

The Georgia 2019 Micro-Enterprise Surveys Data Set

I. Introduction

This document provides additional information on the data collected in Georgia between July and November 2019. The objective of the Enterprise Survey is to gain an understanding of what firms experience in the private sector.

As part of its strategic goal of building a climate for investment, job creation, and sustainable growth, the World Bank has promoted improving the business environment as a key strategy for development, which has led to a systematic effort in collecting enterprise data across countries. The Enterprise Surveys (ES) are an ongoing World Bank project in collecting both objective data based on firms' experiences and enterprises' perception of the environment in which they operate.

The ES currently cover over 190,000 firms in 152 countries, of which 143 have been surveyed following the standard methodology. This allows for better comparisons across countries and across time. Data are used to create statistically significant business environment indicators that are comparable across countries. The ES are also used to build a panel of enterprise data that will make it possible to track changes in the business environment over time and allow, for example, impact assessments of reforms.

In addition to the standard ES that covers formal private sector firms with five or more employees, in Georgia, a separate survey is fielded simultaneously to collect data from micro-enterprises - formal firms with less than five employees.

This report outlines and describes the sampling design of the data, the data set structure as well as additional information that may be useful when using the data, such as information on non-response cases and the appropriate use of the weights.

II. Sampling Structure

The sample for 2019 Georgia Micro-enterprise was selected using stratified random sampling, following the methodology explained in the *Sampling Note*¹. Stratified random sampling² was preferred over simple random sampling for several reasons³:

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

¹ The complete text can be found at

http://www.enterprisesurveys.org/~media/GIAWB/EnterpriseSurveys/Documents/Methodology/Sampling_Note.pdf

² 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).

³ Cochran, W., 1977, pp. 89; Lohr, Sharon, 1999, pp. 95

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.

Two levels of stratification were used in this country: industry and region. The original sample design with specific information of the industries and regions chosen is described in Appendix C.

Industry stratification was designed in the way that follows: the universe was stratified into two manufacturing industries and three services industries- Food and Beverages (ISIC Rev. 3.1 code 15), Other Manufacturing (ISIC codes 16-37), Retail (ISIC code 52), Hospitality and Tourism (ISIC code 55) and Other Services (ISIC codes 45, 50, 51, 60-64, and 72).

There is no further breakdown by firm size for the micro-enterprise survey for sampling purpose and all firms are classified in to one size group, i.e., (1 to 4 employees).

Regional stratification was done across five regions: Tbilisi, East, Adjara, (Guria, Samegrelo, Zemo Svaneti) and Center. For the purposes of achieving the thresholds for representativeness, the ES indicators are calculated with some regions combined. In particular, Adjara and (Guria, Samegrelo, Zemo Svaneti) are combined.

III. Sampling implementation

Given the stratified design, sample frames containing a complete and updated list of establishments as well as information on all stratification variables (number of employees, industry, and region) are required to draw the sample. Great efforts were made to obtain the best source for these listings.

ACT Global was the main contractor that implemented the Georgia 2019 Micro-enterprise survey.

The sample frame used for this survey consisted of listings of establishments obtained from GeoStat.

Table 1: Georgia Micro-enterprise Sample Frame

		Food	Other Manufacturing	Retail	Hospitality and Tourism	Other Services	Grand Total
Tbilisi	Micro (1-4)	580	1307	4214	303	7892	27760
East	Micro (1-4)	254	297	1464	38	1176	
Adjara	Micro (1-4)	127	231	1207	127	1833	
Guria, Samegrelo, Zemo Svaneti	Micro (1-4)	144	127	780	57	873	
Center	Micro (1-4)	313	435	2003	120	1858	
		1418	2397	9668	645	13632	27760

Source: World Bank and GeoStat

Necessary measures were taken to ensure the quality of the frame; however, the sample frame was not immune to the typical problems found in establishment surveys: positive rates of non-eligibility, repetition, non-existent units, etc.

Given the impact that non-eligible units included in the sample universe may have on the results, adjustments may be needed when computing the appropriate weights for individual observations. The percentage of confirmed non-eligible units as a proportion of the total number of sampled establishments contacted for the survey was 30.0% (211 out of 704 establishments)⁴.

Breaking down by industry and size, the following sample targets were achieved (based on the sampling information):

Table 2: Achieved Interviews

		Food	Other Manufacturing	Retail	Hospitality and Tourism	Other Services	Grand Total
Tbilisi	Micro (1-4)	5	7	7	3	13	120
East	Micro (1-4)	9	7	3	3	3	
Adjara	Micro (1-4)	3	3	3	3	3	
Guria, Samegrelo, Zemo Svaneti	Micro (1-4)	3	3	3	3	3	
Center	Micro (1-4)	10	10	4	3	3	
		30	30	20	15	25	120

IV. Data Base Structure:

The structure of the data base reflects the fact that 2 different versions of the survey instrument were used for all registered establishments. Questionnaires have common questions (*core* module) and respectfully additional manufacturing- and services-specific questions. The eligible manufacturing industries have been surveyed using the **Manufacturing** questionnaire (includes the *core* module, plus manufacturing specific questions). Retail firms have been interviewed using the **Services** questionnaire (includes the *core* module plus retail specific questions) and the residual eligible services have been covered using the **Services** questionnaire

⁴ Based on out of target and ineligible contacts

(includes the *core* module). Each variation of the questionnaire is identified by the index variable, *a0*.

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 (some exceptions apply due to comparability reasons). Variable names preceded by the prefix "BM" indicate questions specific to Georgia and other countries in Europe and Central Asia 2018/2019 and Middle East and North Africa 2019, 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.

There are 2 establishment identifiers, *idstd* and *id*. The first is a global unique identifier. The second 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.

There are two levels of stratification: industry and region. Different combinations of these variables generate the strata cells for each industry/region combination. A distinction should be made between the variable *a4a* and *d1a2* (industry expressed as ISIC rev. 3.1 code). The former gives the establishment's classification into one of the chosen industry-strata based on the sample frame, whereas the latter gives the establishment's actual industry classification (four-digit code) based on the main activity at the time of the survey.

All of the following variables contain information from the sampling frame. They may not coincide with the reality of individual establishments as sample frames may contain inaccurate or outdated 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

- a6a*: coded using the same standard for small, medium, and large establishments as defined above.

- a4a*: coded following the stratification by sector as defined above.

The surveys were implemented following a 2-stage procedure. Typically, first a screener questionnaire is applied over the phone to determine eligibility and to make appointments. Then a face-to-face interview takes place with the Manager/Owner/Director of each establishment. However, sometimes the phone numbers were unavailable in the sample frame, and thus the enumerators applied the screeners in person. The variables *a4b* and *a6c* contain the industry and size of the establishment from the screener questionnaire.

Note that there are variables for size (*l1*, *l6* and *l8*) that reflect more accurately the reality of each establishment. Advanced users are advised to use these variables for analytical purposes. Variables *l1* (number of permanent full-time workers at the end of the last complete fiscal year), *l6* (number of full-time seasonal workers employed during last complete fiscal year) and *l8* (average length of employment of full-time temporary employees during last complete fiscal year) 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.

The firms interviewed had several fiscal years. Most firms had January to December 2018 as their last complete fiscal year. Variables *a20m* (starting month of last complete fiscal year) and *a20y* (last complete fiscal year) can be used to obtain the last complete fiscal year for each firm.

For questions pertaining to monetary amounts, the unit is the Georgian lari (GEL).

V. Universe Estimates

Universe estimates for the number of establishments in each cell in Georgia were produced for the strict, weak and median eligibility definitions described below. The estimates were the multiple of the relative eligible proportions.

For some establishments where contact was not successfully completed during the screening process (because the firm has moved, and it is not possible to locate the new location, for example), it is not possible to directly determine eligibility. Thus, different assumptions about the eligibility of establishments result in different adjustments to the universe cells and thus different sampling weights.

Three sets of assumptions on establishment eligibility are used to construct sample adjustments using the status code information.

Strict assumption: eligible establishments are only those for which it was possible to directly determine eligibility. The resulting weights are included in the variable *wstrict*.

$$\text{Strict eligibility} = (\text{Sum of the firms with codes } 1, 2, 3, 4, \& 16) / \text{Total}$$

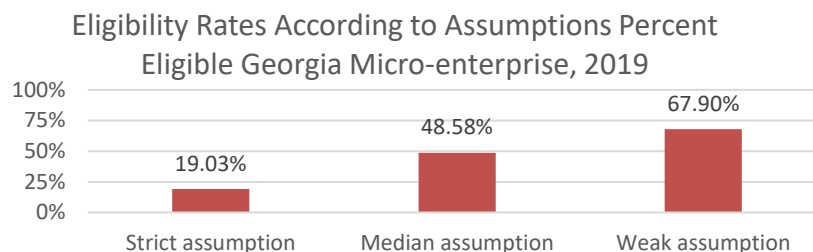
Median assumption: eligible establishments are those for which it was possible to directly determine eligibility and those that rejected the screener questionnaire, or an answering machine or fax was the only response. The resulting weights are included in the variable *wmedian*.

$$\text{Median eligibility} = (\text{Sum of the firms with codes } 1, 2, 3, 4, 16, 10, 11, \& 13) / \text{Total}$$

Weak assumption: in addition to the establishments included in points a and b, all establishments for which it was not possible to contact or that refused the screening questionnaire are assumed eligible. This definition includes as eligible 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. Under the weak assumption only observed non-eligible units are excluded from universe projections. The resulting weights are included in the variable *wweak*.

$$\text{Weak eligibility} = (\text{Sum of the firms with codes, } 1, 2, 3, 4, 16, 10, 11, 13, 91, 92, 93, 94, 12) / \text{Total}$$

The following graph shows the different eligibility rates calculated for firms in the sample frame under each set of assumptions.



Universe estimates for the number of establishments in each industry-region cell in Georgia were produced for the strict, weak and median eligibility definitions. Appendix B shows the universe estimates of the numbers of registered establishments that fit the criteria of the ES.

Once an accurate estimate of the universe cell projection was made, weights for the probability of selection were computed using the number of completed interviews for each cell.

VI. Weights

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 *pw* in Stata.)⁵

Special care was given to the correct computation of the weights. 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, no reply after having called in different days of the week and in different business hours, no tone in the phone line, answering machine, fax line⁶, 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 the cell. Once an accurate estimate of the universe cell (projections) was available, weights were computed using the number of completed interviews.

VII. Appropriate use of the weights

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.

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

⁵ This is equivalent to the weighted average of the estimates for each stratum, with weights equal to the population shares of each stratum.

⁶ For the surveys that implemented a screener over the phone.

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 have the advantage of providing an estimate that is independent of the sample design. This latter point may be quite relevant for the ES 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 use of weighted OLS for a common population coefficient.)⁷

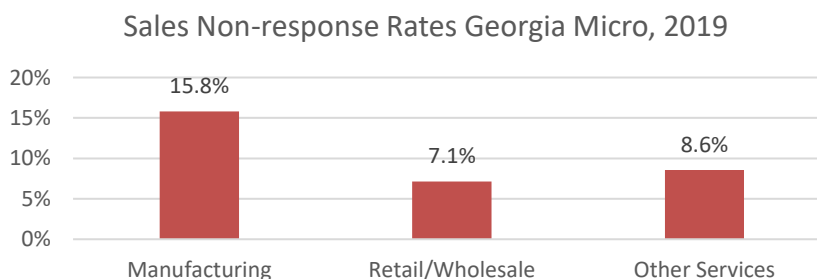
From a more general approach, if the regressions are descriptive of the population then weights should be used. The estimated model can be thought of as the relationship that would be expected if the whole population were observed.⁸ If the models are developed as structural relationships or behavioral models that may vary for different parts of the population, then, there is no reason to use weights.

VIII. Non-response

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.

Item non-response was addressed by two strategies:

- 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 (-8) as a different option from don't know (-9).
- b- 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 sector. Please, note that for this specific question, refusals were not separately identified from “Don't know” responses.



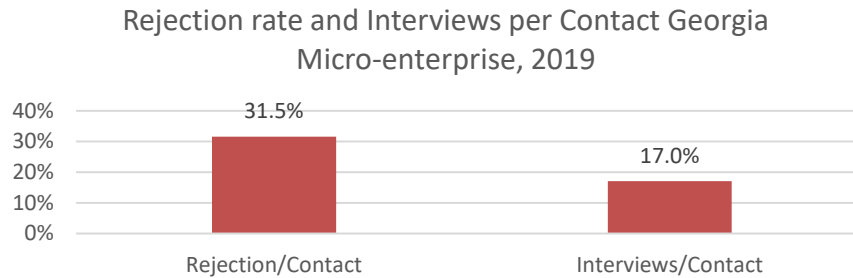
As the following graph shows, the number of interviews per contacted establishments was 0.17.⁹ This number is the result of two factors: explicit refusals to participate in the survey, as reflected by the rate of rejection (which includes rejections of the screener and the main survey)

⁷ 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.

⁸ 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.

⁹ The estimate is based on the total no. of firms contacted including ineligible establishments.

and the quality of the sample frame, as represented by the presence of ineligible units. The share of rejections per contact was 0.32.



Details on the rejection rate, eligibility rate, 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 Georgia. All enterprise surveys suffer from these shortcomings, but in very few cases they have been made explicit.

References:

- Cochran, William G., *Sampling Techniques*, New York, New York: John Wiley & Sons, 1977.
- Deaton, Angus, *The Analysis of Household Surveys*, Baltimore, Maryland: Johns Hopkins University Press, 1998.
- Levy, Paul S. and Stanley Lemeshow, *Sampling of Populations: Methods and Applications*, New York, New York: John Wiley & Sons, 1999.
- Lohr, Sharon L. *Sampling: Design and Techniques*, Boston, Massachusetts: Brooks/Cole, 1999.
- Scheaffer, Richard L.; Mendenhall, W.; Lyman, R., *Elementary Survey Sampling*, Fifth Edition, 1996.

Appendix A

Status Codes Micro-enterprise Survey :

15	Screening in process	14. In process (the establishment is being called/ is being contacted - previous to ask the screener)	15
134	Eligible	1. Eligible establishment (Correct name and address) 134 2. Eligible establishment (Different name but same address - the new firm/establishment bought the original firm/establishment) 0 3. Eligible establishment (Different name but same address - the firm/establishment changed its name) 0 4. Eligible establishment (Moved and traced) 0 16. Eligible establishment (Panel Firm - now less than five employees; this code applies only to panel firms.) 0	
208	Screener refusal	13. Refuses to answer the screener	208
211	Ineligible	5. The establishment has less than 5 permanent full time employees 0 616. The firm discontinued businesses - (Establishment went bankrupt) 2 618. The firm discontinued businesses - (Original establishment disappeared and is now a different firm) 2 619. The firm discontinued businesses - (Establishment was bought out by another firm) 2 620. The firm discontinued businesses - (It was impossible to determine for what reason) 140 621. The firm discontinued businesses - (Other) 6 71. Ineligible legal status: not a business, but private household 26 72. Ineligible legal status: cooperatives, non-profit organizations, etc. 26 8. Ineligible activity: Education, Agriculture, Finances, Government, etc. 7	
0	Out of Target	151. Out of target - outside the covered regions 0 152. Out of target - moved abroad 0 153. Out of target - Not registered with Statistical Authority 0 154. Out of target - establishment is HQ without production or sales of goods or services 0 155. Out of target - establishment was not in operation for the entirety of last fiscal year 0 156. Duplicated firm within the sample 0 157. Out of target - location that is not HQ and does not have financial statements prepared separately 0	
136	Unobtainable	91. No reply after having called in different days of the week and in different business hours 73 92. Line out of order 19 93. No tone 0 94. Phone number does not exist 5 10. Answering machine 0 11. Fax line- data line 0 12. Wrong address/ moved away and could not get the new references 39	
704	Total contacted		

Response Outcomes : Georgia Micro-enterprise Survey 2019 :

Target and totals	Sample target	120
	Sample target completion rate	100.0%
	Total contacts available in frame	27760
	Total contacts issued	890
	Total contacts contacted	704

Screening phase	Screening in process	15
	Eligibles	134
	Screeners refusal	208
	Ineligible + out of target	211
	Unobtainable	136
Interview phase (only if eligible)	Complete interviews without extra module	0
	Complete interviews with extra module	120
	Eligible in process + incomplete interviews	0
	Interview refusal	14

Percent breakdown (relative to total contacted)	Screening in process rate	2.1%
	Screeners refusal rate	29.5%
	Ineligible + out of target rate	30.0%
	Unobtainable rate	19.3%
	Interview conversion rate	17.0%
	Eligible in process + incomplete interviews rate	0.0%
	Interview refusal rate	2.0%

Appendix B: Universe Estimate Based on Sampling Weights

Strict Universe Estimates – Fresh:

		Food	Other Manufacturing	Retail	Hospitality and Tourism	Other Services	Grand Total
Tbilisi	Micro (1-4)	89	274	461	53	943	4430
East	Micro (1-4)	62	98	253	11	222	
Adjara	Micro (1-4)	27	67	181	31	301	
Guria, Samegrelo, Zemo Svaneti	Micro (1-4)	38	46	146	17	179	
Center	Micro (1-4)	75	141	339	33	344	
		291	626	1380	145	1988	4430

Median Universe Estimates – Fresh:

		Food	Other Manufacturing	Retail	Hospitality and Tourism	Other Services	Grand Total
Tbilisi	Micro (1-4)	284	640	2200	160	3625	13275
East	Micro (1-4)	134	157	825	22	583	
Adjara	Micro (1-4)	49	89	497	53	665	
Guria, Samegrelo, Zemo Svaneti	Micro (1-4)	76	67	440	33	433	
Center	Micro (1-4)	147	205	1007	61	822	
		691	1159	4969	329	6128	13275

Weak Universe Estimates – Fresh:

		Food	Other Manufacturing	Retail	Hospitality and Tourism	Other Services	Grand Total
Tbilisi	Micro (1-4)	373	918	2735	203	5066	18555
East	Micro (1-4)	174	222	1011	27	803	
Adjara	Micro (1-4)	85	168	814	88	1222	
Guria, Samegrelo, Zemo Svaneti	Micro (1-4)	103	100	566	43	626	
Center	Micro (1-4)	210	318	1354	84	1242	
		944	1726	6480	445	8960	18555

Appendix C: Original Sample Design
Original Sample Design

		Food	Other Manufacturing	Retail	Hospitality and Tourism	Other Services	Grand Total
Tbilisi	Micro (1-4)	5	7	7	3	13	120
East	Micro (1-4)	9	7	3	3	3	
Adjara	Micro (1-4)	3	3	3	3	3	
Guria, Samegrelo, Zemo Svaneti	Micro (1-4)	3	3	3	3	3	
Center	Micro (1-4)	10	10	4	3	3	
		30	30	20	15	25	120