# Combining microdata:
# the art of matching and joining

*Egon Gerards*

# Matching: personal identifiers

## Datasets contain personal identifiers

### 1. Direct unique personal identifiers:

- &Ł̶Ŧ̶ǔˇꝺ#ꞩˇẞóŁˊˇ#ꞁ°æ ꞏˇẞ
- ’°–ˇ#ꞁ–#ꞏꞆẞ–ª
- ꞁ°æˇ
- $ ꞏꞏẞˇꞁꞁ
- 6ˇÖ
- 3Ŧꞁ–°#ˊꞁ ꞏˇ
- (–ˮ

### 2. Indirect/other personal identifiers

# Matching: personal identifiers

## 1. Direct unique personal identifiers:
 --> replaced by SN linkage keys or removed/adjusted
- PIN: person identification number
- AIN: address identification number

## 2. Indirect/other personal identifiers:
- ẞˇæ Łóˇ ̦#
- ˇı ̋ƁÊø–ˇ ̦

# **Matching: population register**

**The Dutch population register is the key factor in matching!**

- **Cumulative data from 1995 onwards**
- **Every month updated**

| Personal identifiers | Unknown (1995 - 2015) |
|---|---|
| Citizen Service Number (CS) – BSN | <0,0025 % |
| Date of birth | PR contains no unknown, but 1st january and 1st july are often used when exact date is unknown. |
| Sex | 0% |
| Postal code | <0,09% |
| Number of the house | <0,09% |

# Matching: strategy

## Matching strategy
1. Optimizing number of matches
2. Minimizing number of mismatches and missed matches

## Matching keys
1. Registers (i.e.): Citizen Service Number (unique), addresses
2. Surveys (i.e.): Sex, date of birth, address (postal code and house number)

| Matching | Return |
|---|---|
| Citizen Service Number | >95% |
| Postal code6+sex+date of birth | >90% |
| Postal code4+sex+date of birth | >85% |

# Joining: the SSD

**Data from**

**different sources**

**Then:**

**Now:**

# The SSD

## structures and links data

SSD

# The strength of the SSD

**Long period**

**Combine**

**SSD**

**Specific**

**High quality**

# www.CBS.nl

@Statistiekcbs

E.Gerards@cbs.nl

**Egon Gerards**
Head of SSD department

# Privacy

**PRIVACY** audit proof

Checked for disclosure risks

Anonymised data

Limited access

Only output is disseminated, never any data

How does the SSD work?