# Sampling and estimation techniques for the implementation of new classification systems: the change-over from NACE Rev. 1.1 to NACE Rev. 2 in business surveys

Jan van den Brakel

Department of Statistical Methods, Statistics Netherlands

This paper describes some of the methodological problems encountered with the change-over from the NACE Rev. 1.1 to the NACE Rev. 2 in business statistics. Different sampling and estimation strategies are proposed to produce reliable figures for the domains under both classifications simultaneously. Furthermore several methods are described that can be used to reconstruct time series for the domains under the NACE Rev. 2.

**Keywords:** Backcasting, discontinuities, design-based estimators, small area estimation, linking

## 1 Introduction

### 1.1 Discontinuities in series of repeatedly conducted sample surveys

Sample surveys conducted by national statistical institutes are generally conducted repeatedly in time with the purpose of constructing time series that describe the evolution of population parameters of interest. An important quality aspect of these surveys is comparability of the outcomes over time. To maintain consistent time series, the underlying survey process is generally kept unchanged as long as possible. It remains, however, inevitable to change or redesign a survey process from time to time. A major drawback of such redesigns is that it often has systematic effects on the outcomes of the survey, leading to discontinuities in the series. An important aspect of a survey redesign is to minimize this inconvenience for data users. This can be accomplished by quantifying the effect of the redesign on the outcomes of the main parameters. To maintain consistent time series, one might consider to correct the series observed in the past with the observed effects of the redesign. This is sometimes referred to as backcasting.

Van den Brakel et al. (2008) discuss different statistical methods to deal with discontinuities due to survey redesigns. The methods required to quantify a discontinuity depend on the phase of the survey process that is changed. In cases where the underlying sample data remain the same, the differences can be investigated by recalculation. An example is the introduction of a new economic activity classification system in business surveys. When, however, data collection procedures are affected the data are not consistent. In these cases the effect of the change can be quantified by conducting a field experiment where the regular and new survey designs are run concurrently, see e.g. Van den Brakel (2008) for examples and details. Such a parallel run is not always tenable due to budget constraints. In such cases a time series modelling approach can be considered as an alternative. A so-called intervention analysis is described in detail by Van den Brakel and Roels (2010) using state-space models.

This paper describes the statistical methods that can be applied to assess the effect of a new economic activity classification system in business surveys. These methods are also applicable to the implementation of revised versions of the International Standard Classification of Occupations used by the International Labour Organization, the Standard Occupational Classification, used by Federal statistical agencies in the US and Canada, or the International Standard Classification of Education used by the United Nations Educational Scientific and Cultural Organisation.

### 1.2 Introduction of a new economic classification system in business surveys

In all European Union countries the classification of economic activities that is used in the Business Surveys was from 1993 through 2008 based on NACE Rev 1.1. NACE stands for the French abbreviation for European Classification of Economic Activities.[1] Since the economic structure gradually changed, a new classification system, called NACE Rev. 2, was adopted in 2006 by Eurostat, the European national statistical institutes, European trade and business associations, the European Central Bank and United Nations Statistical Division. See Eurostat (2008) for a detailed description of the NACE Rev. 2. This classification system is introduced in the Short Term Statistics (STS) since 2009 and the Structural Business Surveys (SBS) since 2010.

---

Contact information: Jan van den Brakel, Department of Statistical Methods, Statistics Netherlands, P.O. Box 4481, 6401 CZ Heerlen, The Netherlands, e-mail: jbrl@cbs.nl

[1] NACE is derived from the French title: "Nomenclature générale des Activités économiques dans les Communautés Européennes"

*Table 1:* Number of occurrences of different types of transitions

| Transition | Number of occurrences |
|------------|-----------------------|
| 1-to-1     | 196 classes           |
| 1-to-n     | 17 classes            |
| m-to-1     | 87 cases              |
| m-to-n     | 214 cases             |

In a descending order of aggregation, the following levels are distinguished under the NACE Rev. 1.1: sections (one character alphabetic code), subsections (two character alphabetic code), divisions (two digit code), groups (three digit code) and classes (four digit code). Under the NACE Rev. 2, the level of the subsections is dropped. In general terms, the NACE Rev. 2 resulted in a more detailed classification of the activities in Services and less detail in the Industrial activities, reflecting the general development of the economic structure in Europe.

Annex A contains two tables, which summaries the effect of the change-over from the NACE Rev. 1.1 to the NACE Rev. 2. Table A.1 describes the change-over of the sections. This table indicates which sections under NACE Rev. 1.1 are grouped into one section under NACE Rev. 2 and which sections under NACE Rev. 1.1 are divided into two or more sections under NACE Rev. 2. Table A.2 provides further details about the number of subsections, divisions, groups and classes that are distinguished within each section of the NACE Rev. 1.1 and the NACE Rev. 2.

For many classes there are no changes, all business units belonging to a specific class under the NACE Rev. 1.1 transfer to a corresponding class under the NACE Rev. 2 and no other business units join this new class. These are the so-called 1-to-1 transitions. In several cases, however, business units transfer to different classes under the new classification. As a result, business units that are classified to the same class under the NACE Rev 1.1 can transfer to two or more classes under the NACE Rev. 2. These are called the 1-to-$n$ transitions. It is also possible that business units that are classified to different classes under the NACE Rev. 1.1 transfers to the same class under the NACE Rev. 2; the so-called $m$-to-1 transition. The most complex situation is the $m$-to-$n$ transitions, where business units that are classified to $m$ classes under the NACE Rev. 1.1 transfers to $n$ classes under the NACE Rev. 2. To illustrate the importance of the problem, an overview of the number of occurrences of the different types of transitions is given in Table 1.

The implementation of a new classification system results in several methodological challenges. The change-over from NACE Rev. 1.1 to NACE Rev. 2 in the business surveys starts with adding a NACE Rev. 2 code to all units in the Business Register (BR), see Eurostat (2006b) for details. To facilitate a smooth transition from the old to the new classification system, it is recommended to publish figures for a period of one or two years under both classifications for both the SBS and STS, Eurostat (2006a, 2006c). During this period, all units in the BR are preferably double coded under the NACE Rev 1.1 and NACE Rev. 2. It may be necessary to adjust the sample design to produce reliable figures for the domains under both classifications. If it is decided that the sample design is not changed during this period of double reporting, then this will generally require at least an adjustment of the estimation procedure for the domains under the new classification. Sooner or later the sample design needs to be adjusted to this new classification system, since the business statistics will finally be based on the NACE Rev. 2 only. This requires an adjustment of the stratification, determination of the sample size and a reconsideration of the allocation of the business units over the strata.

The results obtained during this period of double coding of the BR and double reporting will be used to reconstruct time series for the NACE Rev. 2 domains starting from the year 2000. This is generally referred to as backcasting. Depending on the available information and resources, a combination of estimation techniques from the classical survey sampling approach and more synthetic adjustment and linking procedures can be applied to construct historical time series for domains under the NACE Rev. 2.

In many national statistical institutes business surveys are based on stratified simple random sampling. Generally the stratification variable is based on the crossing of size class based on employment and publication cells based on the NACE Rev. 1.1. This is for example the situation at Statistics Netherlands where the stratification variables are defined by size based on employment in 10 classes crossed with the primary publication cells (PPC's) under the NACE Rev. 1.1. The PPC's are the most detailed publication cells, which almost correspond one to one with the classes of the NACE Rev. 1.1 at the four digit level.

Taking this situation as a starting point, four approaches are distinguished to produce figures for the domains under both classifications simultaneously. The first three approaches are based on the design-based estimation procedures known from classical sampling theory for stratified simple random sampling using three different stratification schemes. The fourth approach is based on model-based estimation procedures, known from the realm of small area estimation. This will be input for different backcasting procedures.

The paper is organised as follows. Section 2 briefly reviews the procedure used to implement the NACE Rev. 2 in the Dutch BR. The design-based estimation procedures for three different stratification schemes are discussed in sections 3, 4, and 5. The model-based approach is described in section 6. In section 7 an overview of backcasting procedures is provided. The paper concludes with a discussion in section 8.

## 2 Double coding the BR: the Dutch situation

In 2007 and 2008, the NACE Rev. 2 is implemented in the Dutch BR with the purpose to maintain a double coded BR until 2010. The Dutch BR contains about 1,200,000 enterprises. Ninety percent of these enterprises are recoded au-

tomatically using information from the Commercial Register of the Chambers of Commerce. For the remaining ten percent, the NACE Rev. 2 classification is derived manually from the information available from the registration records of the Chambers of Commerce and from the available knowledge of subject matter specialists from Statistics Netherlands. For a small proportion of these enterprises, this information was insufficient to establish their classification under the NACE Rev. 2. For these cases questionnaires were sent to gather the required information. The classification of the BR according the NACE Rev. 2 is checked, and corrected if necessary, with available information from the PRODCOM ("Products of the European Community", the Eurostat system for the collection and dissemination of statistics on the production of manufactured goods) and the SBS. This classification information is further treated as being without error.

## 3 Stratifying to NACE Rev 1.1

The first approach to produce figures under the old and new domains is to draw a stratified simple random sample, with a stratification that is based on the classes of NACE Rev. 1.1 crossed with size class. This implies that the sampling design is kept unchanged and classical design-based estimators are applied for estimating domain parameters under both classifications.

Estimators for the PPC's under the old classification are based on an estimator for a population parameter for stratified simple random sampling since the domains exactly coincide with unions of strata. Due to the 1-to-$n$, $m$-to-1 or $m$-to-$n$ transitions, the domains under the new classification will not necessarily coincide with unions of the strata. Therefore, the estimators for the new classification should be based on a domain estimator that accounts for the possibility that the domains of the publication cells cut through the applied stratification scheme.

Let $\pi_i$ denote the inclusion probability for sampling unit $i$ and $\pi_{ij}$ the joint inclusion probability for the units $i$ and $j$. For stratified simple random sampling it follows that:

$$\pi_i = \frac{n_{g,k}}{N_{g,k}} \quad \text{if} \quad i \in U_{g,k}, \tag{1}$$

$$\pi_{ij} = \begin{cases} \frac{n_{g,k}(n_{g,k}-1)}{N_{g,k}(N_{g,k}-1)} & \text{if} \quad i,j \in U_{g,k} \\ \frac{n_{g,k}n_{g',k'}}{N_{g,k}N_{g',k'}} & \text{if} \quad i \in U_{g,k}, j \in U_{g',k'} \end{cases}, \tag{2}$$

where $U_{g,k}$ denotes the subpopulation or stratum defined by the crossing of size class $g$ and PPC $k$.

Most target parameters of STS's are defined as indices, e.g. Laspeyres indices. Therefore the growth rate is an appropriate variable to illustrate estimation procedures and sample size determination. Parameters under the NACE Rev. 1.1 and the NACE Rev. 2 classification are distinguished with subscripts $k$ and $l$ respectively.

The monthly or quarterly growth rate of the turnover for the $k$-th domain under the NACE Rev. 1.1, for example, is an important target parameter, which is defined as

$$Q_k^{(t)} = \frac{Y_k^{(t)}}{Y_k^{(t-1)}}. \tag{3}$$

In (3) $Y_k^{(t)}$ and $Y_k^{(t-1)}$ denote the total turnover in the $k$-th PPC under the NACE Rev. 1.1 for period $t$ en $t$-1. An estimator for (3) for the PPC's under the NACE Rev. 1.1 is given by

$$\hat{Q}_k^{(t)} = \frac{\hat{Y}_k^{(t)}}{\hat{Y}_k^{(t-1)}}, \tag{4}$$

with

$$\hat{Y}_k^{(t)} = \sum_{g=1}^{G} \frac{N_{g,k}}{n_{g,k}} \sum_{i=1}^{n_{g,k}} y_{i,g,k}^{(t)}. \tag{5}$$

Here $y_{i,g,k}^{(t)}$ denotes the turnover of business unit $i$ that belongs to size class $g$ and PPC $k$ at time period $t$, $N_{g,k}$ the total number of business units in the population of stratum $(g,k)$ and $n_{g,k}$ the sample size in stratum $(g,k)$. An expression for $\hat{Y}_k^{(t-1)}$ is defined analogously to (5) with $y_{i,g,k}^{(t)}$ replaced by $y_{i,g,k}^{(t-1)}$.

In many national statistical institutes, the precision of a direct estimator like (5), is improved by taking advantages of available auxiliary information through the generalised regression (GREG) estimator, see e.g. Särndal et al. (1992). Let $X_k^{(t)}$ denote the vector with the known population totals of the auxiliary information in the $k$-th PPC for period $t$. Then the GREG estimator for $Y_k^{(t)}$ is defined as

$$\hat{Y}_{k;\text{greg}}^{(t)} = \hat{Y}_k^{(t)} + \hat{b}_k^{(t)'}(X_k^{(t)} - \hat{X}_k^{(t)}), \tag{6}$$

with $\hat{X}_k^{(t)}$ a direct estimator for $X_k^{(t)}$ of the form (5), with $y_{i,g,k}^{(t)}$ replaced by a vector with the auxiliary information of the $i$-th business unit belonging to stratum $(g,k)$ for period $t$, say $x_{i,g,k}^{(t)}$. Furthermore, $\hat{b}_k^{(t)}$ denotes the regression coefficient of the regression function of $y_{i,g,k}^{(t)}$ on $x_{i,g,k}^{(t)}$. See formula (6.4.1) in Särndal et al. (1992) for an expression of $\hat{b}_k^{(t)}$. Regression estimator (6) can also be expressed as the weighted sum over the observations obtained in the sample:

$$\hat{Y}_{k;\text{greg}}^{(t)} = \sum_{g=1}^{G} \sum_{i=1}^{n_{g,k}} w_{i,g,k}^{(t)} y_{i,g,k}^{(t)}, \tag{7}$$

where $w_{i,g,k}^{(t)}$ are the so-called regression weights. These weights can be interpreted as the minimally adjusted design weights $d_{i,g,k}^{(t)} = N_{g,k}/n_{g,k}$, under a quadratic loss function, such that the requirement is fulfilled that the weighted auxiliary variables in the sample adds up to the known population totals. See Särndal et al. (1992), section 6.5 for an expression of the regression weights in (7) and Luery (1986), Alexander (1987) or Deville and Särndal (1992) for a more general treatment of the GREG estimator as a special case of the family of calibration estimators.

The notation in (6) suggests that the weighting scheme of the GREG estimator is also stratified according to NACE

Rev. 1.1 classification. This might be preferable, but it is not necessarily required. The weighting scheme might be defined on a larger aggregation level, for example to avoid unstable regression weights.

The ratio estimator can be derived as a special case from the GREG-estimator, Särndal et al. (1992), section 6.4, and is often used in business surveys, for example with value added tax as the auxiliary variable. If, for example, value added tax is used at the level of the PPC's under the NACE Rev. 1.1 classification in a ratio estimator, then (6) simplifies to:

$$\hat{Y}_{k;greg}^{(t)} = \frac{\hat{Y}_k^{(t)}}{\hat{X}_k^{(t)}} X_k^{(t)}, \tag{8}$$

where $\hat{X}_k^{(t)}$ is defined by (5), with $y_{i,g,k}^{(t)}$ replaced by the value added tax of the $i$-th business unit belonging to stratum $(g,k)$ for period $t$, say $x_{i,g,k}^{(t)}$.

An approximately design-unbiased estimator for the variance of (4) is given by (Cochran, 1977, Ch. 6):

$$V\hat{a}r(\hat{Q}_k^{(t)}) = \frac{1}{(\hat{Y}_k^{(t-1)})^2} \sum_{g=1}^{G} N_{g,k}(N_{g,k} - n_{g,k}) \frac{\hat{S}_{g,k}^{(t)2}}{n_{g,k}}, \tag{9}$$

with

$$\hat{S}_{g,k}^{(t)2} = \frac{1}{n_{g,k} - 1} \sum_{i=1}^{n_{g,k}} (\hat{z}_{i,g,k}^{(t)} - \hat{\bar{Z}}_{g,k}^{(t)})^2, \tag{10}$$

$$\hat{z}_{i,g,k}^{(t)} = y_{i,g,k}^{(t)} - \hat{Q}_k^{(t)} y_{i,g,k}^{(t-1)}, \tag{11}$$

and

$$\hat{\bar{Z}}_{g,k}^{(t)} = \frac{1}{n_{g,k}} \sum_{i=1}^{n_{g,k}} \hat{z}_{i,g,k}^{(t)}. \tag{12}$$

In the case of the GREG-estimator, the same formula's can be used for variance estimation. In (9) $\hat{Y}_k^{(t-1)}$ must be replaced by the GREG-estimator $\hat{Y}_{k;greg}^{(t-1)}$ and the residuals in (11) are replaced by:

$$\hat{z}_{i,g,k}^{(t)} = y_{i,g,k}^{(t)} - \hat{\boldsymbol{b}}_k^{(t)'} \boldsymbol{x}_{i,g,k}^{(t)} - \hat{Q}_{k;greg}^{(t)}(y_{i,g,k}^{(t-1)} - \hat{\boldsymbol{b}}_k^{(t-1)'} \boldsymbol{x}_{i,g,k}^{(t-1)}) \tag{13}$$

with $\hat{Q}_{k;greg}^{(t)} = \hat{Y}_{k;greg}^{(t)}/\hat{Y}_{k;greg}^{(t-1)}$. For the example with the ratio estimator where value added tax is used as auxiliary information, the residuals in (11) are defined as:

$$\hat{z}_{i,g,k}^{(t)} = y_{i,g,k}^{(t)} - \frac{\hat{Y}_k^{(t)}}{\hat{X}_k^{(t)}} x_{i,g,k}^{(t)} - \hat{Q}_{k;greg}^{(t)}(y_{i,g,k}^{(t-1)} - \frac{\hat{Y}_k^{(t-1)}}{\hat{X}_k^{(t-1)}} x_{i,g,k}^{(t-1)}). \tag{14}$$

Small stratum sample sizes result in unstable estimates for the stratum population variance $S_{g,k}^{(t)2}$. Stable estimates for $S_{g,k}^{(t)2}$ can be obtained by pooling the within-stratum variance for the strata with assumed equal population variances. Let

$\hat{S}_{g,k,(P)}^{(t)2}$ denote the pooled estimate for the population variance of the strata from size class $g = g_1, ..., g_a$ and PPC's $k = k_1, ..., k_b$. In the case of stratified simple random sampling the following ANOVA-type estimator can be used to pool the within-stratum variances:

$$\hat{S}_{g,k,(P)}^{(t)2} = \frac{1}{\sum_{g=g_1}^{g_a} \sum_{k=k_1}^{k_b} n_{g,k} - M} \sum_{g=g_1}^{g_a} \sum_{k=k_1}^{k_b} \sum_{i=1}^{n_{g,k}} (\hat{z}_{i,g,k}^{(t)} - \hat{\bar{Z}}_{g,k}^{(t)})^2. \tag{15}$$

Here $M$ denotes the number of strata that are pooled. Since the pooled estimator (15) assumes equal within-stratum variances for the strata that are pooled, it is not necessary to account for unequal sampling fractions in the different strata.

The growth rates for the turnover for the PPC's under the NACE Rev. 2 are estimated analogously to (4) as

$$\hat{Q}_l^{(t)} = \frac{\hat{Y}_l^{(t)}}{\hat{Y}_l^{(t-1)}}. \tag{16}$$

The total turnover for PPC $l$ can be estimated with, for example, the following Hájek-type domain estimator:

$$\hat{Y}_l^{(t)} = \frac{\sum_{i \in s} \frac{y_i^{(t)} \delta_i^{(l)}}{\pi_i}}{\sum_{i \in s} \frac{\delta_i^{(l)}}{\pi_i}} N_l = \frac{\sum_{k=1}^{K} \sum_{g=1}^{G} \frac{N_{g,k}}{n_{g,k}} \sum_{i=1}^{n_{g,k}} y_{i,g,k}^{(t)} \delta_i^{(l)}}{\sum_{k=1}^{K} \sum_{g=1}^{G} \frac{N_{g,k}}{n_{g,k}} \sum_{i=1}^{n_{g,k}} \delta_i^{(l)}} N_l. \tag{17}$$

Here $\delta_i^{(l)}$ is an indicator variable taking value 1 if sampling unit $i$ is classified to the $l$-th PPC and zero otherwise:

$$\delta_i^{(l)} = \begin{cases} 1 & \text{if} \quad i \in U_l \\ 0 & \text{if} \quad i \notin U_l \end{cases}. \tag{18}$$

An expression for $\hat{Y}_l^{(t-1)}$ is defined analogously to (17) with $y_{i,g,k}^{(t)}$ replaced by $y_{i,g,k}^{(t-1)}$. Note that (16) can be written as

$$\hat{Q}_l^{(t)} = \frac{\sum_{k=1}^{K} \sum_{g=1}^{G} \frac{N_{g,k}}{n_{g,k}} \sum_{i=1}^{n_{g,k}} y_{i,g,k}^{(t)} \delta_i^{(l)}}{\sum_{k=1}^{K} \sum_{g=1}^{G} \frac{N_{g,k}}{n_{g,k}} \sum_{i=1}^{n_{g,k}} y_{i,g,k}^{(t-1)} \delta_i^{(l)}}. \tag{19}$$

An approximately design-unbiased estimator for the variance of (19) is given by:

$$V\hat{a}r(\hat{Q}_l^{(t)}) = \frac{1}{(\hat{Y}_l^{(t-1)})^2} \sum_{k=1}^{K} \sum_{g=1}^{G} N_{g,k}(N_{g,k} - n_{g,k}) \frac{\hat{S}_{g,k}^{(t)2}}{n_{g,k}}, \tag{20}$$

with $\hat{S}_{g,k}^{(t)2}$ and $\hat{\bar{Z}}_{g,k}^{(t)}$ defined by (10) and (12), respectively with

$$\hat{z}_{i,g,k}^{(t)} = y_{i,g,k}^{(t)} \delta_i^{(l)} - \hat{Q}_l^{(t)} y_{i,g,k}^{(t-1)} \delta_i^{(l)}. \tag{21}$$

In the case of the GREG-estimator, (19) simply reads as

$$\hat{Q}_{l;greg}^{(t)} = \frac{\sum_{k=1}^{K} \sum_{g=1}^{G} \sum_{i=1}^{n_{g,k}} w_{i,g,k}^{(t)} y_{i,g,k}^{(t)} \delta_i^{(l)}}{\sum_{k=1}^{K} \sum_{g=1}^{G} \sum_{i=1}^{n_{g,k}} w_{i,g,k}^{(t-1)} y_{i,g,k}^{(t-1)} \delta_i^{(l)}}, \tag{22}$$

with $w_{i,g,k}^{(t)}$ the regression weights as described below formula (7). The variance of (22) can be estimated using formula

(20), where $\hat{Y}_k^{(t-1)}$ is replaced by $\hat{Y}_{k;\mathrm{greg}}^{(t-1)}$ and the residuals in (21) are defined by

$$\hat{z}_{i,g,k}^{(t)} = [y_{i,g,k}^{(t)} - \hat{\boldsymbol{b}}_k^{(t)'} \boldsymbol{x}_{i,g,k}^{(t)} - \hat{Q}_{k;\mathrm{greg}}^{(t)}(y_{i,g,k}^{(t-1)} - \hat{\boldsymbol{b}}_k^{(t-1)'} \boldsymbol{x}_{i,g,k}^{(t-1)})]\delta_i^{(l)}. \tag{23}$$

The major drawback of this approach is that there is no control over the sample sizes in the PPC's under the new classification, since this variable is not used in the stratification. As a result there will be PPC's with a small number of observations. The design variances will be unacceptably large for these weak domains. One possibility to avoid large design variances is to control the sample size in the PPC's under the old and the new classification by stratifying to the NACE Rev. 1.1 and the NACE Rev. 2. Another possibility is to improve the precision of the domain estimators by using a model-based small area estimator for the domains under the new classification.

The domain estimators (17) and (19) are design unbiased but become unstable in some situations. For example in situations where sample units under the NACE Rev. 1.1. belonging to different PPC's transfer to the same PPC under the NACE Rev. 2 and are selected with different sample fractions. This gives rise to large variation between the design weights. If these sampling units are also heterogeneous, then this will result in unstable domain estimators accompanied by large design variances.

One option is to change the design weights, for example by treating the sample as if it is selected by stratified simple random sampling, where the stratification variable is based on the crossing of size class and the publication cells based on the NACE Rev. 2. This results in more stable estimates and smaller standard errors but it will introduce design bias since the design weights are modified. Subject matter knowledge can be used to judge whether this approach results in an improvement of the estimates. Alternative solutions are drawing additional samples in combination with the correct design-based estimator (section 4) or applying a model-based small area estimator (section 6).

## 4 Stratifying to NACE Rev. 1.1 and NACE Rev. 2

The standard approach to achieve sufficiently reliable estimates for the PPC's under the old and the new classification is to stratify to both domain classifications and calculate the minimum sample size for each domain to guarantee a pre-specified precision. This implies that for the year of double reporting the stratification of the sample design changes to the full crossing of:

- size based on employment in 10 classes (abbreviated with the subscript $g$)
- PPC's based on the NACE Rev. 1.1 (abbreviated with the subscript $k$)
- PPC's based on the NACE Rev. 2 (abbreviated with the subscript $l$)

To settle the sample size and allocation, decisions about the type of allocation and minimum precision requirements must be made. Common allocations are proportional allocation, optimal or Neyman allocation, and power allocations.

Let $\hat{Q}_q^{(t)}$ denote the estimated growth rate of the turnover for period $t$ and domain $q$. These domains are:

- PPC's under the NACE Rev. 1.1 (in which case $q$ equals $k$)
- PPC's under the NACE Rev. 2 (in which case $q$ equals $l$)
- Aggregates of the PPC's under NACE Rev. 1.1 or NACE Rev. 2

Let $d_q^{(t)}$ denote the pre-specified maximum absolute deviation between the real growth rate of the turnover $Q_q^{(t)}$ and its estimate $\hat{Q}_q^{(t)}$, that is $\left|\hat{Q}_q^{(t)} - Q_q^{(t)}\right| \le d_q^{(t)}$. If it is conjectured that $\hat{Q}_q^{(t)}$ is a normally distributed random variable and if it is required that the probability that $\left|\hat{Q}_q^{(t)} - Q_q^t\right| > d_q^{(t)}$ must be smaller then $\alpha$, then it follows that the variance of $\hat{Q}_q^{(t)}$ is bounded by:

$$Var(\hat{Q}_q^{(t)}) \le \left(\frac{d_q^{(t)}}{Z_{(1-\alpha/2)}}\right)^2, \text{with}$$

$$Var(\hat{Q}_q^{(t)}) = \frac{1}{(\hat{Y}_q^{(t-1)})^2} \sum_{k \in q} \sum_{l \in q} \sum_{g \in q} N_{g,k,l}(N_{g,k,l} - n_{g,k,l})\frac{\hat{S}_{g,l,k}^{(t)\,2}}{n_{g,k,l}}. \tag{24}$$

Here $Z_{(\gamma)}$ is the $\gamma$-th percentile point of the standard normal distribution. Generally $\alpha$ is set to 5%, so $Z_{0.975} = 1.96$. The sample size for each $n_{g,l}$ is obtained by assuming optimal allocation within each PPC:

$$n_{g,k,l} = n_q \frac{N_{g,k,l}\hat{S}_{g,k,l}}{\sum_{k \in q} \sum_{l \in q} \sum_{g \in q} N_{g,k,l}\hat{S}_{g,k,l}}. \tag{25}$$

Substituting (25) in () and solving for $n_q$ gives the following expression for the minimum sample size within the $q$-th domain:

$$n_q = \frac{\left(\sum_{k \in q} \sum_{l \in q} \sum_{g \in q} N_{g,k,l}\hat{S}_{g,k,l}^{(t)}\right)^2}{\left(\hat{Y}_q^{(t-1)}\right)^2 \left(\frac{d_q^{(t)}}{Z_{(1-(1/2)\alpha)}}\right)^2 + \sum_{k \in q} \sum_{l \in q} \sum_{g \in q} N_{g,k,l}(\hat{S}_{g,k,l}^{(t)})^2}. \tag{26}$$

In order to calculate sample sizes and allocations under this design, a double coded sample and BR for a preceding period must be available, since parameter estimates, population variance estimates and population totals are required for strata that are based on the NACE Rev. 2 classification.

There are several ways to proceed. One approach is to determine the sample size and allocation at the most detailed publication level, i.e. the PPC's under the old and new classification. Subsequently the precision obtained with this sample size and allocation for aggregates can be checked. Under this approach the precision for the PPC's is controlled. The allocation is not necessarily optimal for aggregates, resulting in insufficient precision for the estimates at an aggregate level.

An alternative approach is to determine the sample size and allocation for publication cells at an aggregate level, e.g. sections at two digits. This approach, however, will result in sub-optimal estimates at the level of the PPC's. The variation between the precision of the PPC's will increase, and as a result the precision of the estimates for some of the PPC's will be insufficient while others will be estimated unnecessarily precise.

Another possibility is to determine the sample size and allocation in two steps. First, a power allocation is applied to the estimates at an aggregate level assuming stratified simple random sampling where PPC's are considered as the strata. Power allocations can be used to find the right balance between the precision requirements for aggregates and strata (Bankier, 1988). After having determined the sample size and allocation over the PPC's, an optimal or a proportional allocation can be applied to the strata within each PPC.

Stratifying to the full crossing of size class, PPC's under NACE Rev 1.1 and NACE Rev. 2 can result in a very detailed stratification. To obtain stable estimates for the population variances within the strata, the pooled variance estimator (15) could be considered. Optimal allocations are in general not very robust for outliers. Therefore it will be necessary to smooth the sample fractions obtained with an optimal allocation manually. An alternative approach to avoid the problems with instable estimates for the population variances, is to base the optimal allocation on an auxiliary variable that is available from a register for the entire population, and correlates well with the target parameter, e.g. value added tax.

Under this stratification scheme, the domains under the NACE Rev. 1.1 and 2 are both controlled. Estimates for both domains are obtained with (4) and (5) or (6). In the case of the GREG estimator, it might be efficient to stratify the auxiliary information in the weighting scheme to the classification of both the NACE Rev. 1.1 and 2. The level of detail depends on the available sample size.

## 5 Stratifying to the NACE Rev. 2

Another approach is to base the stratification on the crossing of size class and the PPC's under the NACE Rev. 2. This stratification will finally be used after the implementation of the NACE Rev. 2. During the period of double reporting, estimates for the NACE Rev. 2 domains are obtained by estimators for stratified simple random sampling, that is (4) and (5) or (6), but now applied to the domains under the NACE Rev. 2. Estimates for the NACE Rev 1.1 domains are now obtained with estimator (19) or (22).

Sample size and allocation is based on stratified simple random sampling where the stratification is based on the crossing of size class and the PPC's under the NACE Rev. 2. The procedure set out in the preceding section can be applied in an equivalent way to this design. An additional complication is that the stratum population variances $S_{g,l}^2$ must be estimated from a sample obtained by stratified simple random sampling where the stratification is based on the crossing of size class and the PPC's under the NACE Rev. 1.1. Sample units that are classified to the same stratum $(g,l)$ can be

selected with unequal selection probabilities, since they originate from different strata under the NACE Rev. 1.1 classification.

A design-unbiased estimator for the population variance $S_{g,l}^2$, that accounts for unequal selection probabilities for the units belonging to stratum $(g,l)$, is given by:

$$\hat{S}_{g,l}^2 = \frac{1}{2N_{g,l}(N_{g,l}-1)} \sum_{i=1}^{n_{g,l}} \sum_{j \neq i}^{n_{g,l}} \frac{(z_i^{(t)} - z_j^{(t)})^2}{\pi_{ij}}. \qquad (27)$$

The joint inclusion probabilities $\pi_{ij}$ are defined by (2).
The proof that (27) is a design-unbiased estimator for the population variances $S_{g,l}^2$ proceeds as follows. Let $a_i$ denote the indicator variable taking value 1 if unit $i$ is selected in the sample and zero otherwise:

$$a_i = \begin{cases} 1 & \text{if} \quad i \in s \\ 0 & \text{if} \quad i \notin s \end{cases}. \qquad (28)$$

Now expression (27) can be written as:

$$\hat{S}_{g,l}^2 = \frac{1}{2N_{g,l}(N_{g,l}-1)} \sum_{i=1}^{N_{g,l}} \sum_{j \neq i}^{N_{g,l}} \frac{(z_{i,g,l}^{(t)} - z_{j,g,l}^{(t)})^2}{\pi_{ij}} a_i a_j. \qquad (29)$$

That (27) is a design-unbiased estimator for $S_{g,l}^2$ follows by taking the expectation with respect to the sample design conditionally on the realised sample and its allocation over the strata. The expectation of the product of two sample membership indicators with respect to the sample design is by definition equal to the joint inclusion probability, that is $E(a_i a_j) = \pi_{ij}$. Since the sample membership indicators $a_i$ are the only random variables with respect to the sample design, it follows that:

$$E(\hat{S}_{g,l}^2) = \frac{1}{2N_{g,l}(N_{g,l}-1)} \sum_{i=1}^{N_{g,l}} \sum_{j \neq i}^{N_{g,l}} \frac{(z_{i,g,l}^{(t)} - z_{j,g,l}^{(t)})^2}{\pi_{ij}} E(a_i a_j)$$

$$= \frac{1}{2N_{g,l}(N_{g,l}-1)} \sum_{i=1}^{N_{g,l}} \sum_{j \neq i}^{N_{g,l}} (z_{i,g,l}^{(t)} - z_{j,g,l}^{(t)})^2$$

$$= \frac{1}{2N_{g,l}(N_{g,l}-1)} \sum_{i=1}^{N_{g,l}} \sum_{j \neq i}^{N_{g,l}} ((z_{i,g,l}^{(t)})^2 + (z_{j,g,l}^{(t)})^2 - 2z_{i,g,l}^{(t)} z_{j,g,l}^{(t)})$$

$$= \frac{1}{2N_{g,l}(N_{g,l}-1)} \left[ 2(N_{g,l}-1) \sum_{i=1}^{N_{g,l}} (z_{i,g,l}^{(t)})^2 - 2 \sum_{i=1}^{N_{g,l}} \sum_{j \neq i}^{N_{g,l}} z_{i,g,l}^{(t)} z_{j,g,l}^{(t)} \right]$$

$$= \frac{1}{2N_{g,l}(N_{g,l}-1)} \left[ 2N_{g,l} \sum_{i=1}^{N_{g,l}} (z_{i,g,l}^{(t)})^2 - 2 \sum_{i=1}^{N_{g,l}} \sum_{j=1}^{N_{g,l}} z_{i,g,l}^{(t)} z_{j,g,l}^{(t)} \right]$$

$$= \frac{1}{(N_{g,l} - 1)} \left[ \sum_{i=1}^{N_{g,l}} (z_{i,g,l}^{(t)})^2 - N_{g,l} \bar{Z}_{g,l}^2 \right]$$

$$= \frac{1}{(N_{g,l} - 1)} \sum_{i=1}^{N_{g,l}} (z_{i,g,l}^{(t)} - \bar{Z}_{i,g,l})^2 = S_{g,l}^2.$$

If the estimates for $S_{g,l}^{(t) \, 2}$ are unstable, the population variance estimates can be pooled. Suppose that the within stratum variances of the strata of size class $g = g_1, ..., g_a$ and PPC's $l = l_1, ..., l_b$ are equal. In this situation a pooled estimator for the population variances is obtained by the weighted average:

$$\hat{S}_{g,l,(P)}^{(t) \, 2} = \sum_{g=g_1}^{g_a} \sum_{l=l_1}^{l_b} \frac{N_{g,l}}{\sum_{g=g_1}^{g_a} \sum_{l=l_1}^{l_b} N_{g,l}} \hat{S}_{g,l}^{(t) \, 2}. \qquad (30)$$

Similar to the stratification proposed in section 4, a double coded sample and BR for a preceding period must be available to calculate the sample size and allocation for this stratification scheme.

For the STS at Statistics Netherlands the stratification for the year of double reporting will be based on size class crossed with the PPC's under the NACE Rev. 2. The samples for the STS are, however, based on a rotating panel. Each year a fraction of about 10% of the businesses in the panel are replaced by a sample of new businesses. In general it takes three or four months before the sample of new businesses that enter the panel has reached an acceptable response level. The major drawback of an optimal allocation under the NACE Rev. 2 is that this results in a large fraction of the businesses in the existing panel to be replaced by new businesses. This will result in an unacceptable loss of accuracy in the first months after the change-over to the new sample. Kish and Scott (1971) discuss sampling techniques to retain a maximum amount of sampling units after changing the stratification scheme of repeatedly conducted survey samples. For the STS at Statistics Netherlands, the following approach is adopted.

In a first step the sample fractions for the new strata are derived from the existing strata. If a stratum under the old classification entirely transfers 1-to-1 to a new stratum or if a stratum splits in two or more new strata (1-to-$n$ transitions), then the sample fractions from the strata under the NACE Rev. 1.1 will be applied to the new strata of the NACE Rev. 2. In the case that two or more existing strata under the NACE Rev 1.1 transfer to 1 new stratum ($m$-to-1 or $m$-to-$n$ transition), then the sample fraction in the new stratum is derived as an average of the sample fractions in the old strata weighted with the population sizes. If A denotes the union of strata under the NACE Rev. 1.1 that are joined in stratum of size class $h$ and PPC $l$ under the NACE Rev. 2, then the sample fraction for this new stratum is given by

$$f_{h,l} = \frac{\sum_{g,k \in A} f_{g,k} N_{g,k}}{\sum_{g,k \in A} N_{g,k}}, \qquad (31)$$

with $f_{g,k}$ the sample fraction in the stratum $(g,k)$. In the case of large deviations from the optimal allocation, the sample fractions are adjusted to guarantee sufficient precision.

To achieve a sample that can be considered as obtained by stratified simple random sampling, sampling units are removed from or added to the existing sample as follows. If $f_{g,k} > f_{h,l}$, then a simple random sample from the sample of stratum $(g,k)$ that transfers to $(h,l)$ is removed such that the sample fraction (approximately) equals (31). If $f_{g,k} < f_{h,l}$, then a simple random sample from the subpopulation of stratum $(g,k)$ that transfers to $(h,l)$ is added to the existing sample, such that the sample fraction (approximately) equals (31).

## 6 Small Area Estimation

The major drawback of stratifying to the NACE Rev 1.1 is that the sample size in the domains of the NACE Rev. 2 are not controlled, which can result in unacceptable large standard errors for some of these domains. The same problem can occur for the domains under the NACE Rev. 1.1 if the NACE Rev. 2 is used as a stratification variable. Instead of drawing additional samples, model-based estimation procedures may be considered to improve the precision of the estimates in the weak domains.

The design-based estimation procedures considered in the preceding sections are widely applied by national statistical institutes. The main advantage of the classical design-based approach is that these estimators are always (approximately) design unbiased. As a result these estimators have a built-in robustness against model-misspecification. These properties also hold for GREG and calibration estimators that incorporate available auxiliary information in the estimation procedure. Another advantage is that only one set of weights needs to be derived to estimate all possible target parameters. This is not only convenient for multi-purpose surveys, but also has the advantage that the various output tables will be consistent. These properties make the design-based estimators very appropriate to apply in a statistical process where there is generally limited time available for the analysis phase.

The major drawback of the design-based approach, however, is the unacceptably large standard errors in the case of small sample sizes. Instead of increasing sample sizes, estimation procedures can be considered that explicitly rely on a statistical model to improve the precision of domain estimates with sample information observed in other domains or preceding time periods. This is the realm of small area estimation. For a comprehensive overview, see Rao (2003). A briefer but very nice overview is given by Pfeffermann (2002).

There is a wide range of methods available in the literature of small area estimation. A potential approach for the STS is the so-called area level model, developed by Fay and Herriot (1979). In this approach the direct estimates for the domains are modelled with a mixed model:

$$\hat{\theta}_q = \theta_q + e_q, \quad e_q \cong N(0, \psi_q), \qquad (32)$$

$$\theta_q = \boldsymbol{\beta}' \boldsymbol{x}_q + v_q, \qquad v_q \cong N(0, \sigma_v^2). \qquad (33)$$

Here $\hat{\theta}_q$ denotes the direct estimator for the unknown domain parameter $\theta_q$ for domain $q$, $e_q$ the sample error, $\psi_q$ the design variance of $\hat{\theta}_q$. The model incorporates available auxiliary information $x_q$ on the level of the domains, for example value added or the monthly growth rate of value added that might be available from tax registers. The domains are linked through the common fixed regression coefficients $\boldsymbol{\beta}$. The unexplained variation between the domains is modelled with the random domain effects $v_q$.

Equations (32) and (33) describe a linear mixed model for the domain parameters. Under this model an empirical best linear unbiased predictor (EBLUP) can be derived to estimate the unknown domain parameters, see Rao (2003), section 6.2 for an expression. This EBLUP-estimator can be expressed as a weighted average of the direct estimator $\hat{\theta}_q$ and the synthetic regression estimator $\boldsymbol{\beta}' \boldsymbol{x}_q$ where the weights are based on the variance estimates of both components:

$$\tilde{\theta}_q = \hat{\gamma}_q \hat{\theta}_q + (1 - \hat{\gamma}_q) \boldsymbol{\beta}' \boldsymbol{x}_q, \qquad (34)$$

$$\hat{\gamma}_q = \frac{\hat{\sigma}_v^2}{\hat{\sigma}_v^2 + \psi_q}. \qquad (35)$$

See Rao (2003) section 7.1 for an estimator of $\boldsymbol{\beta}$ and various methods for estimating $\sigma_v^2$. An appealing property of the area level model for this application is that the direct domain estimates are the input for the model and therefore accounts for the applied sampling design. Moreover, value added is a potentially strong auxiliary variable, but generally not available at the unit level for all business units. The area level model, nevertheless, makes advantage of the available value added information at the domain level.

The direct estimates for domains with large sample sizes will have small design variances. In these cases, the model-based estimates for the domain parameters obtained with (34), are largely based on the direct estimator since $\hat{\gamma}_q$ in (35) tends to one. Direct estimates for the domains with small sample size or large fluctuations in the design weights, will have large design variances. This results in more emphasis on the synthetic regression part of the EBLUP estimator in (34) since $\hat{\gamma}_q$ tends to zero. For domains were no observations are available at all, the EBLUP estimator is completely based on the regression part since the variance estimator for the direct estimator goes to infinity. Therefore the small area estimation approach might provide a solution for the domains that have been out of scope under the NACE Rev. 1.1 and enter the domains of the NACE Rev. 2.

If the auxiliary information is available at the unit level, then it is also possible to specify a multi-level model on the unit level that is originally proposed by Battese, Harter and Fuller (1988). This approach has the advantage that it uses the auxiliary information in a more efficient way and has more degrees of freedom for parameter and variance estimation. The major drawback is that this approach generally assumes self-weighted samples, which is not the case in this application. One solution is to incorporate the sample design into the model, which requires additional work on the modelling part.

## 7 Backcasting procedures

Replacing the NACE Rev. 1.1 by the NACE Rev. 2 results in disrupted time series. A part of the implementation process concerns the reconstruction of historical series for the domains under the NACE Rev. 2. This is generally referred to as backcasting. There are two important reasons for backcasting series for the NACE Rev. 2 domains. First, consistent series are of crucial importance for many users. Second, sufficiently long series are required to construct seasonally adjusted series for the domains under the new classification directly after the change-over to the new classification system. Eurostat (2006d), James (2008), and Buiten et al. (2008) describe various backcasting procedures.

Usually a distinction is drawn between backcasting procedures that operate on the level of business units and on an aggregated level, e.g. strata or publication domains. The first are the so-called micro approaches, while the latter are referred to as the macro approaches.

### 7.1 Micro approach

The micro approach implies that the individual business units in the samples observed in the past, and preferably also the BR are classified with respect to the NACE Rev. 2, resulting in a double coded sample or BR. Subsequently, estimates for the domains under the NACE Rev. 2 are calculated using the same design-based approach described in section 3, for example the domain estimators defined in (17) and (19). If it is not possible to recode the BR, then the Hájek-type domain estimator (17) must be replaced by the Horvitz-Thompson estimator for a domain total. Problems with large design variances due to small sample sizes in weak domains, or instable domain estimators due to extreme variability in the design weights might be overcome with the model-based estimation techniques from the theory of small area estimation. These approaches might also be applied to obtain estimates for the domains that have been out of scope under the old classification system.

The area level model, proposed in section 6, can be used to obtain model-based estimates at each period in time, where sample information from neighbouring domains is used to improve the precision for the estimates in the weak domains. Since time series for the NACE Rev. 2 domains are reconstructed, it will be efficient to apply an estimation approach that combines sample information from different domains with sample information observed in preceding periods. Rao and Yu (1994) extended the area level model with a first order autoregressive component to combine cross-sectional sample information with sample information observed in preceding periods. A different approach is followed by Pfeffermann and Burck (1990) and Pfeffermann and Bluer (1993). They combine time series data with cross-sectional data by modelling the correlation between domain parameters in a multivariate structural time series model. The gen-

eral finding in the literature is that methods based on time-series data result in more precise domain estimates compared to cross-sectional data, Eurarea (2004), Boonstra et al. (2008).

The main advantage of the micro approach is that the estimated series are still based on empirical evidence. As a result, the structural evolution of the economy will be better retained compared to the macro approach that strongly relies on synthetic estimation procedures. The major drawback is that it requires the availability of micro data and more resources for double coding of the sample or the BR in the past. Also the computations are, compared to the macro approach, more intensive.

It is worthwhile to consider the micro approach if the NACE Rev. 2 classification codes are available for the business units in preceding periods or can be derived in a relatively straightforward manner. At Statistics Netherlands, for example, the STS for industry are based on a complete enumeration of the strata with large and intermediate size classes. In this case the number of business units in the sample is relatively small and sufficient retrorespective data are available to derive the NACE Rev. 2 classification for preceding time periods. Therefore the micro approach will be applied in this situation. Such considerations might also apply for panel designs, where sufficient information is available to derive the NACE Rev. 2 codes automatically, or NACE Rev. 2 codes can be imputed through recoding of the main activity using transition or conversion schemes.

### 7.2 Macro approach

In many situations it will not be feasible to apply the micro approach since it is very time costly and often requires the collection of additional retrorespective data to recode the business units for the NACE Rev. 2 classification. In such situations the so-called macro approach can be considered for backcasting purposes. The macro approach can also be used as an alternative for the micro approach, if the direct estimators mentioned in section 7.1 are unstable or have unacceptably large standard errors due to small sample sizes in the weak domains. From this point of view, the macro approach is a synthetic form of small area estimation, based on naive implicit models.

The macro approach implies that estimates for the domains under the NACE Rev. 2 are derived from a linear combination of the estimates for the domains under the NACE Rev. 1.1. For example the total turnover for the $l$-th domain of the NACE Rev. 2 is calculated as

$$\tilde{Y}_l^{(t)} = \sum_k \beta_{k,l} \hat{Y}_k^{(t)}, \qquad (36)$$

where $\hat{Y}_k^{(t)}$ is a direct estimator for the total turnover in the $k$-th domain of the NACE Rev. 1.1, and $\beta_{k,l}$ a conversion factor specifying the fraction $\hat{Y}_k^{(t)}$ that transfers from the $k$-th domain under NACE Rev. 1.1 to the $l$-th domain under NACE Rev.2. The conversion factors are fractions that specify the distribution of $\hat{Y}_k^{(t)}$ over the classes of NACE Rev. 2, that is

$$\sum_l \beta_{k,l} = 1. \qquad (37)$$

The conversion factors can be obtained in several ways and are often derived from so-called transition matrices, Eurostat (2006d), James (2008). The entries for the rows correspond to the NACE Rev. 1.1. classes and the columns to the NACE Rev. 2 classes. The cells of these matrices specify a variable of interest that transfers from class $k$ under the NACE Rev. 1.1 to class $l$ under the NACE Rev. 2 and is denoted by $X_{k,l}$. Possible variables are the number of business units, estimated total turnover from STS or SBS during the year of double reporting, total value added, or number of employees. The conversion factors $\beta_{k,l}$ are easily derived from these matrices by dividing the cells by the column total:

$$\beta_{k,l} = \frac{X_{k,l}}{\sum_l X_{k,l}}. \qquad (38)$$

The advantage of using auxiliary register information to construct conversion factors is the absence of sampling error. The economic structure that is assumed with (38), however, might differ substantially between the various auxiliary variables that are available. As a result, the evolution of the backcasted series mainly depends on the choice of the auxiliary variable and its validity is mainly determined by the correlation between the auxiliary variable and the target variable. A natural choice is to use the same variables that are used as auxiliary information in the calibration estimator for the target variable to be backcasted, James (2008).

As an alternative, direct estimates for the target parameter obtained with STS or SBS can be used to construct the conversion factors. This avoids the choice between different auxiliary variables but may result in unstable estimates for the conversion factors due to sampling error.

If the estimated turnover is used to construct the conversion factors, then a domain estimator like (17) can be used. In this case it follows that

$$X_{k,l} = \hat{Y}_{k,l}^{(T)} = \frac{\sum_{g=1}^{G} \frac{N_{g,k}}{n_{g,k}} \sum_{i=1}^{n_{g,k}} y_{i,g,k}^{(T)} \delta_i^{(l)}}{\sum_{g=1}^{G} \frac{N_{g,k}}{n_{g,k}} \sum_{i=1}^{n_{g,k}} \delta_i^{(l)}} N_{k,l}, \qquad (39)$$

where $T$ refers to the period of double reporting. If the Horvitz-Thompson estimator is used, then it follows that

$$X_{k,l} = \hat{Y}_{k,l}^{(T)} = \sum_{g=1}^{G} \frac{N_{g,k}}{n_{g,k}} \sum_{i=1}^{n_{g,k}} y_{i,g,k}^{(T)} \delta_i^{(l)}. \qquad (40)$$

The advantage of this estimator is that in the year of double reporting it follows that

$$X_{k+} = \sum_l \hat{Y}_{k,l}^{(T)} = \sum_{g=1}^{G} \frac{N_{g,k}}{n_{g,k}} \sum_{i=1}^{n_{g,k}} y_{i,g,k}^{(T)}, \qquad (41)$$

which is equal to the direct estimator (5) for the domains under the NACE Rev. 1.1. As a result it follows that

$$\tilde{Y}_l^{(T)} = \sum_k \hat{Y}_{k,l}^{(T)} = \sum_{k=1}^{K} \sum_{g=1}^{G} \frac{N_{g,k}}{n_{g,k}} \sum_{i=1}^{n_{g,k}} y_{i,g,k}^{(T)} \delta_i^{(l)}. \qquad (42)$$

If the Horvitz-Thompson estimator instead of the Hájek estimator (17) is used for the domains under the NACE Rev. 2, then it follows that the backcasted values equals their direct estimates in the year of double reporting. With this approach, no discontinuities occur in the series of the NACE Rev. 2 domains at the moment that the series change from a backcasting approach to a direct estimation procedure. In other cases some kind of linking procedure might be necessary to deal with this kind of discontinuities (section 7.5).

In formula (36) linear combinations of classes under the NACE Rev. 1.1 are used to backcast the series for the domains under the NACE Rev. 2. Instead of working on a four digit level, it is also possible to work on a more detailed or aggregated level. The lower the level of aggregation, the better the real evolution of the economy is retained. Choosing a low level of aggregation, however, might result in instable estimates for the conversion factors and therefore also for the backcasted domains. Using the SBS for constructing conversion factors has the advantage that the direct estimators are more precise since the sample size of the SBS is generally larger compared to the STS. Small area estimators might also be considered as input for the construction of the conversion factors.

## 7.3 Macro approaches with time dependent conversion factors

The main disadvantage of the macro approach is that it is based on very strong assumptions. Using time independent conversion factors assumes that the economic structure observed in the period to construct the conversion factors is constant over time. Generally, this assumption will not be met. Particularly for new activities this approach easily results in an unrealistic evolution of the backcasted indicators. Therefore it is worthwhile to consider the application of time dependent conversion factors.

One option is to combine the micro and macro approach. The micro approach, for example, can be applied for one or two years in the past, preferably the base years to compile indices. Conversion factors can be constructed for these years. Subsequently the conversion factors for the intervening years can be derived through linear or non-linear interpolation. It is also possible to use the micro approach for the most recent years, or extend the period of double coding and reporting after the change-over to the NACE Rev. 2. This offers the possibility to evaluate the assumption that conversion factors are time independent and to construct time dependent conversion factors that allow for trend or seasonal patterns if necessary.

Subject matter specialists can and should be consulted to judge whether the evolution of a backcasted series seems realistic. Such subject matter knowledge might also be useful to adjust the conversion factors. For example to make

decisions about the moment that innovations and new economic activities are introduced, including realistic interpolation functions for the conversion factors between this moment and the year of double coding or double reporting.

Another approach is to construct transition matrices and conversion factors for the separate years. This might be an option if the BR or the SBS can be double coded in a relative straightforward way, via an automatic procedure.

## 7.4 Backcasting indices

Most target parameters of STS's are defined as indices. One way to proceed is to backcast the underlying series for total turnover. Also the SBS in the base year must be backcasted for the purpose of deriving weights for aggregating the indices from classes to groups, divisions or sections. For this purpose the micro as well as the macro approach, discussed in the preceding sections, can be used. A more detailed discussion is provided by James (2008).

An alternative approach, appealing due to its simplicity, is described in Eurostat (2006d), section 2.2.3. According to this approach, indices are backcasted in two steps. First a transition matrix is constructed for the variable that is used to construct weights for aggregating indices, for example the total value added. This is generally accomplished with the macro approach described in section 7.2 or 7.3, but it is also possible to use the micro approach described in section 7.1. In the second step, the distribution of the total value added over the NACE Rev. 1.1 domains within a domain of the NACE Rev. 2 are calculated, that is

$$\varphi_{k,l} = \frac{X_{k,l}}{\sum_k X_{k,l}}, \qquad (43)$$

with $X_{k,l}$ the total value added that transfers from domain $k$ under the NACE Rev. 1.1 to domain $l$ under the NACE Rev. 2. Formula (43) specifies the distribution of the total value added over the domains under the NACE Rev. 1.1 within a domain of the NACE Rev. 2, so $\sum_k \varphi_{k,l} = 1$. Note the difference with (38), which specifies the distribution over the domains under the NACE Rev. 2 within a domain of the NACE Rev. 1.1. The conversion factors defined by (43) are the weights to be used in (36) to backcast or convert the indices from the NACE Rev. 1.1 to the NACE Rev. 2.

## 7.5 Linking series

Another consequence of applying backcasting procedures is that discontinuities may occur in the series for the domains under the NACE Rev. 2 at the moment that the macro approach changes to the micro approach during the period that a backcasting procedure is used or at the moment of the change-over from the backcasting procedure to the direct estimation approach after the implementation of the NACE Rev. 2 as the regular classification system. A structural time series model with an intervention variable that models both types of change-over could be used to quantify these discontinuities. These models can also be applied as a linking procedure to restore the continuity of these series.

See Van den Brakel et al. (2008) for details and alternative linking procedures, for example based on simple ratios.

## Discussion

In this paper a set of sampling and estimation techniques are reviewed that can facilitate a smooth transition from the NACE Rev. 1.1 to the NACE Rev. 2 in business statistics.

The first step of the transition is the implementation of the new classification system in the BR. Having a double coded BR offers the possibility to produce figures under both classification systems simultaneously. Appropriate domain estimators for the domains under both classifications are available from classical sampling theory if a probability sample is used. Generally the domains are used in the stratification to control the sample size within each domain to meet pre-specified precision requirements. Stratifying to both classifications to meet the precision requirements for the domains under both classifications simultaneously, might result in a substantial increase of the sample size. The traditional design-based domain estimators, on the other hand, may result in unreliable estimates due to small sample sizes in domains under the classification that is not used as a stratification variable in the sample design. Model-based estimation procedures from the realm of small area estimation might be used as an alternative for drawing additional sampling units. The three different stratification schemes in combination with the design- and model-based estimation procedures, discussed in this paper, result in six different sampling strategies for the domains under both classifications during the period of double coding. The pros and cons of these six strategies are summarized in Table B.1 of Annex B.

There is a strong demand for producing historical time series for the domains under the new classification in the past. Many users require consistent series without discontinuities due to the introduction of a new classification system. Also for the purpose of studying cycles and producing seasonally adjusted series it is important to construct series under the NACE Rev. 2 in the past. For this purpose different backcasting procedures are described. The micro approaches, operating at the level of the sampling units, are essentially the traditional domain estimators from classical sampling theory. The advantage is that these approaches are design unbiased. The results can, however, still be unreliable for domains with small sample sizes. Another drawback is that this approach is costly and computation intensive.

Several macro approaches provide alternatives to the micro approach. These procedures operate at an aggregated level and predict the series for a domain under the NACE Rev. 2 as a linear combination from the domain estimates of the NACE Rev. 1.1. These approaches are less computation intense and can result in more stable estimates. They rely, however, on strong and often naive model assumptions, particularly if the transition coefficients are assumed to be constant over time because they are based on one period of double coding or double reporting only. This could result in strongly biased predictions for the domains under the NACE Rev. 2.

It is expected that more accurate predictions for the NACE Rev. 2 domains in the past can be obtained with more advanced model-based estimation procedures that are available from the theory of small area estimation. These procedures borrow sample information from other domains or previous time periods by relying explicitly on a mixed model or time series model. The underlying assumptions are generally more realistic compared to the synthetic procedures that predict the domains under the NACE Rev. 2 as a linear combination from the domain estimates of the NACE Rev. 1.1.

The small area estimation approach provides some useful solutions for problems encountered by the NACE Rev. 2 implementation. Depending on the available auxiliary information, it can be used to improve the precision of estimates for weak domains. These are for example the domains were large design variances occur due to small sample sizes or large fluctuations between sample fractions, resulting in instable parameter estimates. This approach is also useful to obtain synthetic regression estimates for the empty domains that have been out of scope under the old classification. The success of this approach strongly depends on the quality of the available auxiliary information. It can be expected that auxiliary information like value added, available from tax registers, strongly correlates with parameters as turnover.

Model-based estimation procedures require careful model selection and evaluation, since they are not robust for model misspecification. This could hamper the application in a statistical production process, where there is generally a limited amount of time available for the analysis phase to produce timely figures. Since STS generally have a limited set of target parameters, these obstructions may be manageable.

It can be concluded that three different classes of backcasting procedures are distinguished in this paper. The first approach is the micro approach in combination with design-based estimation procedures. The second one is also a micro approach in combination with model-based estimation procedures. The third one is the macro approach, which basically relies on very synthetic model-based procedures. The different properties of these three backcasting approaches are summarized in Table B.2 in Annex B.

## Acknowledgements

# References

Alexander, C. H. (1987). A Class of Methods for Using Person Controls in Household Weighting. *Survey Methodology*, *13*, 183-198.

Bankier, M. D. (1988). Power Allocations: Determining Sample Sizes for Subnational Areas. *The American Statistician*, *82*, 174-177.

Battese, G. E., Harter, R. M., & Fuller, W. A. (1988). An error components model for prediction of county crop areas using survey satellite data. *Journal of the American Statistical Association*, *83*, 28-36.

Boonstra, H. J., Van den Brakel, J. A., Buelens, B., Krieg, S., & Smeets, M. (2008). Towards small area estimation at Statistics Netherlands. *METRON, International Journal of Statistics*, *LXVI*, 21-49.

Buiten, G., Kampen, J. K., & Vergouw, S. (2008). *Theory on the producing of historical time series for Short-term Business Statistics in NACE Rev. 2 with an application in the industrial turnover index in the Netherlands (1995-2008).* Research paper, BPA nr.: DMK-2008-10-02-JKPN, Statistics Netherlands.

Cochran, W. G. (1977). *Sampling Techniques.* New York: Wiley and Sons.

Deville, J., & Särndal, C.-E. (1992). Calibration Estimators in Survey Sampling. *Journal of the American Statistical Association*, *87*, 376-382.

Eurarea. (2004). *Project reference volume, deliverable D7.1.4.* Technical report, EURAREA consortium.

Eurostat. (2006a). *Backcasting handbook. Task Force on the implementation of NACE Rev. 2.* http://circa.europa.eu/irc/dsis/nacecpacon/info/data/en/index.htm.

Eurostat. (2006b). *Handbook on methodological aspects related to sampling designs and weights estimation. Task Force on the implementation of NACE Rev. 2.* http://circa.europa.eu/irc/dsis/nacecpacon/info/data/en/index.htm.

Eurostat. (2006c). *Implementation of NACE Rev. 2 in Business Registers. Task Force on the implementation of NACE Rev. 2.* http://circa.europa.eu/irc/dsis/nacecpacon/info/data/en/index.htm.

Eurostat. (2006d). *Setting up an implementation plan for NACE Rev. 2 in National Statistical Institutes. Task Force on the implementation of NACE Rev. 2.* http://circa.europa.eu/irc/dsis/nacecpacon/info/data/en/index.htm.

Eurostat. (2008). *NACE Rev. 2 – Statistical classification of economic activities in the European Community.* Eurostat Methodologies and Working papers. http://circa.europa.eu/irc/dsis/nacecpacon/info/data/en/index.htm.

Fay, R. E., & Herriot, R. A. (1979). Estimates of income for small places: an application of James-Stein procedures to Census data. *Journal of the American Statistical Association*, *74*, 269-277.

James, G. (2008). *Backcasting for use in Short Term Statistics.* Interim Report from the UK Office for National Statistics.

Kish, L., & Scott, A. (1971). Retaining Units after Changing Strata and Probabilities. *Journal of the American Statistical Association*, *66*, 461-470.

Luery, D. M. (1986). *Weighting Sample Data Under Linear Constraints on the Weights.* Proceedings of the section on Social Statistics, American Statistical Association, pp. 325-330.

Pfeffermann, D. (2002). Small Area Estimation – New developments and directions. *International Statistical Review*, *70*, 125-143.

Pfeffermann, D., & Bleuer, S. R. (1993). Robust Joint Modelling of Labour Force Series of Small Areas. *Survey Methodology*, *19*, 149-163.

Pfeffermann, D., & Burck, L. (1990). Robust Small Area Estimation combining Time Series and Cross-sectional Data. *Survey Methodology*, *16*, 217-237.

Rao, J. N. K. (2003). *Small Area Estimation.* New York: Wiley and Sons.

Rao, J. N. K., & Yu, M. (1994). Small Area Estimation by combining Time-series and Cross-sectional Data. *The Canadian Journal of Statistics*, *22*, 511-528.

Särndal, C.-E., Swensson, B., & Wretman, J. (1992). *Model Assisted Survey Sampling.* New York: Springer Verlag.

Van den Brakel, J. A. (2008). Design-based analysis of embedded experiments with applications in the Dutch Labour Force Survey. *Journal of the Royal Statistical Society*, *Series A, 171*, 581-613.

Van den Brakel, J. A., & Roels, J. (2010). Intervention analysis with state-space models to estimate discontinuities due to a survey redesign. *Annals of Applied Statistics*, *4*, 1105-1138.

Van den Brakel, J. A., Smith, P., & Compton, S. (2008). Quality procedures for survey transitions – experiments, time series and discontinuities. *Survey Research Methods*, *2*, 123-141.

# Appendix A: Overview of the classification of NACE Rev. 1.1 and NACE Rev. 2

*Table A.1:* Change-over of the section of NACE Rev. 1.1 to the NACE Rev. 2

| NACE rev. 1.1 | | NACE Rev. 2 | |
|---|---|---|---|
| Section | Description | Section | Description |
| A | Agriculture, hunting and forestry | A | Agriculture, forestry and fishing |
| B | Fishing | | |
| C | Mining and quarrying | B | Mining and quarrying |
| D | Manufacturing | C | Manufacturing |
| E | Electricity, gas and water supply | D | Electricity, gas, steam and air conditioning supply |
| | | E | Water supply, sewerage, waste management and remediation activities |
| F | Construction | F | Construction |
| G | Wholesale and retail trade: repair of motor vehicles, motorcycles and personal and household goods | G | Wholesale and retail trade: repair of motor vehicles and motorcycles |
| H | Hotels and restaurants | I | Accommodation and food service activities |
| I | Transportation, storage and communication | H | Transportation and storage |
| | | J | Information and communication |
| J | Financial intermediation | K | Financial and insure activities |
| K | Real estate, renting and business activities | L | Real estate activities |
| | | M | Professional, scientific and technical activities |
| | | N | Administrative and support service activities |
| L | Public administration and defence; compulsory social security | O | Public administration and defence; compulsory social security |
| M | Education | P | Education |
| N | Health and social work | Q | Human health and social work activities |
| O | Other community, social and personal services activities | R | Arts, entertainment and recreation |
| | | S | Other services |
| P | Activities of private households as employers and undifferentiated production activities of private households | T | Activities of private households as employers; undifferentiated goods- and services-producing activities of households for own use |
| Q | Extraterritorial organizations and bodies | U | Activities of extraterritorial organizations and bodies |

Sections A and B under NACE Rev. 1.1 are joined into one section A under NACE Rev. 2.
Section E under NACE Rev. 1.1 is divided in two sections D and E under NACE Rev. 2.
Section I under NACE Rev. 1.1 is divided in two sections H and J under NACE Rev. 2.
Section K of NACE Rev. 1.1 is divided in three sections L, M and N under NACE Rev. 2.
Section O under NACE Rev. 1.1 is divided in two sections R and S under NACE Rev. 2.

JAN VAN DEN BRAKEL

*Table A.2:* Overview of the number of subsections, divisions, groups and classes within each section
of the NACE Rev. 1.1 and the NACE Rev. 2.

| NACE rev. 1.1 | | | | | NACE Rev. 2 | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| Section | Subsec. | Divisions | Groups | Classes | Section | Divisions | Groups | Classes |
| A | 1 | 2 | 6 | 14 | A | 3 | 13 | 39 |
| B | 1 | 1 | 1 | 2 | | | | |
| C | 2 | 5 | 13 | 16 | B | 5 | 10 | 15 |
| D | 14 | 23 | 103 | 242 | C | 23 | 95 | 230 |
| E | 1 | 2 | 4 | 7 | D | 1 | 3 | 8 |
| | | | | | E | 4 | 6 | 9 |
| F | 1 | 1 | 5 | 17 | F | 3 | 9 | 22 |
| G | 1 | 3 | 19 | 79 | G | 3 | 21 | 91 |
| H | 1 | 1 | 5 | 8 | I | 2 | 7 | 8 |
| I | 1 | 5 | 14 | 21 | H | 5 | 15 | 23 |
| | | | | | J | 6 | 13 | 26 |
| J | 1 | 3 | 5 | 12 | K | 3 | 10 | 18 |
| K | 1 | 5 | 23 | 39 | L | 1 | 3 | 4 |
| | | | | | M | 7 | 15 | 19 |
| | | | | | N | 6 | 19 | 33 |
| L | 1 | 1 | 3 | 10 | O | 1 | 3 | 9 |
| M | 1 | 1 | 4 | 6 | P | 1 | 6 | 11 |
| N | 1 | 1 | 3 | 7 | Q | 3 | 9 | 12 |
| O | 1 | 4 | 12 | 30 | R | 4 | 5 | 15 |
| | | | | | S | 3 | 6 | 19 |
| P | 1 | 3 | 3 | 3 | T | 2 | 3 | 3 |
| Q | 1 | 1 | 1 | 1 | U | 1 | 1 | 1 |

# Appendix B: Overview of sampling strategies and backcasting procedures

*Table B.1:* Overview of sampling strategies for domain estimation under NACE Rev 1.1 and NACE Rev. 2 during the period of double coding

| Property | Design-based estimation | | | Model-based estimation | | |
|---|---|---|---|---|---|---|
| | NACE Rev 1.1 | Stratification NACE Rev 1.1 and 2 | NACE Rev. 2 | NACE Rev 1.1 | Stratification NACE Rev 1.1 and 2 | NACE Rev. 2 |
| Sample size domains | Controlled for NR. 1.1, not for NR. 2. Empty domains for NR. 2 possible. | Controlled for NR 1.1 and NR. 2. | Controlled for NR. 2, not for NR. 1.1. Empty domains for NR. 1.1 avoided through maximizing sample overlap. | Controlled for NR. 1.1, not for NR. 2. Empty domains for NR. 2 possible. | Controlled for NR 1.1 and NR. 2. | Controlled for NR. 2, not for NR. 1.1. Empty domains for NR. 1.1 avoided through maximizing sample overlap. |
| Total sample size | Constant in time in case of steady precision requirements. | Temporary increase expected, depending on specified precision requirements. | More or less constant in time if precision requirements are comparable with NR 1.1 domains. Increase expected in case of sample overlap requirements. | Constant in time in case of steady precision requirements. | Temporary increase expected, depending on specified precision requirements. | More or less constant in time if precision requirements are comparable with NR 1.1 domains. Increase expected in case of sample overlap requirements. |
| Sample overlap with preceding periods | Maximum overlap during the period of double coding. Temporary decrease when the NR. 2 is used as stratification scheme after the period of double coding. | Controlled at the cost of an increased sample size. | Temporary decrease during the period of double coding. Maximum overlap after the period of double coding, since the NR. 2 is already in use as stratification scheme. | Maximum overlap during the period of double coding. Temporary decrease when the NR. 2 is used as stratification scheme after the period of double coding. | Controlled at the cost of an increased sample size. | Temporary decrease during the period of double coding. Maximum overlap after the period of double coding, since the NR. 2 is already in use as stratification scheme. |
| Change-over to final sample design | Introduction of NR. 2 as stratification scheme delayed. | NR. 2 is smoothly introduced in the stratification during the period of double coding. | NR.2 is directly introduced in the stratification during the period of double coding. | Introduction of NR. 2 as stratification scheme delayed. | NR. 2 is smoothly introduced in the stratification during the period of double coding. | NR. 2 is directly introduced in the stratification during the period of double coding. |
| Manual manipulation | Double coded BR and sample for a preceding period not required for calculation of sample size and allocation | Double coded BR and sample for a preceding period required for calculation of sample size and allocation | Double coded BR and sample for a preceding period required for calculation of sample size and allocation | Double coded BR and sample for a preceding period not required for calculation of sample size and allocation | Double coded BR and sample for a preceding period required for calculation of sample size and allocation | Double coded BR and sample for a preceding period required for calculation of sample size and allocation |

Table B.1: continued.

| Property | Design-based estimation | | | Model-based estimation | | |
|---|---|---|---|---|---|---|
| | Stratification | | | Stratification | | |
| | NACE Rev 1.1 | NACE Rev 1.1 and 2 | NACE Rev. 2 | NACE Rev 1.1 | NACE Rev 1.1 and 2 | NACE Rev. 2 |
| Estimator | Reliability NR. 1.1 domains comparable with preceding periods. Estimates are not available for weak or empty domains under NR. 2. | Reliability NR. 1.1 and NR. 2 domains controlled through pre-specified precision requirements. Depending on the sampleoverlap, estimates might not be available for weak or empty domains of NR. 1.1. | Reliability NR. 2 domains controlled through pre-specified precision requirements. Synthetic estimates for empty domains under NR. 2. Reliability depends on the availability of auxiliary information. | Improvement reliability of the weak domains under NR. requirements. | Not required in case of sufficient precision requirements. | Improvement reliability of the weak domains under NR. 1.1. Weak and empty domains under NR. 1.1. can be partially avoided through sufficient sample overlap requirements. |
| Variance | Controlled for NR. 1.1 domains. Large variances for weak domains under NR. 2. | Controlled for NR. 1.1 and NR. 2 domains. | Controlled for NR. 2 domains. Possibly large variances for weak domains under NR. 1.1 (depending on sample overlap). | Variance improvement through "borrowing strength" depending on pre-specified precision for weak domains of NR. 2. | Generally small variance improvements, depending on pre-specified precision for weak domains of NR. 1.1. | Variance improvement through "borrowing strength" over time or space for weak domains of NR. 1.1. |
| Bias | Approximately design unbiased | | | Size of design-bias depends on the quality of the selected model and availability of auxiliary information. | | |
| Accuracy | Depends on the sample size of the individual domains (see rows above for differences between stratification schemes). | | | Strong improvements possible for weak domains. And relatively small improvements for strong domains (see rows above for differences between stratification schemes). | | |
| Auxiliary information | Available auxiliary information is incorporated through GREG-estimation. Might reduce design-variance and partially corrects for selective non-response. | | | Availability of auxiliary information is crucial for most of the SAE-procedures that rely on models to borrow strength over space (e.g.: unit and area level models). Time series models, used to borrow strength over time, are less dependent of auxiliary information. | | |
| Sample design features | Estimation procedures fully based on the features of the probability sample through first and second order inclusion probabilities. | | | Area level models and time series models account for the sample design, since designbased estimators are used as the input variables for these models. Unit level models require additional explicit modelling of the sample design features. | | |
| Model-misspecification | GREG-estimators are robust for model misspecification in case of sufficiently large sample sizes. Model-misspecification doesn't compromise design consistency but might only result in an increased design variance. | | | Sensitive for model-misspecification, since this might result in biased estimates. | | |
| Computational effort | Relatively minor, since one one set of weights is derived to estimate all target variables. This property also enforces consistency between the marginal totals of different publication tables. Convenient to produce timely releases in a regular statistical production environment. | | | Large computational efforts, since separate models must be derived for separate target variables. Careful model selection and evaluation is required to avoid model-misspecification. Additional adjustments might be required to achieve enforce consistency between the marginal totals of different publication tables. | | |
| Assumptions | GREG-estimator is based on very mild assumptions. It is assumed that, conditionally on the specified weighting scheme, the non-response is not selective. | | | Stronger model assumptions, since SAE-procedures explicitly rely on statistical models to borrow strength over time or space. | | |

*Table B.2:* Overview properties of backcasting procedures

| Property | Micro approach | | Macro approach |
| --- | --- | --- | --- |
| | Design-based estimation | Model-based estimation | |
| Manual manipulation | Intense, due to double coding of the sample and the BR for each time period | | Double coding of the sample and the BR only required for the moment of the change-over. To construct time depending conversion matrices, double coding is required for a limited amount of preceding time periods. |
| Estimator | Reliable estimator for domains with large sample sizes, unreliable for weak domains. Not available for empty domains. | Improvement of the reliability for weak domains through the application of SAE procedures. | Stable estimates. Reliability, however, will in general be affected in a negative way due to the application of naive synthetic models. |
| Variance | Large design-variances for weak domains. | Improvement of the design-variance of the weak domains. | Variance approximations are possible, but do not reflect the bias that is introduced by using highly synthetic estimators. Time independent conversion matrices that operate on a relatively high aggregation level will generally result in stable estimates, at the cost of increased bias (see also row about bias). |
| Bias | Approximately design-unbiased. | Size of design-bias depends on the quality of the selected model and availability of auxiliary information. | High risk of the introduction large bias. Depends on the available auxiliary information, level of detail of the specified conversion matrices and the available information how they evolve over time. The structure of the economy is generally better retained with good auxiliary information, conversion matrices that are specified on a detailed aggregation level and with more realistic time dependent conversion matrices. |
| Auxiliary information | Incorporated through GREGestimation, see Table B.1. | Strongly determines the reliability of the estimates, see Table B.1. | Incorporated through the conversion matrices. Quality, amount of detail and time dependency strongly determines the reliability of the outcomes, see rows above. |
| Model-misspecification | Robust for model misspecification, see Table B.1. | Sensitive for model-misspecification, see Table B.1. | Extremely sensitive for model-misspecification, see rows above. |
| Computational effort | Relatively minor, since one set of weights is required for each time period, see Table B.1. | Large computational efforts, due to model selection and evaluation for separate variables and time periods, see Table B.1. | Minor. The main advantage of this approach is that series can be backcasted in a relatively straightforward way. |
| Assumptions | GREG-estimator is based on very mild assumptions, see Table B.1 | Stronger model assumptions, see Table B.1 | Very strong model assumptions that are generally hard to evaluate. |