This module is part of the

# Memobust Handbook

on Methodology of Modern Business Statistics

26 March 2014

# Theme: Design of Estimation – Some Practical Issues

**Contents**

# General section

## 1.    Summary

This module discusses a set of issues to consider in the design of the estimation of a business survey with emphasis on practical problems rather than on theory. Decisions have to be made together with other design decisions. The issues discussed mostly have neither obvious nor simple solutions. There is some theoretical ground to use, though, together with practical experience. The design decisions should be reasonable from these perspectives. The estimation procedure needs to work in practice also with a pressure on timeliness. It is important to save information about the procedure, including its weak points, in order to be able to improve successively.

There are different types of statistical inference depending on the combination of data and models used. The practical problems include, for instance, skewed distributions, outliers, coverage deficiencies, non-response, organisational changes, early estimates, small domains of estimation, and handling of administrative data.

The requests on quality include, for example, accuracy, coherence, and timeliness. Most business surveys are used both for primary statistics and as input to the National Accounts. In order to achieve coherence similar methods must be used in all business surveys.

Several types of statistical units are used in business surveys. Here the term 'enterprise' is used in most of the cases discussed, out of convenience.

## 2.    General description

Some typical characteristics for a business survey are described in Section 2.1; see also Rivière (2002). Some of the corresponding design considerations are discussed in Sections 2.2–2.13 with emphasis on practical issues. Finally, there is a section on evaluation.

### 2.1    The setting of a business survey and some typical issues

There is a statistical Business Register (BR) providing frames and information about different unit types used for business statistics, see, for instance, the handbook module "Statistical Registers and Frames – Survey Frames for Business Surveys" and other modules in that topic.

The distribution of most continuous variables (like turnover) is strongly skewed. A fairly small or limited number of large enterprises account for a noticeable or considerable portion of the total turnover, the total number of employees etc. See, for instance, figures given by Assoulin (2009, p. 2) and by the Blue-Ets-project (2013, p. 39). It is also typical that enterprises change quickly due to reorganisations, births, and deaths.

It is frequently the case that a statistical office has a special unit or group devoted to data collection from the largest enterprises, a group that keeps up-to-date with the organisational changes of these enterprises, serves as point of contact between each enterprise and many surveys, simplifies their data provision etc. See, for instance, the handbook modules "Data Collection – Design of Data Collection Part 2: Contact Strategies" and "Data Collection – Mixed Mode Data Collection" (about tailoring).

There is an increased interest in reducing direct data collection and in using other, already existing, data, typically administrative data such as tax data; see the handbook module "Data Collection –

Collection and Use of Secondary Data". There are both positive and negative implications of such use: lower costs, lower response burden, and a data set that is close to a census, but also a dependence on the alternative source with its unit type(s), population(s), variables, and reference times. The estimation possibilities become both greater and more restricted than with a sample survey with direct data collection. For instance, estimation for small areas and other small domains may become possible. On the other hand, it may not be possible to influence variables, editing, and production time.

There are some different types of business statistics and surveys in the European Statistical System, see, for instance, the handbook module "General Observations – Different Types of Surveys". There are annual statistics with a high degree of detail, such as the Structural Business Statistics (SBS). There are short-term statistics with a high pressure on timeliness and often with focus on changes over time, such as the STS: Short-Term business Statistics on industry, construction, retail trade and other services. There are further, secondary, statistics like the National Accounts (NA), which build on the primary business statistics just mentioned. There are many requests on the output quality of the primary statistics, stemming from European Regulations, national needs, the NA etc. Different needs often have to be balanced. The pressure on timeliness for short-term statistics leads to a system where there are preliminary statistics, revisions, and final statistics. See, for instance, the handbook module "Repeated Surveys – Repeated Surveys".

## *2.2    Estimation principles*

A statistical estimation procedure starts with some data and arrives at a set of statistics; to do so it uses some principle for the statistical inference. The principle usually includes an element of randomness. Statistical offices have a fairly long tradition of drawing a random sample, collecting data, and making inference to a finite population. Selection probabilities are then an essential ingredient. This is called a design-based method or inference.

The design-based principle may be extended to model-assisted estimation, where auxiliary information is used. Improvement of the accuracy is a major aim. Some form of assumption – implicit or explicit – is included. There may, for instance, be a model for the survey variable(s) in terms of the auxiliary variables. Non-sampling errors, like non-response, may be included in the modelling approach. Still, the sampling design is the foundation for the estimation.

Another type of principle uses a statistical model with its assumptions as the basis for the statistical inference. This is a model-based method. It is sometimes called a prediction approach. The model plays the major role in the inference, and the sampling procedure (if any) is less important. Model-based methods use, for instance, an assumption about a non-sampled part of the population or about relationships between variables or over time. They may be preferred, adequate, or even necessary to use. Small area estimation, early estimates, and non-random samples provide examples.

In general, no or negligible bias and a small variance are the basic demands for the choice of an estimator, possibly summarised in terms of a minimum or small mean squared error (MSE). Obviously, the type of statistical inference must be stated first, in order for bias, variance etc. to have a meaning. Design-based, model-assisted, and model-based methods use different approaches to randomness, and they include somewhat different information. The estimators are different, and so are their properties, since they depend on different assumptions about random elements.

If administrative data are used, some modelling and adjustments may be necessary, for instance to derive or model statistical variables based on the administrative variable definitions and to go from administrative units to statistical units. There are not (yet) generally agreed inference principles for this type of data. It is essential to state clearly the target for the inference, for instance the target population. This is not always done. There may be coverage deficiencies, missing data, and other non-sampling errors to handle in the estimation. It is important to be aware of this and transparent.

The estimation procedure has to be designed together with other parts of the survey design. Important enablers and issues to consider are the BR, further accessible data sources, data collection, sampling method (if any) and questionnaire design (when relevant), accuracy, and other output quality requests. The targets of the statistical output and the statistical inference to be used need to be stated; see, for example, the handbook modules "Overall design – Overall Design" and "Weighting and Estimation – Main Module". There is a lot of literature on estimation principles. A recent description is given by Valliant (2013), who comments on a summary report on non-probability sampling made by Baker, Brick, Bates, Battaglia, Couper, Dever, Gile, and Tourangeau (2013). There are further comments and a rejoinder by Baker et al. (2013). Even if these discussions are more relevant for household statistics, there are many general comments, which are useful also for business statistics and as a starting point to further literature. See also Zhang (2012), who provides a description of register-based statistics.

*2.3*    *A few estimation characteristics*

The statistical output consists of many estimates; there are normally many tables, each of which has many table cells. Hence priorities have to be made among all estimators and accuracy requests expressed, more or less explicitly. There are balances to make, such as follows. Which is more important in the trade-offs below? This should be considered before choosing the estimator – and choosing the sampling method, including the sample allocation.

- The overall table "total", the margins of the table, or the table cells.
- The current level or the change since a previous period.
- The absolute level or the relative level.

There are always differences between the current reality and the information in the BR. This is often due to delays in reporting about births, deaths, and organisational changes. There may also be incorrect information in the BR, for instance about the industry/activity of the enterprise. It is obviously important for the BR to have good and frequently updated sources. The BR will still be a bit behind. The deficiencies and their implications are asymmetric for the survey. The current sample includes over-coverage, which may be possible to detect, whereas outside information is needed to find under-coverage. There are two main ways to obtain and use information that is more up-to-date than the frame information. Firstly, there may be a more recent register – an updated version of the BR or the frame, or a similar register – to use as support in the estimation procedure. Secondly, differences from the frame information may be observed in the survey. Such information often needs to be taken into account, after adequate checking. The cases with negligible or small effects are fairly rare.

Size is mostly an essential factor in the survey design. It is essential to obtain data of good quality from the large enterprises. Small enterprises that have grown in comparison with the information used in the design may be disturbingly influential if straight-forward estimation formulas are used. Hence, care should be taken to have recent and adequate measures of size.

*2.4    Some comments on sampling design*

The sampling and estimation methods should be chosen together with regard to the targets of the survey and the priorities made. Choudhry, Rao, and Hidiroglou (2012) show how the sample can be allocated to satisfy accuracy requests when estimating a set of sub-population means; this is just one example to illustrate the strong link.

The reference time of the target population usually agrees with the reference time of the survey variables, for instance in a monthly survey. If estimates of change are important, it is normally favourable to include the same unit repeatedly in the sample. However, some exchange is necessary, both to include changes in the target population and with regard to response burden for the enterprises selected.

Selecting a new sample from an updated frame leads to more recent information about the population than continuing with the old sample from the old frame. On the other hand it means that enterprises may move in and out of the sample, which is inconvenient for such enterprises and an addition to their response burdens. It means additional work also for the statistical office. The frequency of updates of frames and samples is discussed for example in the handbook modules "Statistical Registers and Frames – Survey Frames for Business Surveys" and "Repeated Surveys – Repeated Surveys". Many countries renew or complement their samples fairly frequently. An addition is typically selected from population parts that would otherwise be under-coverage.

Rather many countries use some form of sample co-ordination, at least over time, but possibly also between surveys. See the handbook module "Sample Selection – Sample Co-ordination". There are two main reasons, as already indicated. Firstly, accuracy is improved through positive co-ordination of samples (which means more overlap than expected in the case with independent samples; typically used over time for each survey and between related surveys at the same time). Secondly, response burden is spread over enterprises through negative co-ordination (less overlap; typically used between unrelated surveys at the same time). Time aspects and a full set of surveys are taken into account. The principal aim is that an enterprise is selected for one or possibly a few sample surveys during a period and then, ideally, has "survey holidays". The success depends on the number of surveys, the number of enterprises to select from, and the methodology chosen. The response burden can, of course, not be spread unless there are a (fairly) large number of enterprises. The positive co-ordination over time is limited with regard to response burden, for instance expressed as an enterprise being selected a targeted number of years.

There is some similarity between a sample survey with positive co-ordination over time and a panel survey. However, the latter has a pre-determined pattern of overlap, whereas the former contains random elements. Therefore, the term panel is normally not used in systems with sample co-ordination.

With panels it may be necessary to consider both cross-sectional and longitudinal weights; see, for instance, the handbook module "Weighting and Estimation – Main Module". As described in the handbook module "Sample Selection – Sample Co-ordination" there are techniques to co-ordinate samples such that each sample can be considered as random, which make the estimation design straight-forward. There is, however, dependence between samples over time, which has to be taken into account in the estimation of variance for estimators of change. See also the handbook module "Repeated Surveys – Repeated Surveys" for a discussion and a few references.

## 2.5 Estimation with non-response and skewed distributions

There are two main methods to handle non-response in the estimation. It is frequently the case that imputation is used for item non-response; see the handbook module "Imputation – Main Module". Similarly, reweighting is often used for unit non-response; see the handbook module "Weighting and Estimation – Main Module", where response propensity is described and furthermore weight adjustments like calibration. This approach with imputation and reweighting is a reasonable starting point for the design of estimation. However, the nature of the distribution of the target variables must also be considered. This is a general aspect, which is necessary to take into account also when it comes to non-response.

For continuous variables which are highly skewed, say turnover, it is of great importance to aim at no unit non-response for the largest enterprises. If there still is unit non-response among the largest enterprises, different treatments of the non-response for different segments of the target population should be considered. The design choice may be a split into two major methods. Non-response among the smaller enterprises is handled by reweighting, which is with or without auxiliary information. Since the largest enterprises normally are unique, it may be wise to consider imputation for them. It needs to be a careful imputation, which builds on information for the particular enterprise, such as previous values. It is often manual. See also the handbook methodological module "Statistical Data Editing – Manual Editing".

It should be observed that non-response may be related to over-coverage or to organisational changes, as further discussed in the next sections.

## 2.6 Estimation with over- and under-coverage

A possibility to reduce coverage problems – in the estimation procedure rather than in the sampling procedure – is to use auxiliary information from an updated frame or from a similar register when building the estimator. Calibration against updated auxiliary information is a way to handle the problems with over-coverage and under-coverage, as far as the new version of auxiliary information goes.

There are, however, some practical drawbacks with the calibration method. Unexpected results may be difficult to penetrate and resolve, since the estimator is rather complicated. Moreover, with quick production rounds in short-term statistics, the difficulties with matching microdata should not be underestimated.

Hence, under-coverage may be reduced through the sampling or estimation procedure or both. If neither is used in short-term statistics and if over-coverage is set to zero, the level of the estimates will decrease gradually as time goes and the effects of over-coverage increase. This may be balanced through some model assumption. ONS (2013, p. 6) and ONS (2012, pp. 41–44) are two examples showing how design weights are adjusted to unit non-response and to births and deaths with an assumption about the births-to-deaths ratio. The basic assumption used in these cases is that the ratio is equal to 1. Large enterprises, which are included with probability one, are a natural exception. A further possibility is to analyse birth and deaths rates and make some extrapolation over time. BLS (2013, p. 37) illustrates variation over time for some variables, and the estimation methodology used for a large employment survey is described in these technical notes.

The coverage issues just discussed are relevant not only for the population as a whole but also for sub-populations, for instance industries/activities within the target population. An erroneous or changed NACE code may be detected in the sample. Again there is a design choice whether to use the "original" NACE code in the frame or to use the correct one. Several issues should be considered, including bias and variance of the possible estimators. The enterprise may have both principal and secondary activities. To "move" such an enterprise between domains of estimation may exaggerate the changes that have occurred.

## 2.7 *Using auxiliary information in both sample design and estimation*

Auxiliary information may be used both in the sample design, for instance in a stratification, and in the estimation, for example through calibration. The auxiliary information should in both cases be highly correlated with the target variable(s).

It is possible to use the same variables as auxiliary information in both stratification and calibration. As an example, consider the situation using the number of employees as stratification variable. The variable is categorised into some classes, and this categorised variable is used as a stratification variable. However, the original variable still contains information, which can be used in calibration in order to reduce non-response bias.

The choice of which variable(s) to use in the sampling and which variable(s) to use in the estimation – the same or different variables – is based on an analysis of the possible auxiliary variables and their correlations with important target variables.

Calibration means that certain estimates are consistent with known totals. This may be a further aim of the calibration; to increase coherence and consistency. It may make it easier for a user to combine several sets of statistics.

## 2.8 *Handling organisational changes in a sample*

### 2.8.1 *General discussion*

The data collection will show that the sampled enterprises differ somewhat from the frame information, and the changes will most likely increase with time. These differences must be handled. It is especially important when using an estimator with no auxiliary information and an "aging" sample. Then other methods than calibration should be used in order to overcome coverage issues. Some such situations are described below together with possible assumptions and methods. There is little literature about this type of methodology and no established practice.

If, for instance, two enterprises merge after the time of the frame information, the resulting enterprise could not have been sampled directly, but it may be sampled via either of the original enterprises. This situation may be regarded as indirect sampling, and it may be handled as such. The weight construction depends in general on the sampling method and the possible routes from the frame and sampling units to the target and collection units. This methodology is used when the unit types are deliberately different. See, for instance, Lavallé and Labelle-Blanchet (2013).

However, such estimation may become quite complex. Here, the sampling units and the target units do not differ systematically but due to certain events, which affect a limited number of units. A simplified, pragmatic, approach may then be considered and possibly chosen. This may be rational if

the causing type of event is fairly rare. For instance, with stratified random sampling, the ordinary point and variance estimators are quite simple, and it may be tempting to essentially keep these estimators. There are further options, such as weight-sharing, as just mentioned above, Lavallé and Labelle-Blanchet (2013). This may be natural, for instance when working with panels. Knottnerus (2011) studies estimators for totals and growth when panels are used, including situations where businesses change, merge, or demerge/split.

Black (2001) discusses different ways to adjust weights and handle changes in numbers of units. He notes that making changes can require a significant amount of processing and possibly create revisions that are unnecessary. With small net effects it may be better not to make changes or to postpone them.

Please observe that the descriptions below are meant as food for thought in the case of an aging sample. In a case where an updated register or frame is used for calibration, there is some but often not full information about changing units. For each unit the current situation may or may not have reached that register or frame. The term enterprise is again used here in a generic sense, out of convenience. From an estimation point of view activities are interesting, not identification numbers and other administrative information. The latter are important as such, for instance for contacts, for comparisons with later versions of registers and frames, including calibration, and also when using administrative data. Lindblom and Nordberg (2004) describe and discuss birth and death rates and also calibration.

### 2.8.2   Births, deaths, and re-constructions

When information about completely new large enterprises is found in media or elsewhere, then such enterprises can be put in a special stratum or group, with design weight one. This method should only apply to very large, new, enterprises. It is an enterprise that almost certainly would have been included in a "take-all-stratum" had it been known at the sampling occasion. For births of smaller enterprises, see Section 2.6 above.

When an enterprise ceases to exist, a common method is to code it as over-coverage and set its variables values to zero. Using this method in a successive set of estimates from the same sample without compensation for under-coverage means a tendency with decreasing estimates, as more and more enterprises die. Another possibility is to code and treat the dead enterprise as non-response. With simple reweighting (no calibration) this corresponds to an assumption that there is a birth rate equal to the death rate; compare the examples given in Section 2.6 above. This is a somewhat dangerous assumption, especially when there are quick movements in the economy (the business cycle with up-and-down movements in economic activity).

If an enterprise is reconstructed – in the sense that it has new owner(s) but is otherwise unaffected – it should be treated as unchanged in the estimation procedure. It is only identification number(s) and some administrative information that change. This is simply a continuity rule.

### 2.8.3   Mergers and splits

If two enterprises merge, it may be a reasonable assumption for the estimation that the larger enterprise has "taken over" (bought) the smaller one: say that A buys B resulting in C. With this simplifying assumption A has changed into C and B no longer exists. The estimation procedure then handles the units involved so. If A belongs to the sample, then C (the new A) is surveyed. If B belongs to the sample, then B is coded as over-coverage with value zero. Hence, the resulting unit C has the

original sampling probability of unit A (belongs to that stratum in case of a stratified random sample). This is not feasible if C is large enough to be an outlier (Section 2.9).

If, instead, an enterprise splits into two (or more) enterprises, this may be considered as a situation where both (all) enterprises are tied to the original sampling unit and its weight. In a simplified approach, as above for mergers, the new enterprises are kept in the sample and surveyed. It may be reasonable to combine the data collected from the parts into data that correspond to the original enterprise. This could be the case with simple variables, like turnover. However, there may be complicating facts with information showing, for instance, that the parts belong to different domains of estimation. The influence on estimates has to be considered. With fairly large enterprises some tailor-made solution is needed.

See also the handbook module "Weighting and Estimation – Main Module", which recommends a strict approach with indirect sampling and weight-sharing, and Lavallé and Labelle-Blanchet (2013).

## 2.9    Outliers in the sample

An outlier is a value that deviates considerably from the "bulk" of observations. It may be due to the skewed distribution, to the current size of the unit differing from the frame information, or to an erroneous value. Especially the last category should be handled and eliminated in the statistical editing. The large enterprises, which are included with probability one, are considered on their own; compare Section 2.5 above. The remaining outliers are usually treated as representative. They may still be too influential in a simple estimation procedure, depending on their weights. There are methods to handle such outliers, for instance by winsorisation of the value or by weight modification. Normally, the variance is decreased by using such methods, but a bias is introduced at the same time. The handbook module "Weighting and Estimation – Main Module" gives an overview of methods, and the method module "Weighting and Estimation – Outlier Treatment" provides more detailed information. There is a recent article by Beaumont, Haziza, and Ruiz-Gazen (2013).

The design should (i) try to prevent or at least reduce the occurrence of outliers and (ii) choose an appropriate method for handling the outliers that still occur. An appropriate and up-to-date (as far as possible) measure of size is essential for the first aim and well worth a study during the design.

For some variables, which are strongly skewed, there will be outliers due to their nature, like investment in buildings, which for a small enterprise may be high at rare occasions. The estimation procedure should foresee such outliers in the design stage and choose an appropriate method to handle them. Especially in a repeated survey there is information about outliers to first collect and then utilise in later production rounds. Another possibility is to use information from annual surveys to make rough estimates for short-term surveys about occurrences and then choose appropriate method(s). There may in both situations be differences, though, between different parts of a business cycle (up-and-down movements in economic activity). Some care may be needed with different levels of detail, for instance for domains of estimation. Higher levels of aggregates are less sensitive than lower levels.

The choice of method has to consider (i) the information needed and available for the estimator being robust and also (ii) the complexity which is added in comparison with the "basic" estimator that would otherwise be used, if there were no outliers.

Some surveys use a technique with a "surprise-stratum", which means that they put found outliers there with weight one. If the outlier really can be considered unique, this handling is reasonable.

Otherwise – when the outlier is representative – such a procedure introduces a bias, and it is not recommended in comparison with the methods described.

The two previously used examples from ONS are different. ONS (2013, p. 6) is a survey using winsorisation. There scaling is used to achieve consistency when one winsorised variable is the sum of two others. ONS (2012, pp. 42–44) describes how outliers are first identified and then treated as non-representative, using also a post-stratification. There are also model-based methods in use; see the handbook module "Weighting and Estimation – Outlier Treatment" for a description and references.

*2.10    Cut-off sampling*

As indicated in the handbook module "Sample Selection – Main Module", it may be motivated not to include a certain part of the population in the sample. Benedetti, Bee, and Espa (2010) provide a framework for such sampling and estimation. They show, for instance, how the estimator can be constructed for the target population, which includes the unobserved population below the cut-off threshold. A model needs to be used, such as an assumption that the share below the threshold of the total of a survey variable is the same as for an auxiliary variable known, for instance, from the frame.

In this case the estimation design determines important parameters, such as the threshold value, the auxiliary variable, and the model assumption to be used. Again, the choice of the size variable is important to get a reasonable model and estimator. Later on, when the BR has been updated, it may be possible to assess the estimator, at least partly.

There is some similarity between this estimation procedure and the use of administrative data for small enterprises, see Section 2.13.

*2.11    Models: early estimates*

A preliminary (early) estimate may be required already at a time when the response rate normally is fairly low. This affects not only the variance (a smaller number of respondents), but also the potential bias, due to the early respondents possibly being different from later ones. The handbook module "Weighting and Estimation – Main Module" describes three different possibilities: to take a random sub-sample with a short response time, to use a design-based estimator based on early respondents, and to use a model-based estimator. The estimators can be compared in terms of assessed possible biases and estimated variances, and – perhaps most important – revision sizes. Again, there may be differences between different parts of a business cycle with its up-and-down movements.

The outcomes of comparisons either when planning (if possible and meaningful) or later – when revisions can be computed – are useful ingredients when choosing an appropriate estimator. If it is a repeated survey, information should be collected, analysed, and used for improvements of the design.

Similar situations may occur when using administrative data, see Section 2.13.

*2.12    Models: small area estimation*

If there are small domains of estimation, design-based estimators will have a large variance, so model-based estimators are often preferred or even necessary. This topic is described in rather much detail in the handbook module "Weighting and Estimation – Main Module" with references to further modules. Hence, there is no specific discussion here. It is important to note, though, that statistical disclosure control may come into play, see handbook module "Statistical Disclosure Control – Main Module".

## 2.13    Use of administrative and other accessible data

The handbook module "Weighting and Estimation – Main Module" has a short section on integration of administrative data, and there are special handbook modules on administrative and secondary data in the topics "Data Collection", "Statistical Data Editing", and "Weighting and Estimation". Especially "Weighting and Estimation – Estimation with Administrative Data" is relevant here.

There are several important issues to consider in the design when considering the possible use of administrative data:

- The administrative and the statistical units.

- The coverage of the administrative source(s) in comparison with the target population.

- Variable definitions and how to derive or model the target variables from the administrative variables.

- Measurement errors and other non-sampling errors (deficiencies).

- Timeliness, which may depend on the size of the unit. (Small legal units may not report frequently to tax authorities, for instance.)

It may be a good solution to use direct data collection for large enterprises and administrative data for medium-sized and small enterprises. This eliminates the response burden for the latter group. There may be problems, though, with some variables and with timeliness. Models need to be introduced, studied, and chosen with respect to quality aspects such as accuracy, coherence, and timeliness.

There is a large set of deliverables from an ESSnet project, see ESSnet on AdminData (2013) with different aims, such as data quality and estimation in different situations, depending on variables and timeliness, for instance. Lewis (2012), Paavilainen (2012), and Brinkley, Preston, and Scott (2012) are three different examples of the use of administrative data in surveys, showing also difficulties in practice and how model assumptions can be included. The two first examples are from the ESSnet project. Lewis (2012) illustrates how to find solutions when variables are not directly available from administrative sources. Paavilainen (2012) describe that administrative data may change, that some of them may be late for short-term statistics with short production time, and how an index can be constructed with mixed sources. Brinkley et al. (2012) is a broad description of the inclusion of administrative data with emphasis on sample design and estimation.

## 2.14    Evaluation of the design

The design of the estimation procedure means work with requests on quality, in particular accuracy. At the end of the production round the achieved quality (for instance accuracy) should be studied and compared with the planned quality (including accuracy). If there are differences, the causes should be analysed and possible actions should be considered. This is particularly the case in a repeated survey, for possible improvements in later production rounds. Experiences may also be shared across surveys.

If changes of the design are motivated, not only the estimator should be considered. It may be easier and better to modify the allocation of the sample, just as a simple example. There may be many further issues to consider, for instance how to handle non-sampling errors in the estimation and in other parts of the production process. See, for example, the handbook modules "Repeated Surveys – Repeated Surveys" and "Overall Design – Overall Design".

**3.     Design issues**


**4.     Available software tools**


**5.     Decision tree of methods**


**6.     Glossary**

For definitions of terms used in this module, please refer to the separate "Glossary" provided as part of the handbook.

**7.     References**

Assoulin, D. (2009), Choosing an Imputation Method for Large Firms. European Establishment Statistics Workshop, Efficient Methodology for Producing High Quality Establishment Statistics (Stockholm, Sweden).

Baker, R., Brick, M., Bates, N., Battaglia, M., Couper, M., Dever, J., Gile, K., and Tourangeau, R. (2013), Summary Report of the AAPOR Task Force on Non-probability Sampling. *Journal of Survey Statistics and Methodology* **1**, 90–105 (and Rejoinder 137–143).

Beaumont, J.-F., Haziza, D., and Ruiz-Gazen, A. (2013), A unified approach to robust estimation in finite population sampling. *Biometrika* **100**, 555–569.

Benedetti, R., Bee, M., and Espa, G. (2010), A Framework for Cut-off Sampling in Business Survey Design. *Journal of Official Statistics* **26**, 651–671.

Black, J. (2001), Changes in Sampling Units in Surveys of Businesses. Federal Committee on Statistical Methodology (FCSM), 2001 FCSM Conference Papers.

BLS (2013), Technical Notes for the Current Employment Statistics Survey. US Bureau of Labor Statistics. (Link from 24-01-2014: http://www.bls.gov/ces/cestn.pdf.)

Blue-Ets-project (2013), *Best practice recommendations on variance estimation and small area estimation in business surveys*. Deliverable 6.2. BLUE-Enterprise and Trade Statistics. (Link from 24-01-2014: http://www.blue-ets.istat.it/fileadmin/deliverables/Deliverable6.2.pdf.)

Brinkley, E., Preston, J., and Scott, A. (2012), Using Administrative Taxation Data to Improve Sample Design and Estimation - An ABS Perspective. Paper presented at the International Conference on Establishments, ICES-IV.

Choudhry, G. H., Rao, J. N. K., and Hidiroglou, M. A. (2012), On sample allocation for efficient domain estimation. *Survey Methodology* **38**, 23–29.

ESSnet on AdminData (2013), *Use of Administrative and Accounts Data for Business Statistics*. Project in three time periods (SGAs 1-3). Deliverables are accessible. (Link from 24-01-2014: http://www.cros-portal.eu/content/use-administrative-and-accounts-data-business-statistics.)

Knottnerus, P. (2011), *Panels – Business Panels*. Statistical Methods (201119), Statistics Netherlands, The Hague/Heerlen.

Lavallé, P. and Labelle-Blanchet, S. (2013), Indirect sampling applied to skewed populations. *Survey Methodology* **39**, 183–215.

Lewis, D. (2012), Methods for using administrative data to estimate survey variables not directly available from administrative sources. Paper presented at the International Conference on Establishments, ICES-IV.

Lindblom, A. and Nordberg, L. (2004), On adjustment for coverage problems in short-term business surveys. Paper contributed to the European Conference on Quality and Methodology in Official Statistics (Q 2004).

ONS (2012), *Annual Business Survey (ABS): Technical Report*. Issued by Office for National Statistics, August 2012.

ONS (2013), *Quality and Methodology Information on Business Register and Employment Survey (BRES)*. Information paper. Issued by Office for National Statistics, 20th November 2013.

Paavilainen, P. (2012), Efficient use of administrative data in the production of economic statistics in Finland. Paper presented at the International Conference on Establishments, ICES-IV.

Rivière, P. (2002), What Makes Business Statistics Special? *International Statistical Review* **70**, 145–159.

Valliant, R. (2013), Comment (on Baker et al., 2013). *Journal of Survey Statistics and Methodology* **1**, 105–111.

Zhang, L.-C. (2012), Topics of statistical theory for register-based statistics and data integration. *Statistica Neerlandica* **66**, 41–63.

# Interconnections with other modules

**8.      Related themes described in other modules**

1. General Observations – Different Types of Surveys

2. Overall Design – Overall Design

3. Repeated Surveys – Repeated Surveys

4. Statistical Registers and Frames – Survey Frames for Business Surveys

5. Sample Selection – Main Module

6. Sample Selection – Sample Co-ordination

7. Data Collection – Main Module

8. Data Collection – Design of Data Collection Part 2: Contact Strategies

9. Data Collection – Mixed Mode Data Collection

10. Data Collection – Collection and Use of Secondary Data

11. Statistical Data Editing – Main Module

12. Imputation – Main Module

13. Weighting and Estimation – Main Module

14. Weighting and Estimation – Estimation with Administrative Data

15. Statistical Disclosure Control – Main Module

**9.      Methods explicitly referred to in this module**

1. Statistical Data Editing – Manual Editing

2. Weighting and Estimation – Outlier Treatment

**10.     Mathematical techniques explicitly referred to in this module**

1.

**11.     GSBPM phases explicitly referred to in this module**

1.

**12.     Tools explicitly referred to in this module**

1.

**13.     Process steps explicitly referred to in this module**

1.

# Administrative section

**14.      Module code**

Weighting and Estimation-T-Design of Estimation

**15.      Version history**

| Version | Date | Description of changes | Author | Institute |
|---------|------|------------------------|--------|-----------|
| 0.0.1 | 16-12-2013 | first plan | Marianne Ängsved and Eva Elvers | Statistics Sweden |
| 0.0.2 | 03-01-2014 | some sections drafted | Eva Elvers | Statistics Sweden |
| 0.0.3 | 06-01-2014 | further sections | Marianne Ängsved | Statistics Sweden |
| 0.0.4 | 08-01-2014 | merging, expanding | MÄ+EE | Statistics Sweden |
| 0.0.5 | 10-01-2014 | further additions | EE+MÄ | Statistics Sweden |
| 0.0.6 | 25-01-2014 | reviews from HU, CH, IT | EE+MÄ | Statistics Sweden |
| 0.0.7 | 28-01-2014 | continued improvements | MÄ+EE | Statistics Sweden |
| 0.1.0 | 07-02-2014 | reviews EB, IT | EE+MÄ | Statistics Sweden |
| 0.1.1 | 12-02-2014 | preliminary release | | |
| 1.0 | 26-03-2014 | final version within the Memobust project | | |

**16.      Template version and print date**

| Template version used | 1.0 p 4 d.d. 22-11-2012 |
|-----------------------|-------------------------|
| Print date | 26-3-2014 13:31 |