# Towards a new integrated uniform production system for business statistics at Statistics Netherlands

## Quality indicators to guide top-down analysis

Anita Vaasen-Otten, Frank Aelen, Sander Scholtus and Wilco de Jong

3-10-2022

# EBN 2.x: Renewal program for business statistics at SN

**Main goals**: towards more efficient, more flexible and futureproof business statistics

**E**fficient, goal-oriented co-creation

**B**alance on innovation, re-use and implementation

**N**ew process, way of working and matching tooling

**2.x**: considered small steps

# Innovation, co-creation and implementations

**Innovation**

- 8 principles
- Applied methodology
- The proof of the pudding... is in a Proof of Concept!

**Co-creation with Business and with Agile teams**

- Re-use of best practices
- As much standardized coherent models & common tools as possible!

**Implementations**

- Small steps instead of a big bang
- Visible results
- Immediate feedback contributes to continuous improvements

# Principles of the new production system*

## Towards real-time processing

1. We process our input automatically and immediately up to provisional output;
2. We measure quality automatically and thus direct the manual work;

## Towards more coherence

3. We make our data consistent as early as possible;
4. We share all our data, right from the start;

## Towards more standardization and re-use of best practices

5. We centrally manage all our (population) frames, which are the basis of our statistics;
6. We have fully standardized our processes, methods, data and IT;
7. Our processes, methods, data and IT are modular;

## Continuous improvement

8. We resolve manual corrections the following iteration in the standard process.

# Quality indicators to guide top-down analysis*

- Focus on score functions that identify potential influential errors in the data;

- Used to prioritize records for manual editing and to quickly zoom in on the part of the record where there may be a problem;

- Can also be used to indicate the expected quality of an aggregate.

**\*Automatic data editing is described in a companion paper**

# Local score for level variable

$$s_{i,j} = \frac{v_i * |y_{i,j} - \tilde{y}_{i,j}|}{|Y_j|}, \quad (j = 1, \ldots, J) \qquad (1)$$

with $v_i$ the sample weight of unit $i$, $y_{i,j}$ the observed value of variable $j$, $\tilde{y}_{i,j}$ a reference value for variable $j$ and $Y_j$ an estimate for the aggregate total for variable $j$

- Relative influence of possible error on the output
- Reference value e.g. t-1, other source, related variable, …

# Additional local scores

- Structure variables;

- Consistency across statistics;

- Non-linear indicators with two or more target variables simultaneously (e.g. production-use ratios for national accounts).

More details available in paper

# Global scores

Global score $s_i$ per unit that is compiled from the underlying local scores:

$$s_i(\alpha) = \left\{ \frac{\sum_{j=1}^{J}\left(w_j * \frac{s_{i,j}}{M_j}\right)^{\alpha}}{\sum_{j=1}^{J} w_j^{\alpha}} \right\}^{1/\alpha} \quad (2)$$

Adjustable weights $w_j$ can be used to indicate that certain target variables, such as totals of revenue or costs, are more important than others (such as their details).

$M_j$ is a measure for the 'maximum acceptable' relative influence per unit of a possible error in the target variable(s) of local scores $s_{i,j}$ on the aggregate in the denominator.

# Aggregate scores

Summary measure based on the scores for all units that contribute to a particular output aggregate $A$. Only scores are counted above a certain threshold $\tau_A$.

This can be done for local scores:

$$S_{j,A} = \sum_{i \in A} \frac{s_{i,j}}{M_j} * I \left\{ \frac{s_{i,j}}{M_j} \geq \tau_A \right\}, \quad (j = 1, \dots, J), \qquad (3)$$

and for global (or other composite) scores:

$$S_A(\alpha) = \sum_{i \in A} s_i(\alpha) * I\{s_i(\alpha) \geq \tau_A\}, \qquad (4)$$

where $I\{.\} = 1$ if the argument is true and else $I\{.\} = 0$.

# Example: single statistic – publication aggregates



Prioritize aggregates – based on aggregate score

# Example: single statistic – overview of units

Aggregatie
KerncelCode

Drempelwaarde
0

Kerncellen | Alle eenheden

Score sorteervariabele
OPBRENG000000    Sorteer

Kerncel commentaar    Opmerking opslaan

Hoofdoverzicht | Procesinfo

Show 25 entries      Search:

| | | | Hoofdvariabelen | | | Pop. Dynamiek | Scores | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| BE_ID | BRANCHE | KerncelCode | OPBRENG000000 | INKWRDE100000 | BEDRLST310000 | • SOM_Bruto_invloed | GLOBAL_SCORE | Personeelsaantallen | Omzet | Overige_Opbrengsten | Inkoopwaarde | Personeelskosten | Bedrijfslasten | Onttrekking_vo |
| | A | All | All | All | All | All | All | All | All | All | All | All | All | All |
| 5..3 | Industrie | 33120 | 249334 | 243005 | 247842 | | 3.16 | 0.01 | 5.21 | 0.01 | 3.56 | 2.02 | 4.74 | |
| 1..9 | Industrie | 33120 | 105747 | 13363 | 112821 | | 2.41 | 0.11 | 3.16 | 0.05 | 4.58 | 0.28 | 10.41 | |
| 7..6 | Industrie | 33120 | 186846 | 176044 | 188738 | | 2.22 | 0 | 7.63 | 0.62 | 10.09 | 0.05 | 0.75 | |
| 2..1 | Industrie | 33120 | 98557 | 5108 | 40611 | | 1.17 | 0 | 0.79 | 0.63 | 0.8 | 1.15 | 1.51 | |
| 6..6 | Industrie | 33120 | 53312 | 21087 | 59173 | | 0.68 | 0.01 | 0.47 | 0 | 0.52 | 0.03 | 0.6 | |
| 1..8 | Industrie | 33120 | 225522 | 100851 | 209809 | | 0.61 | 0.17 | 0.19 | 0.17 | 2.13 | 0.36 | 0.89 | |
| 7..1 | Industrie | 33120 | 52760 | 28676 | 44906 | | 0.42 | 0 | 0.21 | 0.01 | 2.87 | 0.03 | 0.08 | |
| 6..6 | Industrie | 33120 | 290354 | 191240 | 292005 | | 0.41 | 0.18 | 0.25 | 0 | 0.56 | 0.23 | 1.15 | |
| 6..9 | Industrie | 33120 | 311731 | 155441 | 318069 | | 0.34 | 0.28 | 0.54 | 0.61 | 0.12 | 0.21 | 0.53 | |
| 7..7 | Industrie | 33120 | 8008 | 335 | 6507 | | 0.33 | 0.07 | 0.44 | 0 | 0.34 | 0.08 | 0.46 | |
| 7..3 | Industrie | 33120 | 17225 | 895 | 3056 | | 0.3 | 0.03 | 0.33 | 0 | 0.62 | 0.25 | 0.18 | |
| 6..7 | Industrie | 33120 | 52609 | 41499 | 48531 | | 0.27 | 0.05 | 0.65 | 0.07 | 0.66 | 0.13 | 0.04 | |

Showing 1 to 25 of 165 entries      Previous 1 2 3 4 5 6 7 Next

Prioritize units within aggregate – based on global score

# Example: single statistic – individual unit



Prioritize variables within unit – based on local scores

# Example: across statistics – publication aggregates

**Regkol confrontatie cluster**

overview    BE_ID

**Variabele**
Totale_Omzet_2 ▼

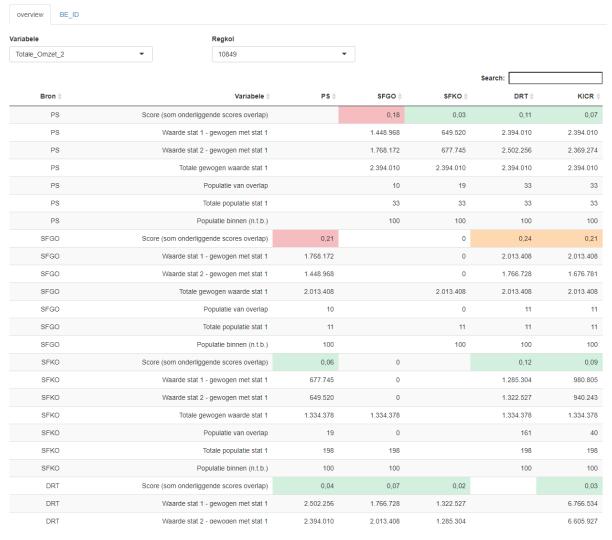**Regkol**
10849 ▼

Search: [_____]

| Bron | Variabele | PS | SFGO | SFKO | DRT | KICR |
|------|-----------|----|------|------|-----|------|
| PS | Score (som onderliggende scores overlap) | | 0,18 | 0,03 | 0,11 | 0,07 |
| PS | Waarde stat 1 - gewogen met stat 1 | | 1.448.968 | 649.520 | 2.394.010 | 2.394.010 |
| PS | Waarde stat 2 - gewogen met stat 1 | | 1.768.172 | 677.745 | 2.502.256 | 2.369.274 |
| PS | Totale gewogen waarde stat 1 | | 2.394.010 | 2.394.010 | 2.394.010 | 2.394.010 |
| PS | Populatie van overlap | | 10 | 19 | 33 | 33 |
| PS | Totale populatie stat 1 | | 33 | 33 | 33 | 33 |
| PS | Populatie binnen (n.t.b.) | | 100 | 100 | 100 | 100 |
| SFGO | Score (som onderliggende scores overlap) | 0,21 | | 0 | 0,24 | 0,21 |
| SFGO | Waarde stat 1 - gewogen met stat 1 | 1.768.172 | | 0 | 2.013.408 | 2.013.408 |
| SFGO | Waarde stat 2 - gewogen met stat 1 | 1.448.968 | | 0 | 1.766.728 | 1.676.781 |
| SFGO | Totale gewogen waarde stat 1 | 2.013.408 | | 2.013.408 | 2.013.408 | 2.013.408 |
| SFGO | Populatie van overlap | 10 | | 0 | 11 | 11 |
| SFGO | Totale populatie stat 1 | 11 | | 11 | 11 | 11 |
| SFGO | Populatie binnen (n.t.b.) | 100 | | 100 | 100 | 100 |
| SFKO | Score (som onderliggende scores overlap) | 0,06 | 0 | | 0,12 | 0,09 |
| SFKO | Waarde stat 1 - gewogen met stat 1 | 677.745 | 0 | | 1.285.304 | 980.805 |
| SFKO | Waarde stat 2 - gewogen met stat 1 | 649.520 | 0 | | 1.322.527 | 940.243 |
| SFKO | Totale gewogen waarde stat 1 | 1.334.378 | 1.334.378 | | 1.334.378 | 1.334.378 |
| SFKO | Populatie van overlap | 19 | 0 | | 161 | 40 |
| SFKO | Totale populatie stat 1 | 198 | 198 | | 198 | 198 |
| SFKO | Populatie binnen (n.t.b.) | 100 | 100 | | 100 | 100 |
| DRT | Score (som onderliggende scores overlap) | 0,04 | 0,07 | 0,02 | | 0,03 |
| DRT | Waarde stat 1 - gewogen met stat 1 | 2.502.256 | 1.766.728 | 1.322.527 | | 6.766.534 |
| DRT | Waarde stat 2 - gewogen met stat 1 | 2.394.010 | 2.013.408 | 1.285.304 | | 6.605.927 |

Prioritize statistics with large inconsistencies between them, for a certain aggregate and mutual variable

13

# Example: across statistics – overview of units



Prioritize units within aggregate

# Implementations

- Past two years we tested and refined ideas in POCs;

- Implemented in generalized R-modules – web service;

- Scores can be tailored to various statistics by means of rules;

- For a limited number of individual statistics, the scores have been implemented and are already being applied in practice;

- Later this year: pilot regarding the top-down analysis of inconsistencies between statistics ➔ gain experience with new roles that are necessary for this new way of working.

# Concluding remarks

- Experiences to date show that the new scores allow analysts to work in a more targeted way than before;

- In the near future we will continue the stepwise developments and implementations and working in an agile manner, we will keep learning from each further step.