

# Estimation of the main variables of the economic account of small and medium enterprises based on administrative sources

Marco Di Zio<sup>1</sup>Ugo Guarnera<sup>1</sup>Roberta Varriale<sup>1</sup>

## Abstract

*The paper describes the imputation procedure of the main variables of small and medium-sized enterprise balance sheet. The procedure is used as part of the project aimed at creating an integrated system for the production of detailed estimates on enterprise economic performance. The variables are imputed using mainly administrative sources as Financial Statements, Studi di Settore and Tax return. The proposed procedure represents an integrated set of different imputation techniques: Predictive Mean Matching, nearest neighbor donor, and a two-step procedure for the treatment of variables characterised by a high presence of zeros. A first evaluation of the procedure is carried out by comparing the estimates based on administrative data with those obtained by the use of sampling weights.*

**Keywords:** Imputation, predictive mean matching, nearest neighbor donor

## Sommario

*Il lavoro descrive la procedura di imputazione delle principali variabili del conto economico delle piccole e medie imprese. La procedura è stata utilizzata nell'ambito del progetto finalizzato alla realizzazione di un sistema integrato per la produzione di stime dettagliate sui risultati economici delle imprese. Le variabili vengono ricostruite utilizzando principalmente le fonti amministrative Bilanci delle Società di Capitale, Studi di Settore e modello Unico. La procedura proposta è un insieme integrato di diverse tecniche di imputazione: Predictive Mean Matching, donatore di minima distanza, ed una procedura a due passi per il trattamento delle variabili caratterizzate da una elevata presenza di zeri. Una prima valutazione della procedura è stata ottenuta confrontando le stime basate su dati amministrativi con quelle ottenute mediante l'utilizzo dei pesi campionari.*

**Parole Chiave:** Imputazione, predictive mean matching, donatore di minima distanza

<sup>1</sup> Istat, Directorate for methodology and statistical process design. email: dizio@istat.it, guarnera@istat.it, varriale@istat.it. The views expressed in this paper are solely those of the authors and do not involve the responsibility of Istat

## 1. INTRODUCTION

In Istat, Structural Business Statistics (SBS) for small and medium enterprises (SME) are traditionally based on sample surveys. In the last years, the increasing availability of information from administrative sources made it possible to take into account the possibility of using administrative data to improve the quality of the produced statistics. Until now this information has been generally used as auxiliary information to treat non-response in survey data and to calibrate the estimates on known aggregates.

The level of maturity in the analysis of these kinds of data lead Istat to use administrative data as a primary source for information to produce SBS statistics. In 2011 for the first time, data from administrative sources as *Financial Statement*, *Studi di settore*, *Tax Return* are used to build a microdata file composed of the main economic variables. The choice of producing a microdata file follows from the difficulty of providing coherent estimates at different level of aggregation, in this regard we remind that these data are also used by National Accounts to build national economic aggregates (Istat, 2014).

Since not all the variables are available in all the data sources, and the sources cover only subsets of the target population, the microdata file is a result of an imputation process. The imputation procedure is based on a combination of different techniques that are introduced to comply with requirements given by constraints, such as statistical relationships among main variables, balance edits, and presence of zero-inflated variables.

Given such a complexity, the assessment of the procedure is not an easy task. A comparison with official estimates based on the SME sample survey data is carried out. The differences are decomposed in terms of sampling and measurement errors. The analysis of the impact of the different error sources may be useful to validate the results and to improve the process of production of statistics in this context.

The paper is structured as follows. Section 2 describes the informative context of SME statistics based on administrative data. The imputation process is described in Section 3, and some results about the evaluation of the estimation procedure are reported and discussed in Section 4.

## 2. Informative context

The administrative data sources are the Financial Statements, *Studi di settore*, and Tax Return data. The units of Financial statements (FS) are the companies, mainly corporate firms, liable to fill in the financial statement. The “*Studi di settore*” (SDS) is a Fiscal Authority survey that aims at evaluating the capacity of enterprises to produce income and at indirectly assessing whether they pay taxes correctly. The units compiling the SDS form, composed of detailed information on costs and income, are the enterprises with a turnover less than 7,500,000 Euros belonging to many activity sectors. Tax return data are mainly based on the fiscal form “*Unico*” and, for a residual part of units representing corporate firms, on “*Irap*” (the Italian regional tax on productive activities).

All the analyses described in the paper have, as a starting point, the quality assessment of the administrative data carried out by the subject matter experts (Curatolo *et al.*, 2015). Although in principle many variables observed in the administrative data sources could be used for the SBS estimates, only some of them can be considered enough reliable both in terms of consistency of definitions with the ones described by the SBS regulation, and in terms of reported values compared to the SME observations. The list of the variables used in the imputation process is reported in Table 1.

**Table 1 - Variables used in the imputation process**

Section	Label	Variable
Revenues	$Y_1$	Income from sales and services (Turnover)
	$Y_2$	Changes in stock of finished and semi-finished products
	$Y_3$	Changes in contract work in progress
	$Y_4$	Changes in internal work capitalized under fixed assets
	$Y_5$	Other income and earnings (neither financial, nor extraordinary)
Costs	$Y_6$	Purchases of goods
	$Y_7$	Purchases of services
	$Y_8$	Use of third party assets
	$Y_9$	Changes in stocks of raw materials and for resale
	$Y_{10}$	Other operating charges
	$PC$	Personnel Costs

It is worthwhile to remark that the variable *personnel costs* is always observed and it is used as auxiliary variable in the imputation procedure. In addition to  $PC$ , some derived variables are used as auxiliary variables in the imputation process. In fact in some cases they are considered more reliable than the variables used for the derivation, this is due to a kind of compensation process that is not easy to model. The derived variables, related to the Cost section, are listed in Table 2.

**Table 2 - Derived variables**

Derived Variable	Transformation	Variable description
CS	$Y_2 - Y_9$	Total Change in Stock
GS	$Y_6 + Y_7$	Purchases of Goods and Services
IC	$GS + Y_8 + Y_{10}$	Total Intermediate Costs

We remark that some variables are observed in more than one data sources, this means that for each of them (generally) different values are available. In this application a hierarchical approach is chosen. It consists in assigning a hierarchy to the administrative sources and consequently values of the variables are chosen according

to this raking. The hierarchy has been established by subject matter experts according to some quality criteria such as coverage of administrative sources with respect to the business register, steadiness of the supply in terms of timing and variable content (Curatolo *et al.*, 2015), and it assigns the first rank to FS, then to SDS, and finally to Tax Return data. Based on the assumed hierarchy, the coverage of the administrative sources for the 2011 is reported in Table 3.

**Table 3 - Coverage of administrative sources for the 2011**

Source	Frequency	Relative frequency (%)
FS	714885	16.1
SDS	2836100	64.0
Unico	714894	16.1
Irap	4201	0.1
NA	162848	3.7
Total	4432928	100.0

The subset of population not covered is small and it is composed of the smaller units in terms of size. In our procedure we also need to take into account that for all the units the number of employees is known from the business register ASIA and that for most of the units in ASIA there is an important information related to turnover (i.e., 'amount of business') and it is a good proxy for the turnover mainly coming from the VAT declarations.

It is however worthwhile to remark that the coverage of each single variable depends on its availability in the different data sources.

In Table 4 the pattern of missing data per variables and data sources is illustrated. The symbols 'X' and '?' stand for observed and missing data respectively. In this table, SDS-F, SDS-G and Unico1-Unico8 refer to the different kinds of SDS and Unico that enterprises have to fill in depending on their legal status. In particular, the units compiling SDS-G and Unico5 are represented by professionals and "minimum taxpayers", respectively.

**Table 4 - Pattern of missing data per variables and data sources**

Source	Y <sub>1</sub>	Y <sub>2</sub>	Y <sub>3</sub>	Y <sub>4</sub>	Y <sub>5</sub>	Y <sub>6</sub>	Y <sub>7</sub>	Y <sub>8</sub>	Y <sub>9</sub>	Y <sub>10</sub>	PC	CS	GS	IC	Coverage rate (%)
FS	X	X	X	X	X	X	X	X	X	X	X	X	X	X	16.13
SDS-F	X	?	X	X	X	X	X	X	?	X	X	X	X	X	50.08
SDS-G	X	X	X	X	X	?	?	?	X	X	X	X	?	X	13.90
Unico1	X	X	X	X	X	?	?	?	X	?	X	X	X	?	0.78
Unico2	X	X	X	X	X	?	?	X	X	?	X	X	X	?	0.04
Unico3	X	?	X	X	X	?	?	?	?	?	X	X	X	?	2.73
Unico4	X	?	X	X	X	X	X	?	?	?	X	X	X	?	0.76
Unico5	X	X	X	X	X	?	?	?	X	?	X	X	?	X	10.86
Unico6	X	?	?	?	?	?	?	?	?	?	X	?	?	?	0.16
Unico7	X	?	?	?	?	?	?	?	?	?	X	?	?	?	0.31
Unico8	X	?	?	?	?	?	?	?	?	?	X	?	?	?	0.49
Irap	X	X	X	X	X	X	X	X	X	X	X	X	X	X	0.09
NA	?	?	?	?	?	?	?	?	?	?	X	?	?	?	3.67

The rate of missing data per variable, taking also into account the information available in the Business Register (ASIA), is reported in Table 5. We notice that the minimum amount of missing data is related to the variable  $Y_1$  (turnover) and it is lower than the minimum of Table 3. This is because in the Business Register there is a variable closely related to the turnover (named *amount of business*) that in some cases can be used to predict  $Y_1$ . On the other hand, a very high rate of missing data (approximately 58%) affects variables  $Y_2$  and  $Y_9$ , associated with the two components of the change in stock. It is worthwhile to mention that however, for a large set of units, the difference  $CS = Y_2 - Y_9$  is known even though the separate components are not observed. This mitigates the impact of the missing data on the estimates of some crucial derived variables such as the *value added*, where only the total change in stock is relevant.

The economic data described in this paragraph are used both by the SBS sector and by National Accounts, requiring many domains of estimation. In order to avoid consistency problems, missing data are imputed to obtain a microdata file.

### 3. The imputation process

The imputation procedure is based on a combination of different techniques.

The entire imputation process is composed by 4 sequential steps:

1. deterministic imputation based on the guidelines of subject matter experts;
2. imputation of the variables  $Y_1$ ,  $Y_6$ ,  $Y_7$  and  $CS$ , through Predictive Mean Matching (PMM);
3. imputation of the variables  $Y_3$ ,  $Y_4$ ,  $Y_8$ ,  $Y_{10}$  and  $Y_5$ , through Nearest Neighbor Donor (NND);
4. imputation of the variables  $Y_9$ ,  $Y_2$  through a two-step procedure composed by a logistic and a linear regression model.

In this paper we focus on the description and evaluation of the imputation process related to the last three steps.

**Table 5 - Rate of missing data per variable**

$Y_1$	$Y_2$	$Y_3$	$Y_4$	$Y_5$	$Y_6$	$Y_7$	$Y_8$	$Y_9$	$Y_{10}$	$CS$	$GS$	$IC$
3.7	58.2	4.7	4.7	4.7	19.2	19.2	19.8	58.3	19.8	6.7	15.6	15.6

The pattern of missing data depicted in Table 4 is the one obtained after the deterministic imputation in step 1.

The steps from 2 to 4 have been carried out inside strata based on the economic divisions Nace2 cross-classified with the two subsets of observations characterised by having or not personnel costs. For the PMM, also the item  $GS$  has been used to define strata, the distinction is made between units either with or without purchases of goods and services.

The enterprises with information coming from SDS-G and Unico5 are represented by professionals and “minimum taxpayers”, and they are considered to behave

quite differently from the rest. Since the imputation process tends to reproduce in the non-observed part of the population the behaviour of the observed units, for this subset of population the imputation is made by resorting to the SME survey, details will be given later in the paper.

The choice of each imputation method for different groups of variables is due to: the percentage of missing values, the variable distribution characteristics (only positive, zero-inflated, etc.), the presence of a (weak/strong) relationship between variables and the presence of balance edits. All these characteristics influence the choice of a statistical model in the imputation process.

### 3.1 Methods

The PMM can be considered as a NND imputation technique based on a distance function where matching variables are weighted through their predictive power with respect to the variables that have to be imputed. In a multivariate context, the PMM is typically applied to match each recipient to the donor having the closest predictive mean with respect to a regression model of the target variables on a set of covariates. Selection of donors is based on the Mahalanobis distance defined in terms of the residual covariance matrix from the regression model. Intuitively, Mahalanobis matrix gives largest weights to the variables with the smallest prediction error. More in detail, when the variables are continuous and in presence of arbitrary patterns of missing items, a typical application of the PMM is the following (Little, 1988).

1. The parameters of a multivariate Gaussian distribution are estimated through the EM algorithm (Dempster et al., 1977) using all the available data (complete and incomplete).
2. Based on the estimates from EM, for each incomplete unit (recipient), predictions of the missing items conditional on the observed ones are computed. The same predictive means (i.e., corresponding to the same missing pattern) are computed for all the complete observations (donors).
3. Each recipient is matched to the donor having the closest predictive mean with respect to the Mahalanobis distance defined through the residual covariance matrix from the regression of the missing items on the observed ones.
4. Missing items are imputed in each recipient by transferring the corresponding values from its closest donor.

The NND method is a common hot deck method, in which a donor is selected from the complete cases in order to minimize some similarity measure, such as the Euclidean distance. In this application, the matching variables used to compute the Euclidean distance are  $Y_1$ ,  $Y_6$ ,  $Y_7$  and  $PC$  (if present), and the variables to be imputed are the ratios of  $Y_3$ ,  $Y_4$ ,  $Y_5$ ,  $Y_8$  and  $Y_{10}$  to  $Y_1$ . The final imputed value is obtained by multiplying the imputed ratios by the size variable  $Y_1$  of the recipient unit, this technique is also known as *ratio hot-deck* (de Waal et al., 2011). This method is preferable to the classical one that imputes directly the value observed in the closest donor, because it ensures that the values of the variables to be imputed are coherent with respect to the value of the reference variable. The reason why  $Y_1$  has been treated

as a size variable, instead of the commonly used *Number of employees*, is that it has both the lowest rate of missing data and the highest quality from a content point of view. When  $Y_1$  is zero the ratio cannot be computed, in this case the standard NND is used.

In this context, both the PMM and NND approaches have the advantage to recover live values from donors. Since the PMM technique relies on a multivariate normal model, it has been used to treat variables having a genuine continuous distribution. On the contrary, the NND method has been used to treat variables with distribution characterized by 0 inflation and a non-linear relation.

Finally, the imputation of  $Y_2$  and  $Y_9$ , representing the two components of *CS* has been carried out through a two-step process, composed by a logistic and a linear regression model. In the first step, we applied a stepwise logistic regression using as covariates  $Y_1$ ,  $Y_3$ ,  $Y_6$ ,  $Y_7$ , *CS*, a modified version of the economic divisions *Nace2* and *PC* (if present) in order to assign each enterprise to one of the 3 subpopulations characterized by the presence or absence of the two components (yes/yes, yes/no, no/yes). The assignment is based on random drawing from a multinomial distribution with parameters corresponding to the probabilities estimated through the logistic model. In the second step, for the enterprises which have been assigned to the subpopulation with only one component, the total value of *CS* has been imputed to such component. For the other enterprises, we estimated the value of the two components through a linear regression model with the same covariates used in the logistic model. This approach has been compared with a NND approach through a simulation study, resulting in a better efficiency (both in terms of time consuming and accuracy of the estimates) of the two-step approach. The difficulty in the imputation of these variables is both in their nature and in the nature of the total *CS* that is semi-continuous and not positive. This means that the value *CS* could be generated by any linear combination of the two components. As an hypothesis, when the variable *CS* is equal to 0, the two components have been imputed to be equal to 0.

For enterprises with information coming from SDS-G and Unico5, all the information on revenues is complete. Costs are imputed through random ratio hot-deck within suitable defined imputation cells. The donors are chosen from the SME survey, this is the same as drawing a vector of ratios from the estimated distribution of the ratios in SME. In particular, we imputed the composition of the costs using as size variable the total costs and transferring the compositional information from the survey data.

#### 4. Evaluation

The complexity of the imputation procedure and the particular nature of administrative data make the evaluation of the accuracy of the estimates a difficult task. A first overall evaluation has been obtained by comparing the estimates based on administrative data with the ones resulting from the classical procedure obtained by means of the SME sample survey data. The comparison has been made by using two different sets of estimation domains: the first is *Nace2* (aggregation of economic sectors), and the second corresponds to the different administrative sources where information is

taken from. While in the former case the domains are aggregations of planned survey domains - thus they are composed of sampling strata - in the latter case the domains - that have not been planned in the survey design phase - are used to analyze possible different levels of discrepancies between administrative data and survey data across the available sources.

For each typology of domain and each analysed variable, relative differences between total estimates based on administrative and sample data are considered. In detail, for a given variable  $Y$  with corresponding population total  $T_y$ , we have computed the indicator:

$$d_y^t = \frac{\hat{T}_y^s - \hat{T}_y^{ad}}{\hat{T}_y^{ad}} \times 100,$$

where  $\hat{T}_y^s$  is the estimate of  $T_y$  obtained with the sample data through the calibration estimator currently used for SME survey, and  $\hat{T}_y^{ad}$  is the estimate computed on the entire archive by summing up all the values. In order to distinguish the source of discrepancies due to the sampling and the measurement error, we have also considered, for each domain, the additional estimate  $\hat{T}_y^{ad,s}$ , that results from using the SME survey estimator on the sampled units, with the replacement of the survey data with the administrative data. As approximate measures of the measurement effect and sample error respectively, we introduce the following two indicators:

$$d_y^m = \frac{\hat{T}_y^s - \hat{T}_y^{ad,s}}{\hat{T}_y^{ad}} \times 100, \quad d_y^s = \frac{\hat{T}_y^{ad,s} - \hat{T}_y^{ad}}{\hat{T}_y^{ad}} \times 100.$$

Thus, the total difference is decomposed into the sum of two differences associated with the two mechanisms:

$$d_y^t = d_y^m + d_y^s. \quad (1)$$

Note that the indicator  $d_y^m$  that evaluates the “measurement effect”, being based on the comparison of different measures only on the sample units, is also affected by sampling error. In particular, a few gross errors may have an high impact on the indicator.

Table 6 reports the indicators  $d^t$  and  $d^m$  for the three variables  $Y_1$ ,  $IC$ , and the Value Added computed as  $VA = \sum_{i=1}^5 Y_i - IC - Y_9$  by source. In Table 7 results are shown for the following economic divisions Nace2: *Manufacture of textiles* (Nace2=13), *Construction of buildings* (Nace2=41), *Wholesale and retail trade and repair of motor vehicles and motorcycles* (Nace2=45) and *Architectural and engineering activities; technical testing and analysis* (Nace2=71). In the tables, the size of domains in the population  $N$  and in the sample  $n$  are also reported.

Results in Tables 6 and 7 show that the largest component in the decomposition (1) is the one associated with the sampling error. This result is encouraging because it implies that the transition from designed based inference to an estimation approach based on administrative sources would result in a significant improvement of the estimation accuracy.



**Table 6 - Discrepancies between sample estimates and estimates based on administrative data for different administrative data sources: total differences ( $d^t$ ) and measurement component ( $d^m$ )**

Source	N	n	$d^t$			$d^m$		
			$Y_1$	TC	VA	$Y_1$	TC	VA
Tot	4432928	74112	-6.5	-8.8	0.2	-0.9	-0.5	-0.9
FS	714885	34284	-2.6	-4.6	3.1	-0.9	-0.7	-2.3
SDS-F	2220050	31732	11.4	11.3	12.7	-0.5	-1.2	1.6
SDS-G	616050	3844	23.1	26.5	26.7	-0.2	13.5	-0.3
Unico1	34570	296	-38.9	-69	-26.7	-1	-3.3	0.3
Unico2	1746	32	-41	-40.6	-40	-1.7	-2.5	0.9
Unico3	120876	756	-70.7	-78.5	-55.7	-0.7	-5.3	8.4
Unico4	33676	356	-69	-81.1	-34.6	-0.5	-5	14.6
Unico5	481517	1434	-58.2	-51.5	-59.1	-1.4	3.3	-1.3
Unico6	7371	151	-70.1	-76.9	-59.5	0.1	-2.5	-6.5
Unico7	13553	361	-68.5	-68.5	-59.4	-0.4	-2.5	13.8
Unico8	21585	381	-63.7	-55.3	-58.2	-2.4	8.7	-1.4
Irap	4201	89	-55.8	-63.1	-61	-0.1	-0.9	4.1
NA	162848	396	-91	-90.9	-89.8	-2.5	-1.3	-4.6

**Table 7 - Discrepancies between sample estimates and estimates based on administrative data for some economic divisions (Nace2): total differences ( $d^t$ ) and measurement component ( $d^m$ )**

Nace2	N	n	$d^t$			$d^m$		
			$Y_1$	TC	VA	$Y_1$	TC	VA
13	15669	1275	8.4	10.1	3.8	0	1.1	-3.5
41	150417	2625	-18	-21.7	-9.2	-0.8	2	1.5
45	118985	2649	0.4	0.3	1.7	1	0.9	-0.7
71	212880	1009	-9.5	-15.3	-4.6	-2.6	1.5	-0.5

An important issue in the evaluation of an estimation procedure is the assessment of the estimate accuracy. According to the estimation approach so far used in Istat, the SBS estimates for SME are based on a sample survey, hence the assessment of their accuracy relies on designed-based inference. As already mentioned, massive use of administrative information requires a change of paradigm. In fact, differently from the context of sample survey, the availability of administrative information is not under control of the researcher, so that some model assumptions are necessary. In particular, one has to think of data as *iid* realizations from a statistical (possibly not explicitly specified) model. This is generally referred to as super-population model. In this framework, the inferential approach is predictive, i.e., the missing values are imputed (predicted) on the basis of the available information. Thus, the uncertainty of the resulting estimates are essentially due to the prediction error associated with the imputation procedure.

Some limitations of the present evaluation methods should be mentioned. First, comparison is performed at one point in time, so that results should be assessed in future occasions. Second, evaluation of measurement component of the total error is based on survey data and thus it is affected by sampling error. It is interesting to note that, because of compensations, substantial differences in the estimates of Turnover

$(Y_1)$  and Total Costs ( $TC$ ) do not result in significant discrepancy for the variable Value Added ( $VA$ ).

If predictions were based on some parametric regression model and the missing patterns were enough simple, standard analytic techniques could be used to evaluate the estimate of the estimator variance (Valliant, 2000). In cases where missing patterns are arbitrary, but imputations are obtained from a unique multivariate normal, the Rubin multiple imputation approach can be (relatively) simply applied to assess the precision of the estimate of any finite population quantity (Rubin, 1987). In the present case, however, the imputation procedure is complex and is composed of many different techniques. This complexity makes it difficult to use standard procedures for the assessment of the uncertainty in the final output. In particular, the assumed super-population model is not explicitly specified and it is only implicitly defined through the imputation procedures that have been used. Because of this characteristic, a replication approach seems to be more appropriate than an analytic approach. However, common univariate techniques for the variance estimation such as jackknife and bootstrap (Wolter, 2007) are difficult to extend to our context and further research is needed.

## References

- Curatolo S., V. De Giorgi, F. Oropallo, A. Puggioni, G. Siesto. 2016. “Quality analysis and harmonization issues in the context of the SBS frame”, *Rivista Statistica Ufficiale*, N.1/2016.
- De Waal T., J. Pannekoek, S. Scholtus. 2011 *Handbook of Statistical Data Editing and Imputation* Wiley
- Dempster A.P., N.M. Laird, D.B. Rubin. 1977. “Maximum likelihood from incomplete data via the EM algorithm”, *Journal of the Royal Statistical Society*, B 39; 1-38.
- Istat. 2014. I nuovi conti nazionali in SEC 2010. “Innovazioni e ricostruzione delle serie storiche (1995-2013)”. *Nota informativa Istat*, Ottobre 2014.
- Little R.J.A. 1988. “Missing-Data Adjustments in Large Surveys”. *Journal of Business & Economic Statistics*, 6, 3; 287-296.
- Rubin D.B. 1987. *Multiple Imputation for Nonresponse in Surveys*, Wiley
- Valliant R., A.H. Dorfman, R.M. Royal. 2000. *Finite Population Sampling and Inference: A Prediction Approach*, Wiley
- Wolter K.M. 2007. *Introduction to Variance Estimation*, Springer