

Introduction: World Population Prospects 2021 Upgrade - Towards a more open and reproducible WPP

United Nations Expert Group Meeting on Methods for the World Population Prospects 2021 and Beyond



Session 1: Monday 6 April 2020

- Evolution of previous and current revisions for WPP estimates:
 - Reconstruct internally coherent time series of population and demographic components since 1950
 - 5x5 cohort-component data model + post-facto 1x1 interpolations for selected outputs
 - Incorporate the demographic impact of the HIV/AIDS epidemic
 - More systematic compilation of empirical data since 1950 into internal SQL database (DemoData)
 - Greater interactive internal validation of WPP estimates with country data
 - More systematic documentation of data used/considered and generation of metadata (country "short notes")
 - Use of probabilistic projection methods for fertility and mortality components in addition to deterministic projection scenarios

- Challenges and pending issues/requests for WPP estimates:
 - More transparent and open access to all underlying country empirical data not just the metadata documentation of the data sources and estimation methods used
 -> SQL database + front-ends (DemoData and DataPortal)
 - Greater documentation and explanations of the various methods used to derive demographic estimates for each demographic components and the reconciliation with population estimates -> WPP method protocol
 - More replicable set of criteria/decisions used to derive estimates that can be more easily communicated and potentially replicated, or updated as new data become available
 - Better capacity to use annual time series (upon data availability and reliability)
 - Better capacity to use single age data (upon availability and reliability)
 - More efficient capacity to update and revise country estimates as new data become available and with greater resilience to staff turnover and staff shortage

Guidelines for Accurate and Transparent Health Estimates Reporting (<u>GATHER</u>)

ltem #	Checklist item (completed in blue text with hyperlinks, pending in red)	Status						
Objecti	Dbjectives and funding							
1	Define the indicator(s), populations (including age, sex, and geographic entities), and time period(s) for which estimates were made.							
2	2 List the funding sources for the work.							
Data In	ata Inputs							
For al	data inputs from multiple sources that are synthesized as part of the study:							
3	Describe how the data were identified and how the data were accessed.	 ✓ 						
4	Specify the inclusion and exclusion criteria. Identify all ad-hoc exclusions. 20							
5	Provide information on all included data sources and their main characteristics. For each data source used, report reference information or contact name/institution,							
	population represented, data collection method, year(s) of data collection, sex and age range, diagnostic criteria or measurement method, and sample size, as relevant.							
6	Identify and describe any categories of input data that have potentially important biases (e.g., based on characteristics listed in item 5). 2021							
For do	For data inputs that contribute to the analysis but were not synthesized as part of the study:							
7	Describe and give sources for any other data inputs.	 ✓ 						
For all data inputs:								
8	8 Provide all data inputs in a file format from which data can be efficiently extracted (e.g., a spreadsheet rather than a PDF), including all relevant meta-data listed in item 5.							
	For any data inputs that cannot be shared because of ethical or legal reasons, such as third-party ownership, provide a contact name or the name of the institution that							
	retains the right to the data.							
Data ar								
9	Provide a conceptual overview of the data analysis method. A diagram may be helpful.	•						
10	Provide a detailed description of all steps of the analysis, including mathematical formulae. This description should cover, as relevant, data cleaning, data pre-processing,	2021						
	Describe how condidate models were evaluated and how the final model(s).	2021						
11	Describe now candidate models were evaluated and now the final model(s) were selected.	2021						
12	Provide the results of an evaluation of model performance, if done, as well as the results of any relevant sensitivity analysis.	2021						
13	Describe methods for calculating uncertainty of the estimates. State which sources of uncertainty were, and were not, accounted for in the uncertainty analysis.	2021						
14 Bosulto	14 State now analytic or statistical source code used to generate estimates can be accessed. 202							
	and Discussion Dravide nublished estimates in a file format from which data can be officiently avtracted							
15	Penert a quantitative measure of the upportainty of the estimates (e.g. upportainty intervals)	┼┷┤						
17	Interpret results in light of existing evidence. If undeting a previous set of estimates, describe the reasons for changes in estimates	+						
10	Discuss limitations of the estimates, include a discussion of any modelling assumptions or data limitations that affect interpretation of the estimates.	+						
18	Discuss limitations of the estimates. Include a discussion of any modelling assumptions or data limitations that affect interpretation of the estimates.							

- Long term goals for WPP 2025 and beyond:
 - More open access to both methods, underlying data and analytical steps
 - More collaborative with potential greater country engagement/interactions
- Short term goals for WPP 2021 (planned release May/June 2021)
 - Upgrade production system into 1x1
 - Streamline/harmonize steps used to prepare country data and WPP estimates
 - Provide access to both WPP estimates and underlying empirical data for key demographic indicators -> <u>Data Portal</u> + <u>Demo Data</u> + <u>Data Archive</u>
 - More GATHER compliant
- Constraints and limitations
 - Limited time between upgrade of methods, IT infrastructure and data system
 - Fixed/limited resources (staff, skills, computing, budget, etc.)

DataCatalog, DataArchive and DemoData

- <u>DataCatalog</u>: comprehensive inventory for each country of primary data sources (censuses, demographic surveys, etc.) providing data on demographic processes (fertility, mortality, population structure and dynamics, marital status and family planning) for all countries and areas, as of March 2020, more than 6750 entries
- <u>DataArchive</u>: a virtual repository of documents, tabular datasets and reports (potentially) for each data source, as of March 2020, more than 27,000 files
- <u>DemoData</u>: SQL database to store in a structured and standardized way empirical data and demographic estimates (with meta-information) on population, fertility, mortality and migration data from as many sources as possible

🥶 UNI	TED	NATION	S Page	Catalog & Archive dation Division				
Home Catal	a Arthur	Collected Date						
 Ountries 	O Major Areas							
NDUA (HE) X		1945	End W	an 🖬 Solect Data Process Types	Res	et 27		
Data Catalo	g country or are	na of primary data sou	rces (cer	sures, surveys, etc.) that collected data on demographic topics	A topol	Ver Color O	Q Search o	in Norne or S
country !	17.01			NAME	SHORT NAME	R.C. 848100		
India	Census	Census		India 1951 Census	1951 Census	1951	1951	1951
India	Survey	Survey	0	India 1953-1954 National Sample Survey (rural)	1953-1954 NSS	1953-1954	1963	1954
India	Survey	Survey	D	India 1957-1958 National Sample Survey (rural)	1957-1958 NSS	1957-1958	1967	1958
India	Survey	Survey	D	India 1960-1961 National Sample Survey (urban)	1950-1951 NSS	1960-1961	1960	1961
India	Census	Census		India 1961 Census	1951 Census	1961	1961	1961
India	Survey	Survey	•	India 1965-1966 National Sample Survey	1955-1956 NSS	1965-1966	1965	1956
India	Survey	Survey	•	India 1970-1971 National Family Planning Survey	1970-1971 NFPS	1970-1971	1970	1971
		Census		India 1971 Census	1971 Census	1971	1971	1971
India	Census							
india India	Census Eurvey	Survey	•	India 1972 Fertility Survey	1972 FS	1972	1972	1972
India India India	Census Burvey Survey	Burvey Survey	6 6	India 1972 Fertility Burvey India 1979 Survey on Infant and Child Mortality	1972 F0 1979 SICM	1972 1979	1972 1979	1972 1979

UNITED NA	ATIONS DemoData: D Population Divi	ataBrowser sion					
The Data Browser provides access to the online database of empirical demographic data and selected tabulations (DemoData) on population by age and sex, fertility and mortality data compiled and used by the UN Population Division to derive estimates and projections from censuses, population and vital registers, surveys and other sources - going back to the 1950's depending upon availability.							
<u></u>	<u>jenni</u>						
Custom searches based on indicators and country selections	Browse & filter data series and tabulations	Export data series and tabulations	Input data or make edits to existing data				
based on indicators and country selections	data series and tabulations	data series and tabulations	or make edits to existin data				

Data Portal: dissemination of estimates/projections and empirical data



Country data -> provisional WPP estimates -> final by demographic component for WPP2021



WPP estimation process for each country/area





An open suite of R packages and functions

- New 1x1 cohort-component population projection computational engine designed to work with a standard set of 1x1 inputs and outputs
 - R implementation for deterministic projections/simulations
 - C implementation for probabilistic projections/simulations
- [TBD: B3-type of robust time trend modelling for TFR/ASFR, and adult mortality]
- <u>DDSQLTools</u>: set of functions to query <u>DemoData</u> SQL database with API
- <u>DemoTools</u>: set of functions to evaluate, transform and adjust counts or rates
- <u>DDM</u> and <u>FertEstR</u>: set of functions to evaluate and adjust mortality and fertility data
- <u>Ungroup</u> and <u>MortalityLaws</u>: mortality graduation and extension at older ages
- <u>SVDcomp</u>: new expanded set of model life tables (including HIV and ART)
- <u>Calibrated Splines</u>: graduation of fertility age patterns
- [TBD: additional packages/functions to operationalize WPP method protocol]
- <u>popReconstruct</u>: probabilistic demographic estimation and population reconstruction
- <u>BayesTFR</u>, <u>BayesLife/BayesLifeHIV</u>, <u>MortCast</u>, <u>BayesPop</u>: probabilistic projections

Time table

- Feb-August 2020: upgrade of existing IT/data system
- 6-8 April 2020: EGM about WPP2021 and beyond
- Sept.-Dec. 2020: update of country data and WPP estimates using new 1x1 approach
- Jan. 2021: QA data checking and error detections
- Feb. 2021: probabilistic projections
- March 2021: computation of all derived outputs
- April-May 2021: preparation of reports + Data Portal & media launch
- June 2021: public launch

Strategy = Five year plan: 3 tiers of countries into 3 steps...

 Following the UN Statistics Division (UNSD) classification for the reliability of population estimates and vital statistics officially reported by national statistical authorities and used since the 1950s for the Demographic Yearbook (DYB)

"Reliability of data: *Reliable mid-year population estimates are those that are based on a complete census (or a sample survey) and have been adjusted by a continuous population register or on the basis of the calculated balance of births, deaths and migration" (DYB 2017, p. 51)*

- Two key dimensions for the most recent population estimates:
 - <u>Nature of the base measurement of the population</u> (register, census, etc.) and <u>Time elapsed</u> since the last measurement
 - <u>Method and quality of time adjustment</u> by which the base figure was brought up to date (registers, vital registration (VR) + migrations or not, availability of multiple censuses with intercensal period ≤ 10-15 years, etc.)

Note: both census in last 10 year and VR coverage (100% birth registration, and 80% death registration) are used for SDG indicator 17.19.2

Based on these UNSD DYB criteria for the reliability of official population estimates for the period 1950-2020, countries/areas can be grouped into 3 tiers:

- Tier 1: countries/areas (n<50) with population registers or regular censuses and VR ≥ 99% since 1950
- Tier 2: countries/areas (n<100) with regular censuses and VR ≥ 60% since 1950 (not included in Tier 1)
- Tier 3: countries/areas (n<90) with no regular censuses and/or no VR or < 60% since 1950

[note the cut/off for VR is only for illustrative purpose]

WPP revision	Tier 1: Countries with population registers or regular censuses and VR ≥ 99%	Tier 2: Countries with regular censuses and VR ≥ 60%	Tier 3: Countries with no regular censuses and/or no VR or < 60%
2021	<u>Probabilistic</u> estimates	Deterministic estimates and reconciliation informed by statistical modeling of time trends	Deterministic estimates and reconciliation informed by statistical modeling of time trends
2023	<u>Probabilistic</u> estimates	<u>Probabilistic</u> estimates (conditional on staff/resources availability)	Deterministic estimates and reconciliation informed by statistical modeling of time trends
2025	<u>Probabilistic</u> estimates	Probabilistic estimates (conditional on staff/resources availability)	Probabilistic estimates (conditional on staff/resources availability)