
                                        local cValLblNew_i_ts = "`r(xR)'"
                                        local cValLblOld_i2 = subinstr("`cValLblOld_i'", "_", " ", .)

                                        local nIsTermFound = 0
                                        foreach cValLblOld_i_t of local cValLblOld_i2 { // e.g. 'mother' 'recall'

                                                        if strpos(" `cValLblNew_i_ts' ", " `cValLblOld_i_t' ") > 0 {
                                                                        local nIsTermFound = `nIsTermFound' + 1
                                                                        continue, break
                                                        }
                                        }
                                        local nTermsFound = `nTermsFound' + `nIsTermFound'
                        }
        }
} // If this point is reached, there are value labels to handle (cValLbl is not missing)
```

---

| NewVal | NewLab | OldVal | OldLab | Dataset | OldVar | NewVar |
|---|---|---|---|---|---|---|
| 1 | Female | 2 | feminino | mozambique-mics_3-ch | hl4 | Sex |
| 1 | Female | 2 | female | afghanistan-mics_4-ch | hl4 | Sex |
| 1 | Female | 2 | femenino | cuba-mics_4-ch | hl4 | Sex |
| 1 | Female | 2 | feminin | centralafricanrepublic-mics_3-ch | hl4 | Sex |

NO DUPLICATES: N = 1,101

---

| | Dataset | OldVar | OldVal | OldLab | NewVar | NewVal | NewLab |
|---|---|---|---|---|---|---|---|
| 2158 | democraticrepublicof | 1m3d3d | 0 | pas reçue | DTP3d | 0 | No |
| 2159 | ke-dhs_52-ch | h7d | 98 | dk | DTP3d | -7 | DK |
| 2160 | sl-dhs_51-ch | h7d | 97 | inconsistent | DTP3d | -8 | Inconsistent |
| 2161 | djibouti-mics_4-ch | 1m4cd | 44 | marquee sur la carte | DTP3d | 44 | Mark |
| 2162 | democraticrepublicof | 1m3d3d | 66 | reportee par la mere | DTP3d | 66 | Recall |
| 2163 | togo-mics_4-ch | 1m3d3d | 44 | marquee sur le carne | DTP3d | 44 | Mark |
| 2164 | sn-dhs_60-ch | h7d | 98 | dontknow_dk | DTP3d | -7 | DK |
| 2165 | nigeria-mics_4-ch | 1m2d3d | 0 | not g4en | DTP3d | 0 | No |
| 2166 | bd-dhs_61-ch | h7d | 97 | inconsistent | DTP3d | -8 | Inconsistent |
| 2167 | zimbabwe-mics_3-ch | 1m4cd | 0 | not g4en | DTP3d | 0 | No |
| 2168 | gy-dhs_51-ch | h7d | 0 | not g4en | DTP3d | 0 | No |
| 2169 | nigeria-mics_3-ch | 1m4cd | 44 | marked on card _ nod | DTP3d | 44 | Mark |
| 2170 | suriname-mics_3-ch | 1m4cd | 66 | reported by mother | DTP3d | 66 | Recall |
| 2171 | afghanistan-mics_4-c | 1m3d3d | 66 | reported by mother | DTP3d | 66 | Recall |
| 2172 | zimbabwe-mics_3-ch | h7d | 98 | dk | DTP3d | -7 | DK |
| 2173 | suriname-mics_4-ch | 1m3d3d | 44 | marked on card _ nod | DTP3d | 44 | Mark |
| 2174 | suriname-mics_4-ch | 1m3d3d | 99 | missing | DTP3d | -9 | Missing |
| 2175 | vanuatu-mics_3-ch | 1m4cd | 66 | reported by mother | DTP3d | 66 | Recall |
| 2176 | vietnam-mics_4-ch | 1m3d3d | 66 | reported by mother | DTP3d | 66 | Recall |
| 2177 | am-dhs_61-ch | h7d | 98 | dontknow_dk | DTP3d | -7 | DK |
| 2178 | ng-dhs_52-ch | h7d | 98 | dk | DTP3d | -7 | DK |
| 2179 | kh-dhs_51-ch | h7d | 98 | dontknow_dk | DTP3d | -7 | DK |
| 2180 | cm-dhs_60-ch | h7d | 97 | inconsistent | DTP3d | -8 | Inconsistent |
| 2181 | bu-dhs_61-ch | h7d | 97 | inconsistent | DTP3d | -8 | Inconsistent |
| 2182 | vanuatu-mics_4-ch | 1m4cd | 99 | missing | DTP3d | -9 | Missing |
| 2183 | tl-dhs_61-ch | h7d | 98 | dk | DTP3d | -7 | DK |
| 2184 | democraticrepublicof | 1m3d3d | 44 | marque sur le carne | DTP3d | 44 | Mark |
| 2185 | zimbabwe-mics_3-ch | 1m4cd | 44 | marked on card _ nod | DTP3d | 44 | Mark |

ALL RECORDS: N = 7,891

---

**3. Rename and recode**

When all cleared (= 1)

**RENAME**

**RECODE**

Value labels dataset

**WAY FORWARD**

1. Database of surveys
2. Database of coverage from survey reports and web-based estimates
3. Harmonisation platform integrated into WUEIC engine
4. Use variables distribution
5. Availability to third parties / cross-platforms

**World Health Organization** | **Global Immunization News (GIN)** | **January 2014**

**WHO and UNICEF** working group on monitoring national immunization coverage development session

Xavier BOSCH-CAPBLANCH, Swiss Tropical and Public Health Institute

Location:      Basel, Switzerland

Date:          12-14 November 2013

Participants:  United Nations Children's' Fund, World Health Organization, Swiss Tropical and Public Health Institute.

---

**Thanks for your infinite patience!**

42

11