

SWISS STRUCTURAL BUSINESS STATISTICS: DATA HARMONIZATION FOR THE CONSTRUCTION OF FULL-TIME EQUIVALENTS

Prepared by Desislava Nedyalkova (Desislava.Nedyalkova@bfs.admin.ch) and Daniel Assoulin
(Daniel.Assoulin@bfs.admin.ch), Swiss Federal Statistical Office

I. Summary

The construction of full-time equivalents (FTE) for the Swiss Structural Business Statistics is based on the integration of register and survey data. Register data comes from the Social security register (SR) and from the business register (BR). FTE are not collected in SR. They need to be constructed on the basis of a model fitted on matched data (survey-register). This data set shows differences between the variables number of employees reported in the survey and in the register. These differences are treated in order to make FTE coming from the survey coherent with data obtained from the registers. A first method which treats the differences by a simple ratio adjustment was applied. Our analyses has shown that this is not the optimal solution. Another approach was thus developed. This method permits to treat the differences by taking into account information about employment.

II. Introduction

The Swiss Business Census (BC) was held for the last time in 2008. It played an important role for producing various statistics on the structure of the Swiss economy. For the reference year 2011 it was replaced by an integrated system called STATENT (Swiss Structural Business Statistics). STATENT is mainly based on the business register (BR), the social security register (SR) and complementary surveys like the Quarterly Survey of Employment (JobStat).

The transition from the BC to STATENT induces several changes in definition and methodology. The principal differences concern the covered units, definition of employment¹ and periodicity. For instance, the BC was conducted every 3-4 years whereas STATENT appears each year. The BC referred to an exact date, whereas STATENT refers to the last month of the year.

Another major difference between BC and STATENT is the way full-time equivalents (FTE) are calculated. In the past, FTE were derived from the information in the BC (employment rates) which is not available in the STATENT. For this reason the construction of such a variable is an important challenge for STATENT.

For enterprises not included in a complementary survey, FTE for STATENT are constructed using a linear prediction model based on explanatory variables coming from the register. This model is fitted on matched data coming from the register and some complementary surveys. In case that an enterprise has FTE collected from a complementary survey, these FTE will be in principle used in STATENT. The integration of data coming from different sources reveals the existence of inconsistencies regarding the number of employees. In such cases the FTE from the survey has to be adapted in a way that it reflects the employees according to the register.

¹The BC counted employees that worked at least 6 hours per week in an enterprise or an establishment. In STATENT all persons working in an enterprise (as wage-earner or as independent) and paying their mandatory contributions to the SR for a minimum annual wage of CHF 2 300 are counted for (criteria for reference year 2011)

We begin by a description of the data and of the FTE model based on matched data available in survey and registers. Next, we present two different methods for treating inconsistencies between the different sources. We describe a first approach based on a ratio adjustment and show how this approach is employed in the FTE model. Then, we present a second approach in which differences are treated by taking into account information about employment. Finally we show how FTE for STATENT are constructed.

III. FTE model

For the construction of the model, we used survey data (fourth quarter of 2011) matched with register data on the enterprise level. These are mainly single-establishment enterprises (EUNT) for which we know FTE, annual standardized wages and some other variables like NUTS2 region and NACE. The model is estimated separately for the subpopulation of men and women in each of the two economic sectors of activity 2 and 3. For the sake of simplicity, the same notation is used to describe the estimated models.

On the basis of the information contained in the survey about occupation levels and on the salary distribution in the register, we construct, for each NACE section, four salary classes. These classes form the basis for the construction of the explanatory variables of the model. The variable of interest is the number of FTE.

The model we want to estimate is the following:

$$y_i = \alpha_1 \cdot V_{i1} + \sum_{j=2}^4 \alpha_{jkl} V_{ij} + \epsilon_i, \quad (1)$$

where:

- y_i , is the number of FTE of an enterprise i ,
- V_{ij} , the number of employees of enterprise i in the salary class j ($j = 1, \dots, 4$) ($\sum V_{ij} = \text{EMPTOT_R}$)
- α_1 , regression coefficient for V_{i1} ,
- α_{jkl} , regression coefficient for V_{ij} in NUTS 2 k ($k = 1, \dots, 7$) and NACE section ℓ ,
- ϵ_i , residual with $E(\epsilon_i) = 0$ and $\text{Var}(\epsilon_i) = \sigma^2 \text{EMPTOT_R}_i$.

IV. Harmonization of variables

Matched data on which FTE are estimated present certain differences between the number of employees from the survey (EMPTOT_S) and this from the register (EMPTOT_R). There exist different methods for variable adjustment which can be used to treat these differences, e.g. *prorating* and *generalized ratio adjustment* (Panekoek, 2011; Panekoek, 2014). For instance, the *prorating* method represents a simple multiplicative adjustment which is applied on variables employed in control rules.

A. A first approach for treating the differences

Based on the methods described above, for each enterprise i , we define a new variable FTE_R (by gender) as follows:

$$\text{FTE_R}_i = \eta_i \text{FTE_S}_i, \quad (2)$$

where $\eta_i = \text{EMPTOT_R}_i / \text{EMPTOT_S}_i$ and FTE_S_i is the survey FTE.

This new variable *harmonized with the register* will be used for modelling. In this way predicted FTE will be consistent with the values of EMPTOT_R. If the inconsistencies are treated according to Equation (2), we can rewrite the equation of Model (1) as follows:

$$\eta_i \text{FTE_}S_i = \alpha_1 \cdot V_{i1} + \sum_{j=2}^4 \alpha_{jkl} V_{ij} + \epsilon_i, \quad (3)$$

or

$$\text{FTE_}S_i = \alpha_1 \cdot \frac{V_{i1}}{\eta_i} + \sum_{j=2}^4 \alpha_{jkl} \frac{V_{ij}}{\eta_i} + \frac{\epsilon_i}{\eta_i}, \quad (4)$$

In the case $\eta_i > 1$, we can interpret Equation (4) as follows: The adjustment between EMPTOT_R and EMPTOT_S is done uniformly in the four salary classes by reducing the number of employees by η_i . This procedure can be justified only in the case where the inconsistencies are independent of the salary classes.

B. A new approach for treating the differences

We present an alternative of Model (2) in which inconsistencies in the variables number of employees are not treated uniformly. We will examine the problem for the following cases:

- Case 1: EMPTOT_R > EMPTOT_S.
- Case 2: EMPTOT_S > EMPTOT_R.
- Case 3: EMPTOT_S = EMPTOT_R.

Let $\text{diff_ab} = \text{EMPTOT_R} - \text{EMPTOT_S}$ (by gender) and $\text{diff_ba} = \text{EMPTOT_S} - \text{EMPTOT_R}$ (by gender). Note that an enterprise for which the variables EMPTOT present extreme differences either in the subpopulation of men or in the subpopulation of women will be treated as if survey data is missing.

B.1. Treatment of Case 1.

We suppose that $\text{EMPTOT_R} > \text{EMPTOT_S}$. Knowing the number of employees in each salary class, we estimate the coefficients of the following model (by gender and economic sector of activity):

$$\text{diff_ab}_i = \sum_{j=1}^4 \beta_j V_{ij} + \epsilon_i,$$

under the hypothesis $\text{Var}(\epsilon_i) = \sigma^2 \text{EMPTOT_}R_i$. This is not done with the aim to estimate the difference between EMPTOT_R and EMPTOT_S. We are rather interested in the estimated coefficients, $\hat{\beta}_j$, which can be seen as estimation of the proportion of persons in the salary class j which are in the register but not in the survey.

	$\hat{\beta}_1$		$\hat{\beta}_2$		$\hat{\beta}_3$		$\hat{\beta}_4$	
	Estimator	StdErr	Estimator	StdErr	Estimator	StdErr	Estimator	StdErr
m/s2	0.659	0.063	0.562	0.016	0.153	0.019	0.039	0.004
m/s3	0.808	0.025	0.624	0.012	0.329	0.012	0.088	0.003
w/s2	0.724	0.023	0.495	0.015	0.241	0.013	0.051	0.006
w/s3	0.693	0.012	0.390	0.009	0.125	0.006	0.129	0.004

TABLE 1. Values of $\hat{\beta}_j$ and the estimated standard errors

Table 1 contains the estimated coefficients, $\hat{\beta}_j$, as well as their standard errors, obtained with the ROBUSTREG procedure in SAS with weights proportional to $1/\text{EMPTOT_}R_i$. It indicates, for example,

that the coefficients for the class of employees having the smallest wages (salary class 1) are larger than the coefficients for the salary class 2. Thus, at least in this subpopulation (Case 1), a uniform treatment using Equation (2) seems not to be appropriate.

Initial idea - probability proportional to size (PPS) sampling of fixed size.

We suppose that the persons which are not in the survey come from a sample s_i of fixed size $n_i = \text{diff_ab}_i$ of probabilities proportional to $\hat{\beta}_j$ (see Table 1). The probability that a person d belonging to the set of wage-earners in the salary class j is removed from this set is given by:

$$P(d \in s_i) = n_i \frac{\text{mos}_d}{\text{mos}_i} = \frac{n_i \hat{\beta}_j}{\sum_{j=1}^4 \hat{\beta}_j V_{ij}} = \hat{\beta}_j \frac{n_i}{\hat{n}_i} = \hat{\beta}_{ij}^*,$$

where $\text{mos}_d = \hat{\beta}_j$ and $\text{mos}_i = \sum_{d \in s_i} \text{mos}_d = \sum_{j=1}^4 \hat{\beta}_j V_{ij} = \hat{n}_i$.

If $\hat{\beta}_{ij}^* \geq 1$, then the person will be automatically removed (Särndal et al. (1992, p.89)).

We should note that $\hat{\beta}_{ij}^*$ can be seen as a $\hat{\beta}_j$ adjusted so that $\sum_{j=1}^4 \hat{\beta}_{ij}^* V_{ij} = n_i$.

Calculation of the average number of persons which should be removed in each salary class.

Some inconveniences of using PPS sampling of persons are its random aspect with a potential impact on comparability over years and the difficulty of implementation in production (matched data is on enterprise and not on person level). This has led us to develop a general procedure to replace the random PPS sample. Instead of drawing a sample we calculate the expected number of persons in class j to be selected in the sample s_i . This number is given by:

$$E\left(\sum_{d \in \mathcal{V}_{ij}} 1(d \in s_i)\right) = V_{ij} P(d \in s_i) = V_{ij} \hat{\beta}_{ij}^*. \quad (5)$$

where \mathcal{V}_{ij} denotes the set of employees of enterprise i in the salary class j .

As in the case of PPS sampling, our procedure first removes all persons for which $\hat{\beta}_{ij}^* \geq 1$. Next, we calculate the mean number of persons which have to be eliminated according to Equation (5). At the end of this iterative procedure we obtain new variables:

$$\tilde{V}_{ij} = V_{ij} - V_{ij} \hat{\beta}_{ij}^*,$$

such that $\sum_{j=1}^4 \tilde{V}_{ij} = \text{EMPTOT_S}_i$. These new variables will replace the variables V_{ij} in the estimation of Model (1) where y_i will be given by FTE_S_i .

B.2. Treatment of Case 2.

We suppose that $\text{EMPTOT_S} > \text{EMPTOT_R}$. Knowing the number of employees working at part time III (T_{i1}), part time II (T_{i2}), part time I (T_{i3}) and full time (T_{i4}) from the survey data, we estimate the following model (by sex and economic activity sector):

$$\text{diff_ba}_i = \sum_{j=1}^4 \gamma_j T_{ij} + \epsilon_i,$$

The used procedure is PROC ROBUSTREG in SAS with weights proportional to $1/\text{EMPTOT_S}_i$. The estimated coefficients and their standard errors are given in Table 2. These estimated coefficients can be seen as estimation of the proportion of persons working at part time III, for example, that are in the survey but not in the register. It can be seen from the table that the coefficients for the persons working at part time III are larger than the coefficients for the persons working full time.

	$\hat{\gamma}_1$		$\hat{\gamma}_2$		$\hat{\gamma}_3$		$\hat{\gamma}_4$	
	Estimator	StdErr	Estimator	StdErr	Estimator	StdErr	Estimator	StdErr
m/s2	0.469	0.039	0.657	0.041	0.256	0.032	0.071	0.002
m/s3	0.448	0.014	0.237	0.011	0.203	0.014	0.123	0.003
w/s2	0.548	0.025	0.329	0.019	0.063	0.018	0.065	0.007
w/s3	0.442	0.010	0.251	0.008	0.002	0.007	0.140	0.005

TABLE 2. Values of $\hat{\gamma}_j$ and the estimated standard errors

Adaptation of the harmonization procedure to Case 2.

We know that in order to eliminate the differences between the employment variables EMPTOT_S and EMPTOT_R and make them consistent we have to eliminate, for each enterprise, a fix number of persons, $n_i^* = \text{diff_ba}_i$. Using the coefficients $\hat{\gamma}_j$, we apply the same procedure as for the Case 1, with the required modifications. In this way we obtain the new variables \tilde{T}_{ij} such that $\sum_{j=1}^4 \tilde{T}_{ij} = \text{EMPTOT_R}_i$.

Let suppose that FTE_S can be modeled as follows:

$$\text{FTE_S}_i = \sum_{j=1}^4 \delta_j T_{ij} + \epsilon_i, \quad (6)$$

where $\text{Var}(\epsilon_i) = \sigma^2 \text{EMPTOT_S}_i$. This model, by sex and economic activity sector, is estimated using PROC GLM from SAS with weights proportional to $1/\text{EMPTOT_S}$. The estimated coefficients and their standard errors are given in Table 3.

	$\hat{\delta}_1$		$\hat{\delta}_2$		$\hat{\delta}_3$		$\hat{\delta}_4$		R^2
	Est.	StdErr	Est.	StdErr	Est.	StdErr	Est.	StdErr	
m/s2	0.110	0.006	0.300	0.007	0.632	0.005	0.998	0.000	0.999
m/s3	0.078	0.004	0.278	0.003	0.645	0.004	0.998	0.001	0.997
w/s2	0.095	0.008	0.293	0.006	0.654	0.006	0.994	0.002	0.997
w/s3	0.091	0.004	0.281	0.003	0.660	0.003	0.988	0.002	0.992

TABLE 3. Values of $\hat{\delta}_j$ and the estimated standard errors

Using the estimated coefficients of Model (6), we calculate a new adjusted variable FTE_S, denoted by FTE_R, which is coherent with the variable EMPTOT_R:

$$\text{FTE_R}_i = \min \left(\text{FTE_S}_i \frac{\sum \hat{\delta}_j \tilde{T}_{ij}}{\sum \hat{\delta}_j T_{ij}}, \text{EMPTOT_R}_i \right). \quad (7)$$

We can explain the minimum in Equation (7) by the fact that mean occupational level of harmonized data, $\text{FTE_R}/\text{EMPTOT_R}$, should be bounded by 1.

C. Consequences of data harmonization for the estimation of the FTE model

Table 4 presents the variables used in the estimation of Model (1) in Case 1, 2 and 3, respectively. It can be seen that in Case 1 it is the explanatory variables that are adjusted in order to correspond to the FTE recorded in the survey. In Case 2, where the total number of employees is larger for the survey than for the register, the explanatory variables remain unchanged but the dependent variable for the model is adjusted.

Case	Variable of interest	Explanatory variables
1	FTE_S	\tilde{V}_{ij}
2	FTE_R	V_{ij}
3	FTE_S	V_{ij}

TABLE 4. Variables used in the model

D. Consequences for calculating the FTE in STATENT

We have already explained the consequences of the harmonization on the variables of the model. Now we show how we calculate FTE for enterprises for which survey data is available.

Case 1: Let us denote $FTE_R_{i,model} = \sum_{j=1}^4 \hat{\alpha}_j V_{ij}$ the FTE calculated using the estimated FTE model parameters and the variables V_{ij} and $FTE_S_{i,model} = \sum_{j=1}^4 \hat{\alpha}_j \tilde{V}_{ij}$ the FTE calculated using the FTE model parameters and the variables \tilde{V}_{ij} . Then, the ratio of these variables is applied to the FTE_S in order to produce the harmonized FTE, denoted by FTE_R. Thus,

$$FTE_R_i = \min \left(FTE_S_i \frac{\sum_{j=1}^4 \hat{\alpha}_j V_{ij}}{\sum_{j=1}^4 \hat{\alpha}_j \tilde{V}_{ij}}, EMPTOT_R_i \right).$$

The minimum is explained by the fact that mean occupation level of harmonized data should be bounded by 1. In Case 2 the harmonized FTE_R is calculated by applying a multiplicative or enhanced ratio adjustment as defined in Equation (7). For case 3 we have that $FTE_R = FTE_S$.

V. Conclusion

This document presents the methodology used to treat the inconsistencies between the different data sources used for the development and the estimation of the FTE model for the construction of the Swiss Structural Business Statistics. The analyses put into question the application of a simple ratio adjustment according to Definition (2). The divergences seem rather to be due to low wages or small occupation levels. The presented harmonization based on PPS sampling takes into account information on the employment type (salary class, occupation level). However, this procedure has the inconvenience to be random and difficult to apply in practice. If we use expected sample sizes instead of random sampling for adjusting the number of employees in the different salary classes, we can overcome this inconvenience. This new approach was tested and used for STATENT 2011.

References

- [1] Pannekoek, J. (2011), Models and algorithms for micro-integration, dans Final Report on WP2: Methodological developments, ESSNET on Data Integration, <http://www.cros-portal.eu/content/wp2-development-methods>
- [2] Pannekoek, J. (2014) Method: Reconciling Conflicting Microdata, dans *Memobust Handbook on Methodology of Modern Business Statistics*, <http://www.cros-portal.eu/content/reconciling-conflicting-microdata-method>
- [3] Särndal, C.E., Swensson, B. et Wretman, J. (1992), *Model Assisted Survey Sampling*, Springer.