**Republic of Bulgaria Enterprise Surveys Data Set**

## 1. Introduction

1.      This document provides additional information on the data collected in the Republic of Bulgaria during the calendar years 2007 and 2008. The reader will find a description of the sampling design of the data, the data set structure and additional information that may be useful when using the data, such as information on non-response cases and the appropriate use of the weights.

## 2. Sampling Structure

2.      The sample for the Republic of Bulgaria was selected using stratified random sampling, following the methodology explained in the Sampling Manual. Stratified random sampling[1] was preferred over simple random sampling for several reasons[2]:
        a. To obtain unbiased estimates for different subdivisions of the population with some known level of precision.
        b. To obtain unbiased estimates for the whole population. The whole population, or universe of the study, is the non-agricultural economy. It comprises: all manufacturing sectors according to the group classification of ISIC Revision 3.1: (group D), construction sector (group F), services sector (groups G and H), and transport, storage, and communications sector (group I). Note that this definition excludes the following sectors: financial intermediation (group J), real estate and renting activities (group K, except sub-sector 72, IT, which was added to the population under study), and all public or utilities-sectors.
        c. To make sure that the final total sample includes establishments from all different sectors and that it is not concentrated in one or two of industries/sizes/regions.
        d. To exploit the benefits of stratified sampling where population estimates, in most cases, will be more precise than using a simple random sampling method (i.e., lower standard errors, other things being equal.)
        e. Stratification may produce a smaller bound on the error of estimation than would be produced by a simple random sample of the same size. This result is particularly true if measurements within strata are homogeneous.
        f. The cost per observation in the survey may be reduced by stratification of the population elements into convenient groupings.

3.      Three levels of stratification were used in this country: industry, establishment size, and oblast (region). The original sample designs with specific information of the industries and regions chosen are included in the attached Excel file (Sampling Report.xls.)

---

[1] A stratified random sample is one obtained by separating the population elements into non-overlapping groups, called strata, and then selecting a simple random sample from each stratum. (Richard L. Scheaffer; Mendenhall, W.; Lyman, R., "Elementary Survey Sampling", Fifth Edition).
[2] Cochran, W., 1977, pp. 89; Lohr, Sharon, 1999, pp. 95

4.      Industry stratification was designed in the way that follows: the universe was stratified into 4 manufacturing industries, 2 services industries -retail and IT-, and one residual sector as defined in the sampling manual. Each industry had a target of 120 interviews. For the manufacturing industries sample sizes were inflated by about 25% to account for potential non-response cases when requesting sensitive financial data and also because of likely attrition in future surveys that would affect the construction of a panel.

5.      Size stratification was defined following the standardized definition for the rollout: small (5 to 19 employees), medium (20 to 99 employees), and large (more than 99 employees). For stratification purposes, the number of employees was defined on the basis of reported permanent full-time workers. This seems to be an appropriate definition of the labor force since seasonal/casual/part-time employment is not a common practice, except in the sectors of construction and agriculture. The former is represented by 22 establishments out of 1015 establishments in the total sample and the latter is not included in the sample frame.

6.      Regional stratification was defined in 4 regions. These regions are Sofia, Plovdiv, Varna, and Burgas. These regions contain the main cities and industrial clusters of the country. These areas comprise 40% of the country's population and 23% of the territory.

## 4. Sampling implementation

7.      Given the stratified design, sample frames containing a complete and updated list of establishments for the selected regions were required. Great efforts were made to obtain the best source for these listings. However, the quality of the sample frames was not optimal and, therefore, some adjustments were needed to correct for the presence of ineligible units. These adjustments are reflected in the weights computation (see below.)

8.      The source of the sample frame was the Bulgarian National Statistical Institute (2007).

9.      The quality of the frame was assessed at the onset of the project. The frame proved to be useful though it showed positive rates of non-eligibility, repetition, non-existent units, etc. These problems are typical of establishment surveys, but given the impact these inaccuracies may have on the results, adjustments were needed when computing the appropriate weights for individual observations. The percentage of confirmed non-eligible units as a proportion of the total number of contacts to complete the survey was 34% (951 out of 2,768 establishments).

## 3. Data Base Structure:

10.     The structure of the data base reflects the fact that 3 different versions of the questionnaire were used. The basic questionnaire, the Core Module, includes all common questions asked to all establishments from all sectors (manufacturing, services and IT). The second expanded variation, the Manufacturing Questionnaire, is built upon the Core Module and adds some specific questions relevant to the sector. The third expanded

variation, the Services Questionnaire, is also built upon the Core Module and adds to the core specific questions relevant to either retail or IT. Each variation of the questionnaire is identified by the index variable, *a0.*

11.     All variables are named using, first, the letter of each section and, second, the number of the variable within the section, i.e. *a1* denotes section *A*, question *1*. Variable names preceded by a prefix "*BG*" are specific to the Republic of Bulgaria and, therefore, they may not be found in the implementation of the rollout in other Countries. All other suffixed variables are global and are present in all country surveys over the world. All variables are numeric with the exception of those variables with an "x" at the end of their names. The suffix "x" denotes that the variable is alpha-numeric.

12.     There are 3 establishment identifiers, *idstd*, *idu*, and *id*. The first is a global unique identifier. The second is a regional unique identifier, and *the* third one is a country unique identifier.  The variables *a2* (sampling region), *a6a* (sampling establishment's size), and *a4a* (sampling sector) contain the establishment's classification into the strata chosen for each country using information from the sample frame. The strata were defined according to the guidelines described above.

13.      As noted above, there are 3 levels of stratification: industry, size and region. Different combinations of these variables generate the strata cells for each industry/region/size combination. The variable *strata* identifies each cell. A distinction should be made between the variable *a4a* and *NACE*. The former gives the establishment's classification into one of the chosen industry-strata, whereas the latter gives the actual establishment's industry classification in the sample frame.

14.     All of the following variables contain information from the sampling frame and were defined with the sampling design. They may not coincide with the reality of individual establishments as sample frames may contain inaccurate information. The variables containing the sample frame information are included in the data set for researchers who may want to further investigate statistical features of the survey and the effect of the survey design on their results.
>    -*a2* is the variable describing sampling regions (oblasts)
>    -size_sample: coded using the same standard for small, medium, and large establishments as defined above. The code *-9* was used to indicate units for which size was undetermined in the sample frame.
>    -ind_sample: coded using ISIC codes for the chosen industries for stratification. These codes include most manufacturing industries (15 to 36), and retail, and IT for services (52, and 72 respectively). All establishments within the residual stratum were coded with ind_sample=2.
>    -strata: unique stratum identifier. This variable is important in Stata when setting the data set as a survey data set.
>    -isic: original ISIC classification from the sample frame

15.     The surveys were implemented following a 2 stage procedure. In the first stage a screener questionnaire was applied over the phone to determine eligibility and to make

appointments; in the second stage, a face-to-face interview took place with the Manager/Owner/Director of each establishment. The variables *a4b* and *a6b* contain the industry and size of the establishment from the screener questionnaire. Variables *a8* to *a11*contain additional information and were also collected in the screening phase.

16.     Note that there are additional variables for location (*a3x*), industry (*d1a2*), and size (*l1*, *l6* and *l8*) that reflect more accurately the reality of each establishment. Advance users are advised to use these variables for analytical purposes.

17     Variable *a3x* indicates the actual location of the establishment. There may be divergences between the location in the sampling frame and the actual location, as establishments may be listed in one place but the actual physical location is in another place.

18.     Variable *d1a2* indicates the actual ISIC code of the main output of the establishment as answered by the interviewee. This is probably the most accurate variable to classify establishments by activity.

19.     Variables *l1*, *l6* and *l8* were designed to obtain a more accurate measure of employment accounting for permanent and temporary employment. Special efforts were made to make sure that this information was not missing for most establishments. *l1*, *l6* and *l8* are missing for 1, 7, and 7 establishments, respectively. That is, *l1*, *l6* and *l8* are not missing for at least 99% of the dataset.

**3. Weights**

20.     Since the sampling design was stratified and employed differential sampling individual observations should be properly weighted when making inferences about the population. Under stratified random sampling unweighted estimates are biased unless sample sizes are proportional to the size of each stratum. With stratification the probability of selection of each unit is, in general, not the same. Consequently, individual observations must be weighted by the inverse of their probability of selection (probability weights or *pa* in Stata.)[3]

21.     Special care was given to the correct computation of the weights. Considering the varying quality of the sample frames, it was imperative to accurately adjust the totals within each region/industry/size stratum to account for the presence of ineligible units (the firm discontinued businesses or was unattainable, education or government establishments, establishments with less than 5 employees, no reply after having called in different days of the week and in different business hours, out of order, no tone in the phone line, answering machine, fax line, wrong address or moved away and could not get the new references) The information required for the adjustment was collected in the first stage of the implementation: the screening process. Using this information, each stratum cell of the universe was scaled down by the observed proportion of ineligible units within

---

[3] This is equivalent to the weighted average of the estimates for each stratum, with weights equal to the population shares of each stratum.

the cell. Once an accurate estimate of the universe cell (projections) was available, weights were computed using the number of completed interviews.
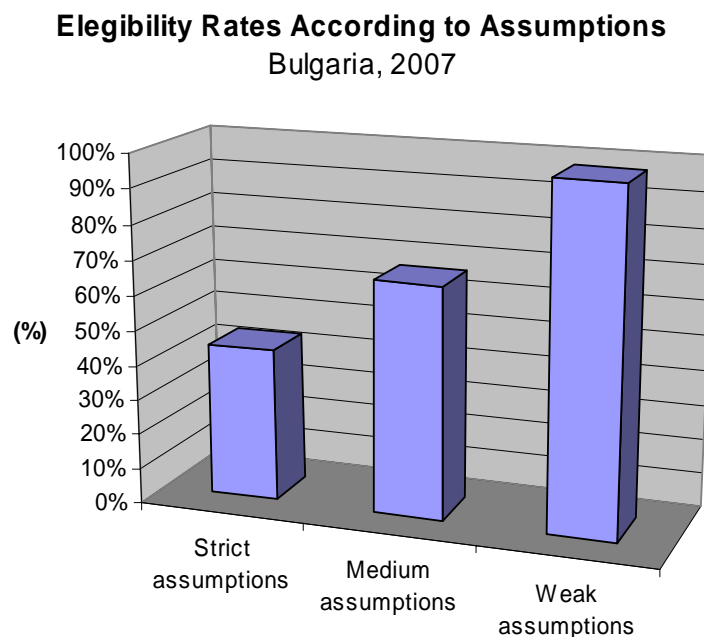
22.    For some units it was impossible to determine eligibility because the contact was not successfully completed. Consequently, different assumptions as to their eligibility result in different universe cells' adjustments and in different sampling weights. Three sets of assumptions were considered:
        a- Strict assumption: eligible establishments are only those for which it was possible to directly determine eligibility. The resulting weights are included in the variable *w_strict*.
        b- Median assumption: eligible establishments are those for which it was possible to directly determine eligibility and those that rejected the screener questionnaire or an answering machine or fax was the only response. The resulting weights are included in the variable *w_median*.
        c- Weak assumption: in addition to the establishments included in points a and b, all establishments for which it was not possible to finalize a contact are assumed eligible. This includes establishments with dead or out of service phone lines, establishments that never answered the phone, and establishments with incorrect addresses for which it was impossible to find a new address. The resulting weights are included in the variable *w_weak*. Note that under the weak assumption only observed non-eligible units are excluded from universe projections.
        The following graph exhibits the different eligibility rates under each set of assumptions.

**Elegibility Rates According to Assumptions**
Bulgaria, 2007



23.    Within each of these assumptions regarding eligibility a pair of weight sets was calculated. The first set of estimates calculated proportions using the raw sample count

for each cell. However, the achieved sample numbers in many cells were small. Hence, those eligibility rates, and the adjusted universe cells projections, are subject to relatively large sampling variations. Therefore a second set of more robust estimates (collapsed weights) was also produced. These estimates made use of the multiples of the relative eligibility rates for each industry, size, and region. Those relative rates were based on much larger samples than the individual cells and thus produced values with smaller sampling variations. The data sets include only these robust weights.

## 4. Appropriate use of the weights

24.    As discussed above, under stratified random sampling weights should be used when making inferences about the population. Any estimate or indicator that aims at describing some feature of the population should take into account that individual observations may not represent equal shares of the population.

25.    However, there is some discussion as to the use of weights in regressions (see Deaton, 1997, pp.67; Lohr, 1999, chapter 11, Cochran, 1953, pp.150). There is not strong large sample econometric argument in favor of using weighted estimation for a common population coefficient if the underlying model varies per stratum (stratum-specific coefficient): both simple OLS and weighted OLS are inconsistent under regular conditions.  However, weighted OLS has the advantage of providing an estimate that is independent of the sample design. This latter point may be quite relevant for the Enterprise Surveys as in most cases the objective is not only to obtain model-unbiased estimates but also design-unbiased estimates (see also Cochran, 1977, pp 200 who favors the used of weighted OLS for a common population coefficient.) [4]

26.    From a more general approach, if the regressions are descriptive of the population then weights should be used. The estimated model can be thought of as the relationship that would be expected if the whole population were observed[5]. If the models are developed as structural relationships or behavioral models that may vary for different parts of the population, then, there is no reason to use weights.

## 5. Non-response

27.    Survey non-response must be differentiated from item non-response. The former refers to refusals to participate in the survey altogether whereas the latter refers to the refusals to answer some specific questions. Enterprise Surveys suffer from both problems and different strategies were used to address these issues.

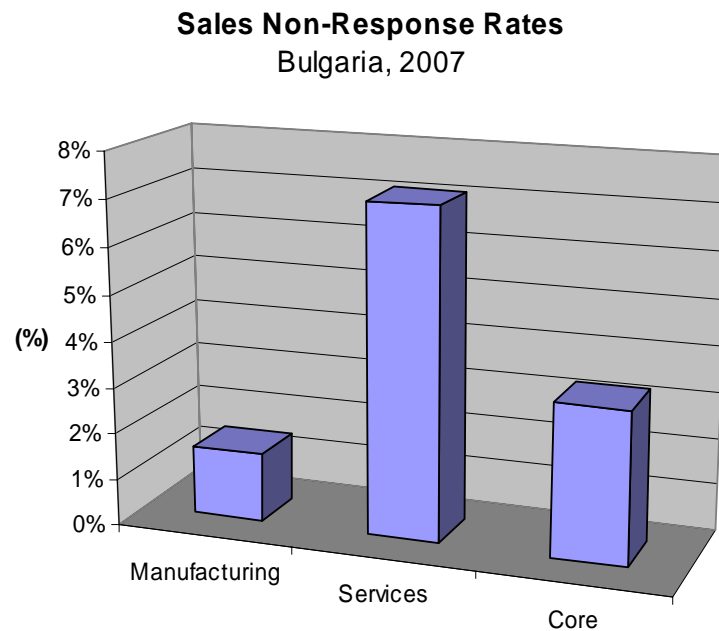28.    Item non-response was addressed by two strategies:

---

[4] Note that weighted OLS in Stata using the command regress with the option of weights will estimate wrong standard errors. Using the Stata survey specific commands svy will provide appropriate standard errors.
[5] The use weights in most model-assisted estimations using survey data is strongly recommended by the statisticians specialized on survey methodology of the JPSM of the University of Michigan and the University of Maryland.

a- For sensitive questions that may generate negative reactions from the respondent, such as corruption or tax evasion, enumerators were instructed to collect the refusal to respond as a different option from don't know (-7).
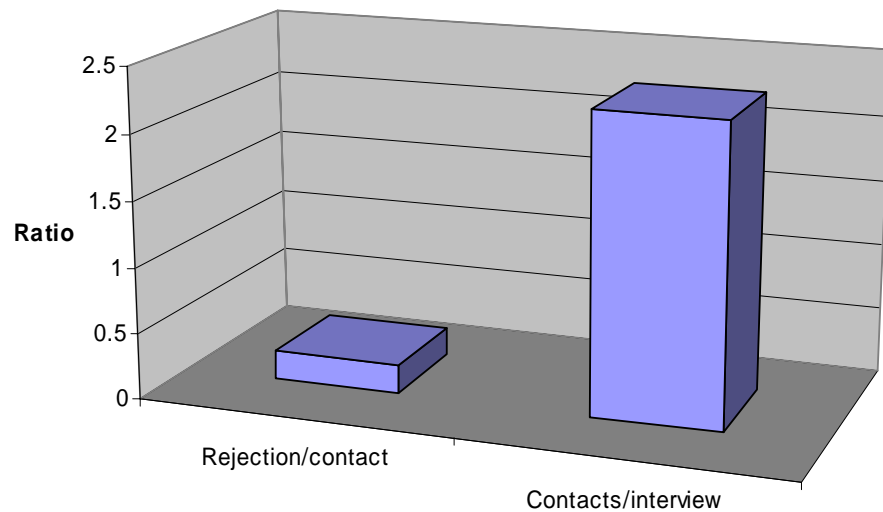
b- For information that establishments may consider too private to share such as financial information sample sizes were inflated by 25% to account for a margin of non-response. Special attention was paid to the variables needed to assess performance at the establishment level. Establishments with incomplete information were re-contacted in order to complete this information, whenever necessary. However, there were clear cases of low response. The following graph shows non-response rates for the sales variable, *d2,* by type of questionnaire. Non-response was kept below a 10% threshold.

**Sales Non-Response Rates**
Bulgaria, 2007



29.     Survey non-response was addressed by maximizing efforts to contact establishments that were first selected in the sample and by trying to keep a tight control over the process of substitutions. Up to 4 attempts were made to contact the establishment before it was replaced. However, non-response of the complete survey was faced in at some degree.

30.     As the following graph shows, the number of contacted establishments per realized interview was 2.28. This number is the result of two factors: explicit refusals to participate in the survey, as reflected by the rate of rejection (which includes rejections of the screener and the main survey) and the quality of the sample frame, as represented by the presence of ineligible units. The relatively low ratio of contacted establishments per realized interview (2.28) suggests that the main source of error in estimates in the Republic of Bulgaria may be selection bias and not frame inaccuracy.

**Rejections Rate and Contacts per Interview**
Bulgaria, 2007



31.     Details on rejections rates, eligibility rates, and item non-response are available at the level strata. This report summarizes these numbers to alert researchers of these issues when using the data and when making inferences. Item non-response, selection bias, and faulty sampling frames are not unique to the Republic of Bulgaria. All enterprise surveys suffer from these shortcomings but in very few cases they have been made explicit.

## References

Cochran, William G., Sampling Techniques, 1977.

Deaton, Angus, The Analysis of Household Surveys, 1998.

Levy, Paul S. and Stanley Lemeshow, Sampling of  Populations: Methods and Applications, 1999.

Lohr, Sharon L. Samping: Design and Techniques, 1999.

Scheaffer, Richard L.; Mendenhall, W.; Lyman, R., Elementary Survey Sampling, Fifth Edition, 1996