



SAVE THE DATE

ONLINE AFRICAN SCHOOL *on Migration Statistics*

2nd edition


8–10 June 2021



Estimation Methods for Migration Statistics

Some examples of methodological challenges, based on the experiences from the work on the 3rd edition of the report Labour Migration Statistics in Africa (LMSA-3)

Hans Pettersson

- Estimation – the final step to get the data in order for tabulation and analysis in a sample survey
- Use the statistical measures (averages, proportions, totals) from a group of units (households, individuals) to make a “good conjecture” about these measures in a larger group (“population”)
- Group characteristic (average household income in the sample)

- Estimate of population characteristic (average household income in the country)

- We can confidently say something about the characteristic for a large group of units (all households in the country), although we are missing data for most of the units.
- Intuitive:
- The larger the group from which we have data, the better is the estimate.
- Put in another way: The less units (data) missing from the survey, the better is the estimate

- Data (units) can be missing for two reasons.
 - Missing-by-design: data (units) are "missing-at-random"
 - Missing due to (poor) execution of a survey (imperfect sample frame, data collection problems, nonresponse)

- Use the data you have for a group of units (part of the population).
- Compensate for the units missing-by-design. This is done by applying sampling weights to the data. The weights are the inverted value of the inclusion probabilities. (All units should have a nonzero probability of being included in the sample).
- Compensate for units missing for other reasons (noncoverage, nonresponse, not surveyed,). This can be done by various methods like imputation and "modeling".

“Estimation” in LMSA

- In the presentation I will use a wider definition of estimation
- The “survey” is the labour migration statistics system for LMSA
- The survey population is all countries in all the years 2010 to 2019 (units= country*year)
- The characteristics (variables) are number of migrants (totals), LFPR (proportions) and many other.
- No sample, it is a census. But a lot of missing data.
- Estimation: compensate for the missing data (imputation, models)

Presentation

1. Briefly about sample survey methodology. I will go through some of the central concepts. Good as a background when we discuss quality issues with the survey conducted for LMSA
2. Estimation methods in LMSA-3
3. Experiences from LMSA-3
4. LMSA-4 and beyond

- Initial planning: State the purpose; define/clarify concepts; agree on content.....

- Decide on mode of data collection
- Questionnaire design
- Interviewer recruitment and training
- Field supervision
- Data processing procedures



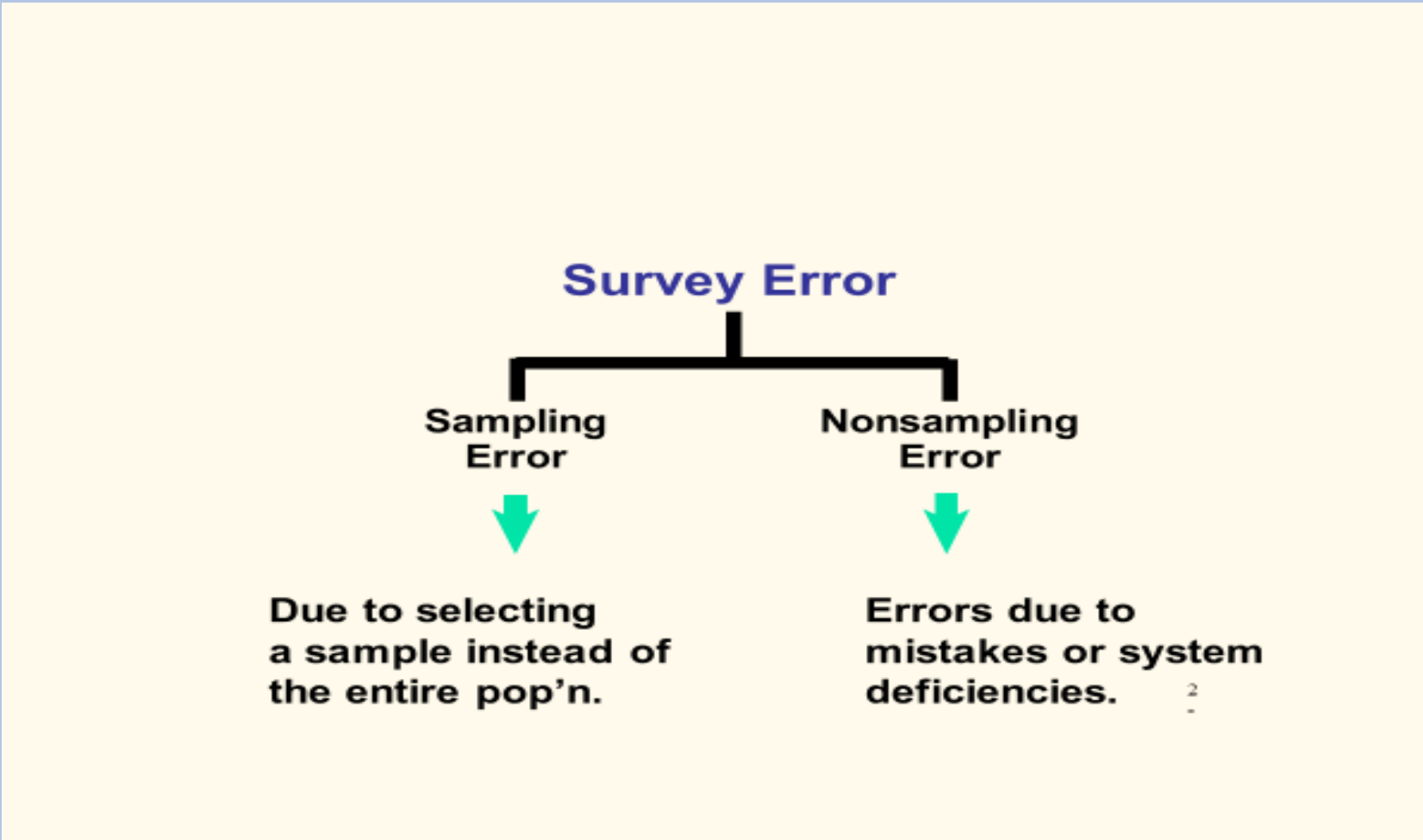
Non-sampling errors

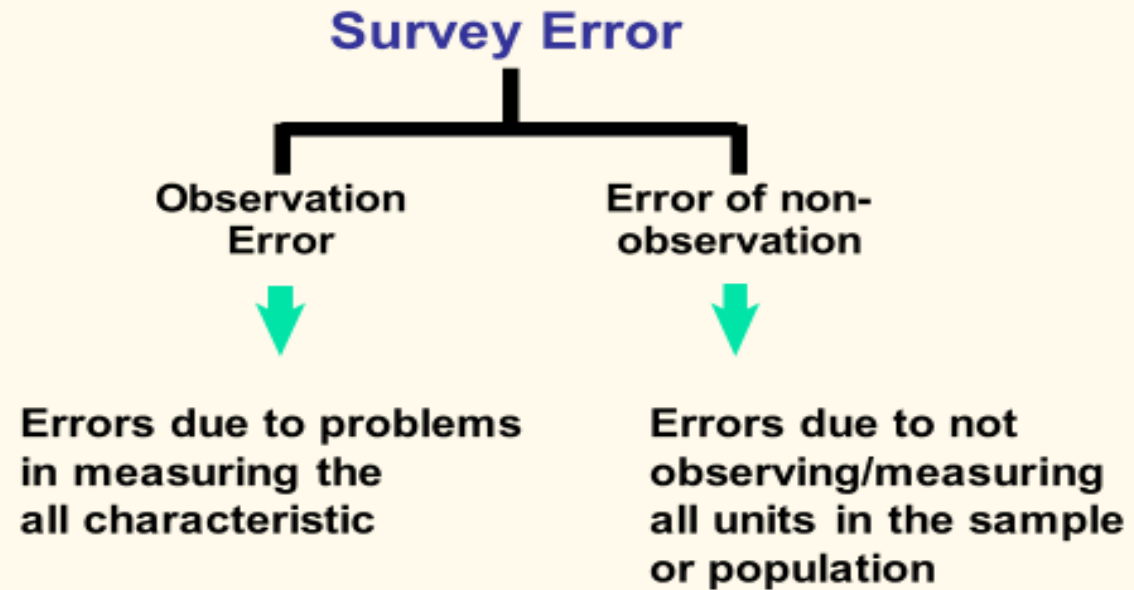
- Sample design



Sampling errors

- *Sampling errors* is one type of error that occurs in surveys. Error due to sampling, not all units are surveyed (missing data).
- Frame errors (missing data due to imperfect sampling frame)
- Non-response errors (missing data due to no observation/measurement)
- Measurement errors (error in data due to error in observation/measurement)
- Processing errors (error in results due to mistakes in the data processing)





Errors of non-observation:

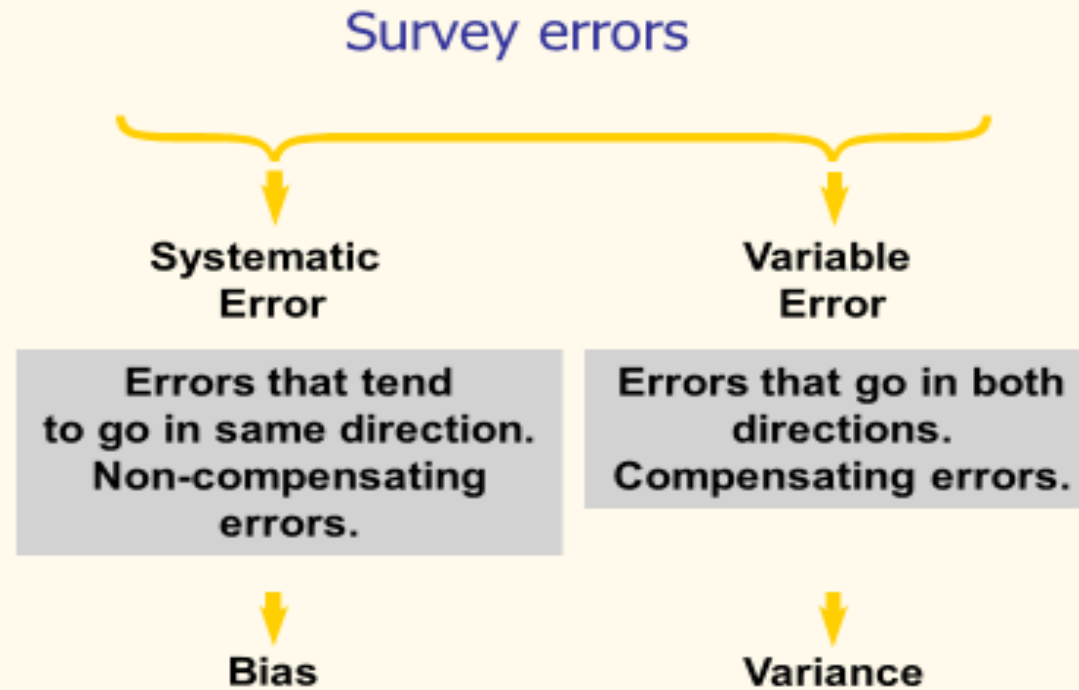
Frame errors (coverage errors). Errors because the frame doesn't cover all units in target population. Households not in the sampling frame have no chance of being included in the sample (undercoverage). Are they different from covered households?

Non-response errors. No data collected from some households. Are they different from responding households?

Sampling errors. Only a small part of the population is surveyed. Error occurs because the sample households (as a group) differ somewhat from the population.

Observation errors:

Measurement errors. Measure not correct (questionnaire)? Do we measure what we want to measure (validity)?




Sampling error

Variance of the mean =

$$Var(\bar{x}) = S_{\bar{x}}^2 = \frac{S_x^2}{n} = \frac{1}{n(n-1)} \sum_1^n (x_i - \bar{x})^2$$

Standard error = $s.e(\bar{x}) = \frac{S_x}{\sqrt{n}}$

$$s.e(\bar{x}) = \frac{S_x}{\sqrt{n}}$$


The standard error of a statistic depends on the **sample size**. The larger the sample size, the smaller the standard error.

As the sample size increases, the standard error decreases, and for a census or complete enumeration (where $n=M$), the standard error is zero.

Confidence intervals indicate the uncertainty:

95 % confidence interval: $\bar{x} \pm 1.96 \cdot s.e(\bar{x})$

“The confidence interval covers the true population mean with 95% confidence”

“If we repeat the sampling 1000 times and calculate the mean for each sample, the mean will be within the interval 950 times out of 1000”

From Cambodia Socioeconomic Survey: Monthly household expenditure per capita (,000 Riel)

Phnom Penh: $212 \pm 1.96 \cdot 6$ [200 – 224]

Other urban: $166 \pm 1.96 \cdot 6$ [154 – 178]

The confidence intervals **do not overlap**, so we conclude that there is a significant difference in household expenditure between PP and other urban

Estimates of totals from multistage samples

A *total* could be estimated from the sample by the estimator :

$$\hat{Y} = \sum_{h=1}^L \sum_{i=1}^{n_h} \sum_{j=1}^{m_{hi}} w_{hij} y_{hij}$$

Sampling weight

- The sample implementation and the data collection are never perfect
- Missing values due to unit nonresponse, noncoverage
 - If "missing at random" – adjust the sampling weights with the factor n/n^* where n^* is the number of responding households and n the number of sampled households
 - If not "missing at random" – post-stratification, calibration
- Missing values due to item nonresponse
 - Imputation (ex: "nearest neighbour", mean value)

- Imputation can be seen as a form of estimation.
- Mean imputation is like sampling weight adjustment (or sampling weight adjustment is like mean imputation).
- Nearest neighbour is estimation by a model: the neighbour household is similar to the missing (nonresponding) household

Part 2: Estimation methods in LSMA-3

- The “survey” is the labour migration statistics system for LMSA
- The survey population is all countries in all the years 2010 to 2019 (units= country*year)
- The characteristics (variables) are number of migrants (totals), LFPR (proportions) and many other.
- No sample, it is a census. But a lot of missing data.
- Estimation: compensate for the missing data (imputation, models)

- We are leaving the sample survey context and look at estimation with the presence of missing data in a more general setting where data may be collected from other sources (census, admin data, international data bases, accounting systems)
- Specifically, we look at the situation where we want to present annual statistics but the data are only collected intermittently, every 5 or 10 years, and data will come from different sources.
- The data set has a lot of "holes" that need to be filled before the data set can be used for estimation of population characteristics

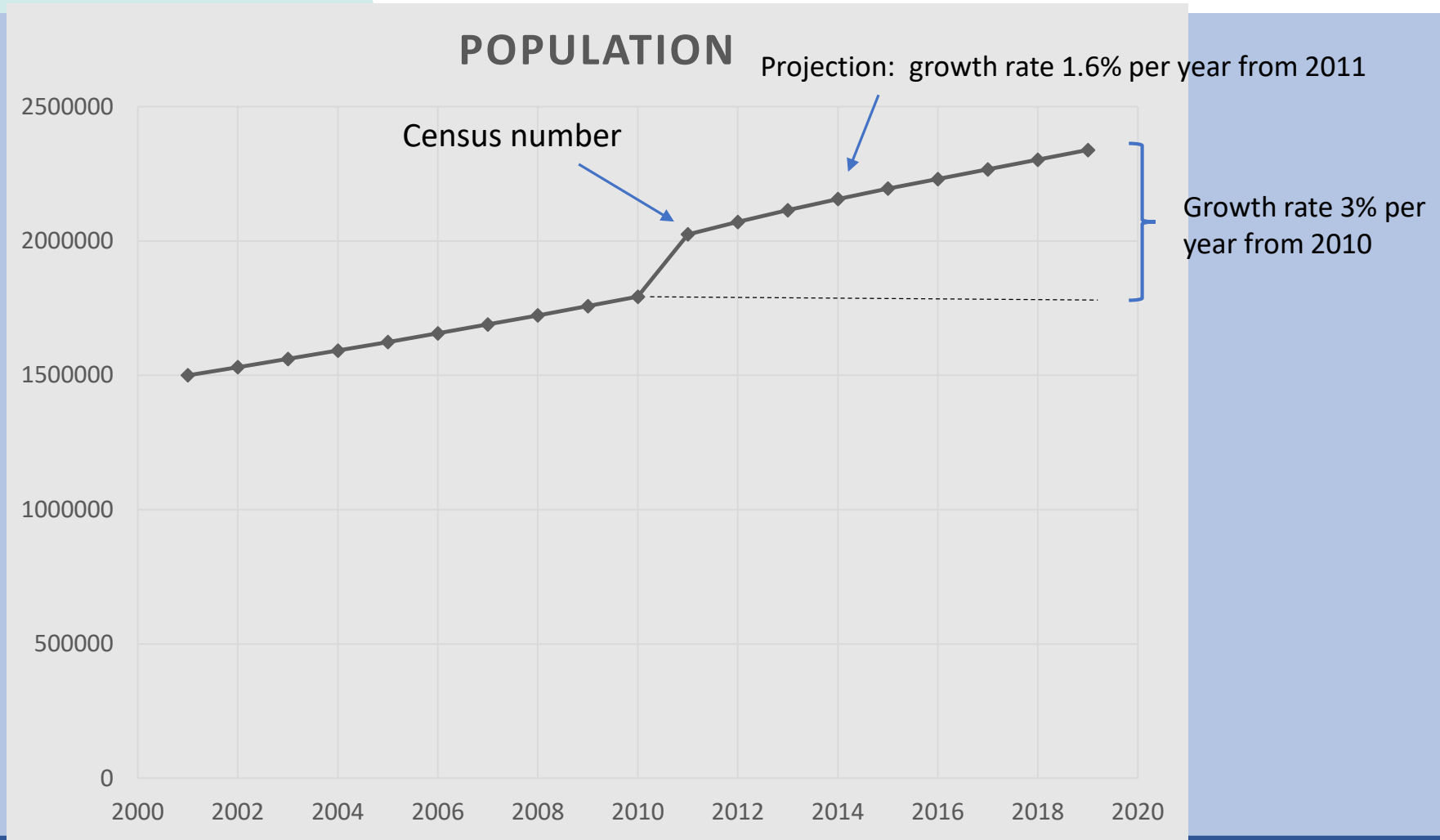
- **Missing data in time series**
- Imputation by using the mean= Interpolation – set the missing value equal to the mean of the previous and the following value
- Imputation by using the trend = If several subsequent missing values! - Interpolation/extrapolation – set the missing values equal to the trend values
- Imputation by using a model = use a stable relationship with another variable and assume this model is valid for the missing data point

- **Imputation by using a model**
- Example: We have data on number of people 15+ of age for the years 2010 – 2019. And we have data on number of people in the labour force for the years 2010 and 2012 (from LFS).
- Model= the proportion of 15+ in the labour force is constant. Therefore, calculate the mean proportion based on the two values for 2010 and 2012. Impute the mean proportion in all missing data points

- **Imputation by using a model**
- Another example: We have data on number of people 15+ of age for the years 2010 – 2019. And we have data on number of people in the labour force for the years 2010 and 2019 (from LFS). Can we assume a another model?
- Model= the proportion of 15+ in the labour force is slowly changing over the years. Therefore, calculate the trend in the proportion based on the two values for 2010 and 2019. Impute the trend values of the proportion in all missing data points 2011-2018
- Beware of the uncertainty in the trend due to sampling errors!

- **Imputation by using a model**
- Another example: We have data on number of people 15+ of age for the years 2010 – 2019. But we have no data on number of people in the labour force.
- Model= Borrow data from other "similar" countries. The proportion of 15+ in the labour force is equal to the mean value for a group of (neighbouring) countries that resemble our country in important aspects. Therefore, calculate the mean of the proportion for the other countries. Impute the trend values of the proportion in all missing data points 2010-2019.
- Somewhat dangerous. Can we really rely on "similarity"? But it is done.

- Situation: A country has conducted a census in 2001 and prepared demographic projections for number of people 15+ of age for 2002 to 2030. In 2011 a new census is carried out and it is found that the projection for 2011 from census 2001 falls short of the 2011 census result. That means that the projection for 2010 also is below the true value.
- We want to present data for the period 2010 – 2019. What to do?
 - Keep the inaccurate value for 2010 and use the census value for 2011 and new projections for 2012-2019?
 - Calculate a trend based on the census values 2001 and 2011 and impute the trend value for 2010?



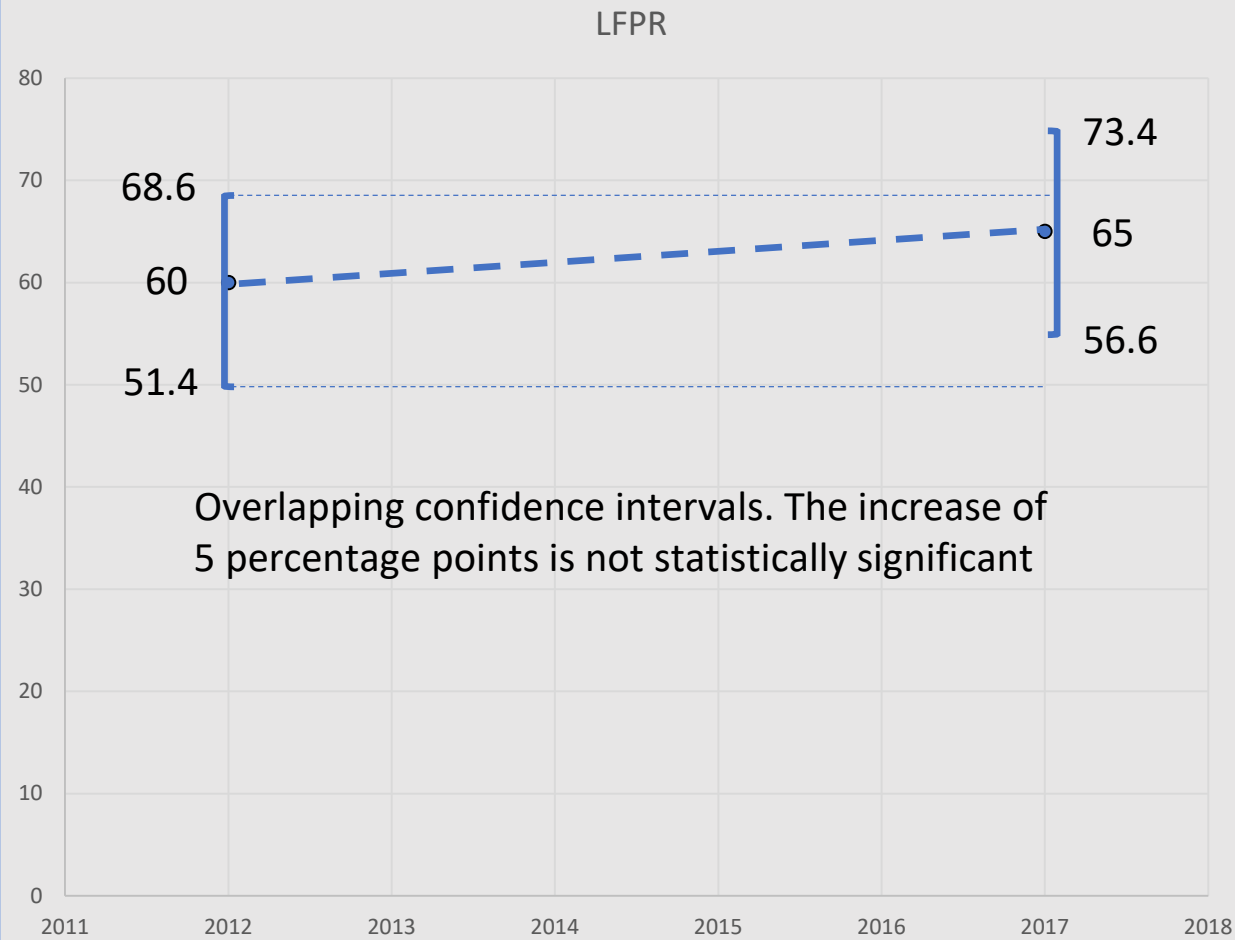
- Situation: A country has conducted two labour force surveys, in 2012 and 2017.
- We want to present data on the number of people in the labour force for the period 2010 – 2019. What to do?
 - Calculate the trend 2012 – 2017 based on the survey estimates and extrapolate the trend to 2010-2011 and 2018-2019?
- Beware of the uncertainty due to sampling error

- A labour force survey has a sample size of 15,000 individuals 15-64. The migrant population 15-64 in the country constitutes 1% of the population 15-64.
- We therefore expect to get around 150 migrants in the sample (if oversampling of migrants has not been designed).
- The estimate of LFPR is 60%. Then the standard error (s.e.) will be:

- $$s.e. = \sqrt{\frac{60 \cdot (100 - 60) \cdot 1.2}{150}} = 4.4$$

- Confidence interval: $60 \pm 1.96 * 4.4$ [51.4 -- 68.6]

- The next labour force survey (5 years later) gets the estimate of LFPR at 65%. Is there a trend?
- The standard error (s.e.) will be:
- $s.e. = \sqrt{\frac{65 \cdot (100 - 65) \cdot 1.2}{150}} = 4.3$
- Confidence interval: $65 \pm 1.96 * 4.3$ [51.6 -- 68.4]



- What is the confidence interval for the trend? (expressed as the annual increment)
- Trend: $(65-60)/5 = 1$ percentage point per year
- The standard error (s.e.) for the trend will be:
- $s.e. = \sqrt{(4.4^2 + 4.3^2)/5} = 2.8$
- Confidence interval for trend: $1 \pm 1.96 * 2.8$ [-4.5 ---- 6.5]

- We calculated the standard error and confidence interval of the LFPR for migrants. Due to the small sample we got a very wide confidence interval. What is the interval for the general population where the sample size is 15,000 individuals?
- The estimate of LFPR is 60%. Then the standard error (s.e.) will be:
- $s.e. = \sqrt{\frac{60 \cdot (100 - 60) \cdot 1.2}{15000}} = 0.44$
- Confidence interval: $60 \pm 1.96 * 0.44$ [59.1 -- 60.9]
- Large sample >>>> narrow confidence interval

- confidence interval for the trend? (expressed as the annual increment)
- Trend: $(65-60)/5 = 1$ percentage point per year
- The standard error (s.e.) for the trend will be:
- $s.e. = \sqrt{(0.44^2 + 0.43^2)/5} = 0.28$
- Confidence interval for trend: $1 \pm 1.96 * 0.28$ [0.4 ---- 1.6]

- The estimates of number of migrants in the labour force (“working migrants”) often have very wide confidence intervals because of small sample sizes.
- As a consequence, the estimate of the LFPR will also have a very wide confidence interval.
- Also, some countries may not have carried out a survey to measure migrant LFPR
- What can be done instead of using very poor or non-existing estimates?
- Imputation by using a model:
 - $LFPR(migr) = LFPR(pop) \times C$
- C is a model constant calculated for a group of countries

Part 3: Experiences from LMSA-3

- Describe levels and trends of labour participation among general population and migrant population.
- Analyse/study LFPR and employment status among subgroups of the general population and migrant population. (Can be seen as a part of a bigger effort by AU to develop a system of migration statistics).
- Users of the report findings and the data:
 - “Describers” – many (Governments, AU, RECs, int organizations)
 - “Researchers/ modelers” – not so many (lack of accurate detail)

- **Existing data and reports**
- National: Census, surveys, admin data
- International and regional: UN, WB, ILO, AU.....
- LMSA-3 Strategy: get the data directly from the countries. Countries are the data owners. Country NSOs were the primary data providers

- **Data collection:**
- Country data on labour participation were collected by an e-mail survey where the country NSO (or other entity) responded to a standardized questionnaire. The questionnaire *International Labour Migration Questionnaire* (ILMQ) has been designed within the Joint Labour Migration Programme (JLMP).
- 17 Excel tables
- 47 countries responded

- **“Estimation”:**
- Imputations (mean, model), inter- and extrapolations to fill in missing data.

- Picture of survey design

- Picture of survey design - alternative

Resident population, by sex and labour force participation - total and youth (total & migrant population)

	TOTAL POPULATION						TOTAL LABOUR FORCE						TOTAL EMPLOYED					
	All ages			Working age (15+)			Labour Force (15+)			Youth (15-35)			Employed (15+)			Youth (15-35)		
	Total	Men	Women	Total	Men	Women	Total	Men	Women	Total	Men	Women	Total	Men	Women	Total	Men	Women
2010	30 712 974	14 988 862	15 724 112	15 095 387	7 154 581	7 940 806	13 464 118	6 412 272	7 051 847	8 719 382	4 087 131	4 632 251	13 009 003	6 261 095	6 747 909	8 443 771	4 007 121	4 436 650
2011	31 013 084	15 135 325	15 877 759	15 385 944	7 292 293	8 093 651												
2012	31 316 126	15 440 662	15 875 464	15 682 093	7 661 890	8 020 203	12 700 751	6 323 438	6 377 313	7 763 682	3 786 467	3 977 215	12 431 618	6 209 880	6 221 738	7 538 930	3 697 741	3 841 189
2013	34 092 136	16 528 950	17 563 185	16 678 789	7 873 699	8 805 091	14 293 470	6 990 393	7 303 077	8 688 431	4 264 058	4 424 373	14 015 112	6 888 255	7 126 857	8 462 979	4 181 753	4 281 226
2014	34 634 650	16 897 849	17 736 801	18 585 420	9 058 841	9 526 579												
2015	35 491 000	17 338 500	18 152 500	18 620 800	8 725 200	9 895 600												
2016	37 673 800	18 190 178	19 483 622	19 345 900	9 012 508	10 333 393	16 190 463	7 765 531	8 424 932	9 826 842	4 703 111	5 123 730	15 275 372	7 461 270	7 814 102	9 090 863	4 458 032	4 632 831

	MIGRANT POPULATION						MIGRANT LABOUR FORCE						TOTAL EMPLOYED MIGRANTS					
	MIGRANTS (all ages)			Working age Migrants (15+)			Total Migrant Labour Force (15+)			Youth (15-35)			Total Employed Migrants (15+)			Youth (15-35)		
	Total	Men	Women	Total	Men	Women	Total	Men	Women	Total	Men	Women	Total	Men	Women	Total	Men	Women
2010																		
2011																		
2012	270,051	110,229	159,822	230,403	99,030	131,373	197,312	90,367	106,945	103,809	36,874	66,935	192,537	89,130	103,407	99	35,994	63,397
2013																		
2014	337 696	160 950	176 746	232 940	106 548	126 392												
2015																		
2016																		

- **Data from the countries**
- Genuine, primary data from census and surveys
- Modeled data (projections, extrapolations, pro-rata calculations)

- **Projections done by NSO's**
- **The general population:**
 - Some countries have "official" demographic projections.
 - Many other countries have done simple ad hoc projections.
 - Some countries have no projections, only one or two data points.
- **The migrant population:**
 - Some countries have done simple ad hoc projections, usually assuming the same growth rate as what was used for the projection of the general population

- Under certain assumptions, it is possible to calculate the average “age” of the genuine data used for an estimate.
- We can look at an estimate - an average or a total - for Africa calculated for the year 2018. We assume that the estimate is entirely based on data from the most recent population census for each country.
- The oldest census data used for the survey is from 2003 (C.A.R.). The age of the genuine data is in this case 15 years.
- For a country that conducted the census in 2014, the data age is 4 years.

- A calculation shows that the average age of the genuine data points for a 2018 estimate for Africa is 7.5 years.
- Put in another way: the average length of the projections (extrapolations) up to year 2018 is 7.5 years. That's a rather long projection, the uncertainty is substantial.
- The projections (along with imputations, inter- and extrapolations) will not capture the year-to-year variation in the migrant stock. A substantial change in numbers from one year to another will not be reflected in the statistics. (Flow data needed).

- **Ways to check the data: Data confrontation**
 - UN DESA: population, migrant stock
 - WB: remittancies
 - ILO: labour force, employed, migrant workers
 - ECOWAS: migrant stock, migrant workers

■ A real challenge

- Lots of methodological considerations and quality issues. Responsibility ultimately with the countries. Varying quality.
- Surveying migrants more difficult than surveying households and individuals in households (household budget, labour force, health,).

Part 4: LMSA-4 and beyond

■ Improving the survey:

- Better editing. Use ILO models for data validation.
- Revise the ILMQ. Abandon the idea of getting flow data from the countries in ILMQ (in the short term)
- Put more effort in checking the data for the very big countries. They have a big impact on the aggregates (regions, the continent). Small error in big country is worse than big error in small country

■ Development work:

- Migrant surveys needed (module in LFS or separate survey)
- Use of admin data (border controls, immigration office....) to estimate flows

- **Final words:**
- LMSA is a work in progress
- Also a learning process for country NSOs and STATAFRIC
- Too many "holes" in the data sets. Labour migration statistics in Africa – a story about missing data?