

1 **A consistent dataset for the net income distribution for 190**
2 **countries and aggregated to 32 geographical regions ~~and the~~**
3 **~~world~~ from 1958-2015**
4

5 Kanishka B. Narayan¹, Brian C. O'Neill¹, Stephanie Waldhoff¹ and Claudia Tebaldi¹

6 ¹Joint Global Change Research Institute (JGCRI), Pacific Northwest National Lab (PNNL)

7 *Correspondence to:* Kanishka B. Narayan (kanishka.narayan@pnnl.gov)

8 **Abstract**

9 Data on income distributions within and across countries are becoming increasingly important to
10 inform analysis of income inequality and to understand the distributional consequences of climate
11 change. While datasets on income distribution collected from household surveys are available for
12 multiple countries, these datasets often do not represent the same income concept and therefore
13 make comparisons across countries, over time and across datasets difficult. Here, we present a
14 consistent dataset of income distributions across 190 countries from 1958 to 2015 measured in
15 terms of net income. We complement the observed values in this dataset with values imputed from
16 a summary measure of the income distribution, specifically the GINI coefficient. For the
17 imputation, we use a recently developed principal components-based approach that shows an
18 excellent fit to data on income distributions compared to other approaches. We also present another
19 version of this dataset aggregated from the country level to 32 geographical regions ~~and the world~~
20 ~~as a whole.~~ ~~Our aggregation method takes into account both within-country and across-country~~
21 ~~income inequality when aggregating to the regional level.~~ ~~Our dataset is developed for the purpose~~
22 ~~of calibrating models such as Integrated human-Earth system models with detailed data on income~~
23 ~~distributions aggregate regions.~~ This dataset will enable more robust analysis of income
24 distribution at multiple scales.

25 **1. Introduction**

26 Data on income distributions are important to understand trends in global and regional income
27 inequality. These data are also routinely used to train models that project income distributions
28 into the future (Fujimori et al., 2020; Hallegatte & Rozenberg, 2017; Hughes et al., 2009;
29 Hughes, 2019; Soergel et al., 2021; Van der Mensbrugge, 2015). In the climate literature, long-
30 term projections of within-country income distribution have been used to inform analyses of how
31 the impacts of climate change may affect inequality and poverty (Hallegatte & Rozenberg, 2017;
32 Jafino et al., 2020). Income distribution data are generally collected through national and local
33 household surveys. The most prominent sources of national-level income distribution data are
34 the datasets presented by the World Bank through the PovCal tool (Bank, 2015) and the income
35 distribution datasets available from the Luxembourg Income Study (LIS) (Ravallion, 2015;
36 Smeeding & Grodner, 2000). Both these datasets present useful time series of income
37 distribution for income groups such as deciles, based on multiple household surveys.

1 While these datasets have been widely used, they are subject to certain limitations. The definition
2 of income in these datasets is often not the same, making comparisons across countries and
3 datasets difficult (Smeeding & Latner, 2015). For example, the PovCal dataset has mixed
4 observations for net income and consumption for the same country in different years. Such
5 inconsistencies can occur because the underlying surveys in different years might have been
6 conducted to measure different income concepts. The two income concepts that these data tend
7 to use are:

8 i) *Post tax income or disposable income or net income* - This measure is defined as employee
9 income plus income from firms (self-employment) plus income from rentals (excluding any
10 payments), property income (these are generally capital gains and include dividends) plus current
11 transfers received (these include insurance benefits, employer contributions) less transfers paid
12 (taxes paid and employee contributions). This is the concept of income recommended by the
13 Canberra group for the international comparison of incomes (Europe, 2011).

14 ii) *Consumption* - This measure is the sum of food consumption plus non-food consumption plus
15 durable goods purchases (expenditure value minus cost of repairs) plus housing expenditures
16 (rent, mortgage payments) less any payments made (taxes, loan payments, asset purchases, etc).
17 This is the concept of income recommended by Deaton & Zaidi (2002) for welfare measurement.

18 Temporal and spatial coverage of the data are another issue. The LIS dataset provides consistent
19 data on the net income distribution. However, these data are only available for 50 countries from
20 1980 to 2016. The PovCal dataset provides data for a considerably higher number of countries
21 (165) compared to the LIS. However, the data are a combination of net income and
22 consumption-based observations (net income distribution data for 73 countries and consumption
23 distribution data for 118 countries).

24 Previous studies that have made use of these datasets for analysis or for modelling income
25 distributions have treated these income concepts as interchangeable (Rao et al., 2019; Sauer et
26 al., 2020). Moreover, for countries where no survey data on income distributions are available,
27 studies have used simple methods such as using a summary measure of income distribution such
28 as the GINI coefficient in combination with a parametric functional form such as a lognormal
29 distribution to impute the within country or within-region income distribution (Fujimori et al.,
30 2020; Rao et al., 2019; Shorrocks & Wan, 2008; Soergel et al., 2021).

31 There have been efforts to generate consistent datasets of the income distribution. However,
32 these efforts have been limited to local or regional data. For example, Frank (2009) generated a
33 consistent dataset of income distribution metrics for a single income concept for the fifty US
34 states. That particular study builds on previous studies that have compiled data for the US states
35 (Piketty & Saez, 2003). At the national level, there have been some efforts to produce
36 standardized datasets of income inequality, but they have generally been limited to summary
37 metrics of the income distribution such as the GINI coefficient (Babones & Alvarez-Rivadulla,
38 2007). Lanker and Milanovic (2013) developed a useful time series of income deciles across
39 countries which is a combination of data from the LIS, PovCal and other sources. However, this

1 dataset is still a combination of different income concepts and has a limited temporal time series
2 (the dataset only extends to the year 2013).

3 In this study we present a consistent dataset on national income distributions that represents a
4 single income concept namely, net income. This dataset contains a total 8522 data points of
5 income deciles across 190 countries. This dataset is constructed by first choosing net income
6 decile data observations from all available sources for all available countries (1191
7 observations). For countries that only have consumption distribution data, we impute the net
8 income distribution using a regression-based approach (494 observations). For countries and
9 years where no data on income distribution is available, we impute income deciles using the
10 GINI coefficient combined with a principal component analysis (PCA) based method that
11 provides a better fit to data than existing methods (6837 observations). This PCA-based method
12 was recently developed as a non-parametric approach to projecting income distribution (Narayan
13 et al., 2023). We note that the PCA based imputation provides the maximum number of
14 observations in the dataset.

15 One intended use of this dataset is to initialize income distribution variables in the Global
16 Change Analysis Model (GCAM) (Calvin et al., 2019). GCAM is a global, integrated model of
17 the energy, land, water, climate, and socioeconomic systems that produces projections for several
18 economic, climatological and physical systems variables for 32 geopolitical regions. Hence, we
19 also present income distributions for these 32 aggregated regions in addition to the 190 countries.
20 We use an aggregation method that takes into account cross-country inequality within a
21 region in addition to within-country inequality.
22

23 This dataset can be used to train projection models for income distribution across different scales
24 and, given the consistent income concept represented, can also be used to understand trends
25 within and across countries and regions. While these data are generated to enable modelling of
26 the income distributions in GCAM, they can be used to train any model for projecting income
27 distributions.

28 **2. Dataset construction**

29 We explain our approach for the dataset construction in detail in the sections below. To
30 summarize, we used the following steps:

- 31 a. We first identified observations by country and year of net income deciles from all
32 available datasets (LIS, PovCal, and individual research studies). In doing so, we
33 prioritized the LIS dataset over all other datasets given its high data quality on the net
34 income distribution. Our selection process is explained in **section 2.1 and 2.2** below.
- 35 b. For countries/years in which there were no net income data, but consumption data was
36 available, the net income distribution was imputed from the consumption distribution
37 using a regression-based approach. This is explained in **section 2.3**.
- 38 c. Where there were no net income or consumption data, but the GINI coefficient, a
39 summary metric of the income distribution, i.e., was available, we imputed the net

income distribution from the summary measure using a PCA-based approach. This is explained in **section 2.4**.

Note that point c. in the above yields the maximum number of data points in our final dataset. Table 1 below summarizes the coverage of our dataset-

Type of data	country-year observations
Original data on net income (<i>Explained in section 2.2</i>)	1198
Imputed based on original data on consumption (<i>Explained in section 2.3</i>)	4394
Imputed from GINI coefficient (using PCA algorithm) (<i>Explained in section 2.4</i>)	6857
Total	8522

Table 1: Summary of data points covered in our data set

2.1 Literature review and data selection from available household survey data

We first conducted a literature review to identify sources of national-level data on income distributions for as many countries as possible. There are three main datasets available, from the Luxembourg Income Study (LIS)(Ravallion, 2015; Smeeding & Grodner, 2000) the World Bank (whose data on income distributions are available through the PovCalNet tool) (Bank, 2015) and UNU WIDER (which compiles data from different sources including the LIS, PovCal and other research studies) (WIDER, 2008). Each dataset contains income distribution data for different income concepts such as net income and consumption, based on nationally representative surveys that may also represent sub-groups of the population (e.g., Urban vs Rural). These data are sometimes supplemented with data from research studies, and they use different equivalence scales to convert from household to per capita income. We first evaluated data availability for net income deciles based on these criteria (income concept, scale, temporal coverage, and spatial coverage).

In Table 1, we summarize these datasets differentiated by these criteria. Since the UNU WIDER dataset is a compilation of data sources (i.e., LIS, PovCal or others), we also identified

Formatted: Font: Bold, Italic

Formatted: Font: (Default) Times New Roman, 12 pt

Formatted: Normal, No bullets or numbering

Formatted: Font: Bold, Italic

Formatted: Font: Bold, Italic

Formatted: Font: Not Bold, Italic

Formatted: Font: Not Bold, Italic

- 1 the number of observations (country-year) in the UNU WIDER data derived from each source.
- 2 **SI Table 1** of this document summarizes some of the other studies which were used in the
- 3 collection of data for the UNU WIDER database. We are primarily interested in decile-level
- 4 income distributions derived from household surveys.

Source	Income concept	Scale of survey	Countries	Years (range)	Observations (n)
Luxemburg income study	Net income	National	50	1980-2016	347
	Consumption	National	25	1980-2016	209
PovCalNet	Net Income	National	73	1981-2018	1644
		Urban/Rural	3	1981-2018	37
	Consumption	National	114	1981-2018	2341
		Urban/Rural	3	1983-2018	54
UNU WIDER	Net Income	National	163	1979-2017	1707 347 from LIS 533 from other sources 827 from PovCal
		Urban	22	1961-2018	315 51 from PovCal 264 from other sources
		Rural	20	1950-2017	215 3 from PovCal 212 from other sources
	Consumption	National	66	1973-2018	1030 116 from LIS 779 from PovCal 135 from other sources

		Urban	5	1975-2017	52 45 from PovCal 7 from research studies
		Rural	5	1975-2017	50 46 from PovCal 4 from research studies

1

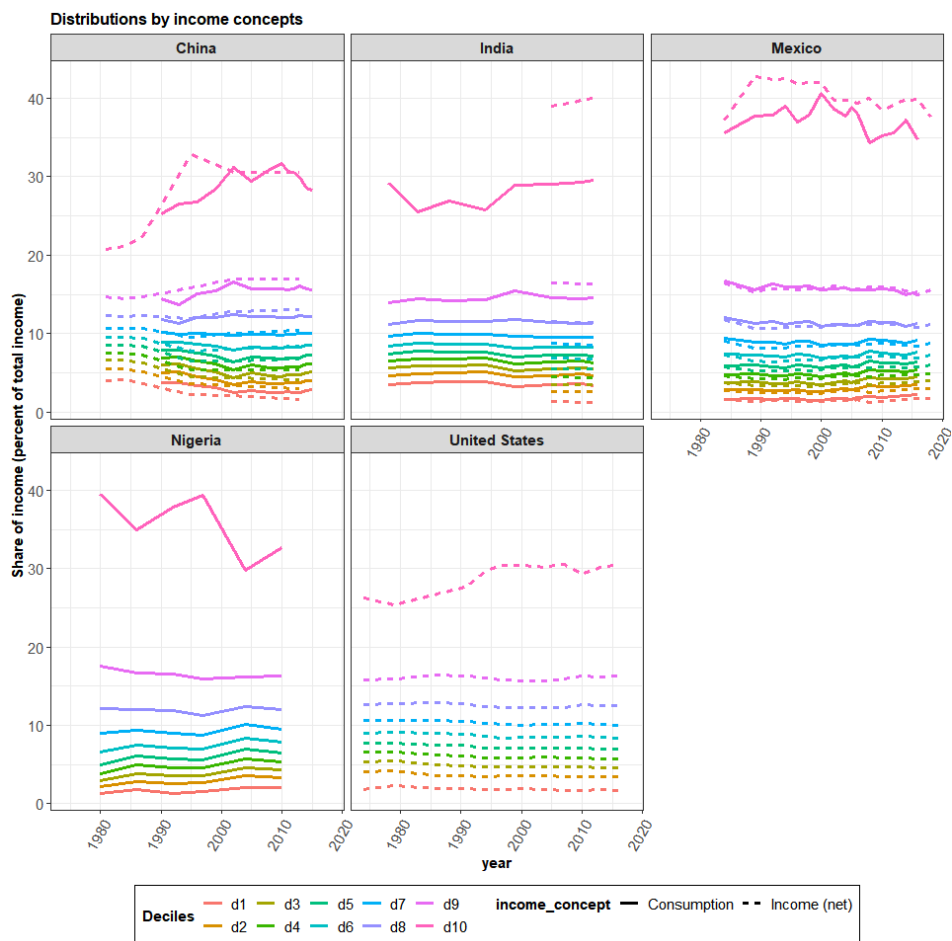
2 *Table 2†: Summary of coverage by data source*

3 We also evaluated access to microdata (i.e., underlying household-level data from household
4 surveys) for each of these datasets, since detailed microdata allows us to validate and understand
5 how the different income distributions for different income concepts were arrived at. Of all
6 datasets evaluated, we found that the LIS database has the most access to microdata via the
7 METIS tool (<https://www.lisdatacenter.org/frontend>).

8 The PovCal database maintained by the World Bank has the highest coverage geographically and
9 temporally in terms of observations. PovCal uses the disposable income data from LIS for high-
10 and middle-income countries and uses household survey data for consumption and disposable
11 income for low-income countries. The scales of the surveys are mostly national other than India,
12 China, and Indonesia where distribution data from separate rural and urban surveys are available.
13 Mean and median values of the income concepts are available in 2011 USD PPP converted using
14 country-specific conversion factors.

15 PovCal sometimes combines data of different types even within countries, e.g., for China,
16 PovCal uses income data in early years up to 1990 and then switches to consumption data.
17 Moreover, the micro-data for PovCal are not readily available.

18 UNU WIDER releases quality scores of individual datasets. It classifies the LIS database as
19 “High quality”, due especially to the availability of metadata, and classifies the PovCal dataset as
20 “Average quality”. Figure 1 below shows the income distributions by deciles for different
21 countries for different income concepts from the UNU-WIDER dataset.



1
 2 Figure 1: Income distributions across countries (facets) for different deciles (color) for different income concepts (line types) from
 3 the UNU WIDER dataset

4 **2.2 Selection of income concept and scheme for selection of data points**

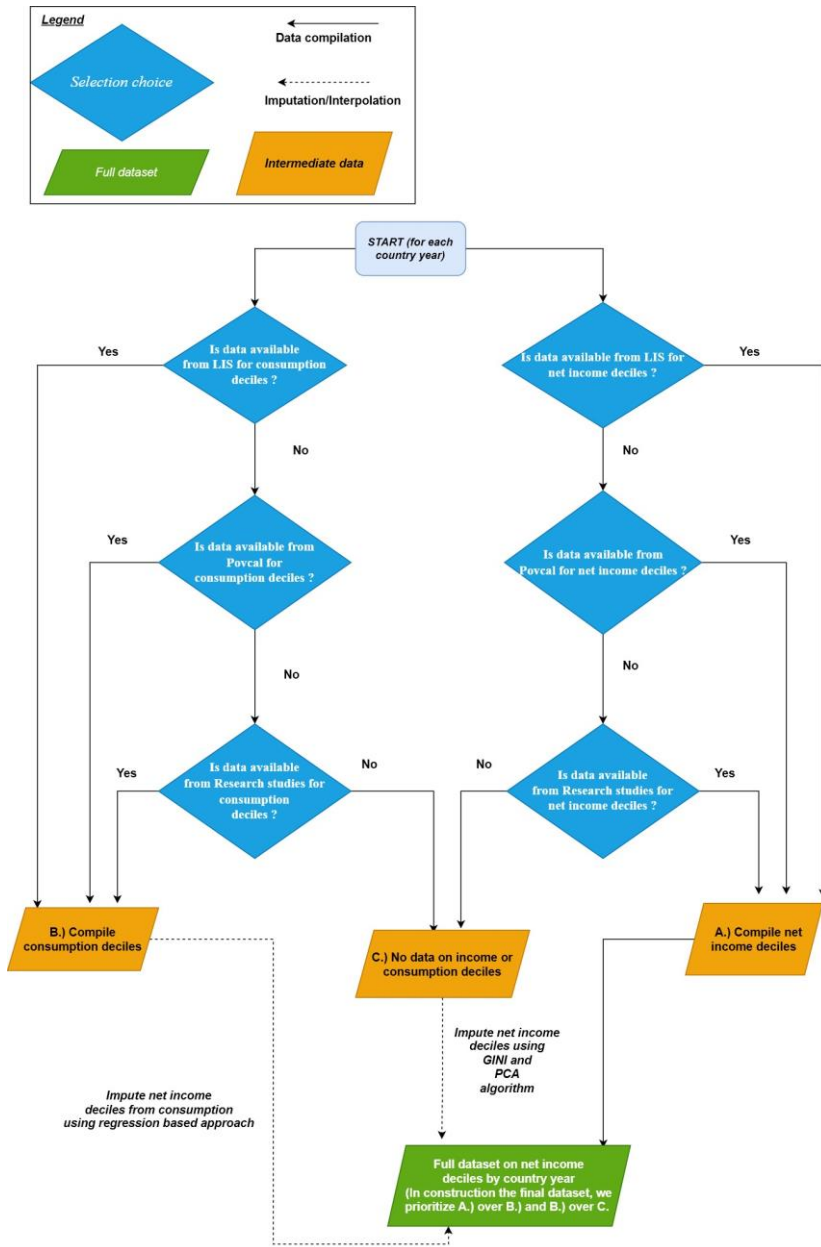
5 We construct a dataset that represents solely net income based on the same per-capita
 6 equivalence scale. The per capita equivalence scale is calculated using total household income
 7 divided by the household size assuming equal sharing of income. Our process, summarized in
 8 Figure 2, improves upon other attempts to construct income distribution datasets from different
 9 sources (Rao & Min, 2018; Rao et al., 2019), since the previous studies used the income concept
 10 from different datasets interchangeably. We primarily select observations for net income deciles
 11 across countries from the LIS, given the high quality of data available from that dataset. We

1 begin by compiling separate datasets of the income distribution for net income and consumption.
2 In construction of both these datasets, we prioritize data points from the LIS. If no data were
3 available from the LIS for a country-year, we selected an observation of net income or
4 consumption from the PovCal database. Finally, if data were not available from that database, we
5 rely on income distribution data from other research studies available from the UNU WIDER
6 database. Note that when selecting values across multiple research studies we select values based
7 on the rating assigned by the UNU WIDER database to the studies. All data are selected for the
8 equivalence scale applied in the WIDER dataset, in which household income was converted to
9 per capita units by dividing the household income by the household size assuming equal sharing
10 of income.

11 Thus, at this stage, we compiled two different data sets, one that represents net income
12 distribution across countries across time and another that represents consumption for the same
13 countries. Now, we prioritize the selection of net income distribution values over consumption
14 for each country-year.

15 Where data are only available for the consumption distribution, we convert the consumption data
16 to net income data (as explained in section 2.3 below), using a regression approach to generate a
17 harmonized dataset of net income deciles. Where necessary, we aggregated data sources across
18 different survey scales (urban vs. rural) using a population-weighted average.

19 Figure 2 summarizes our data selection approach.



1
2 Figure 2: Summary of data selection approach for each country, year observation

1

2 Based on the above, we evaluated data coverage for the 229 countries we are targeting. The
3 geographical boundaries of the 32 GCAM regions are defined based on these 229 countries
4 (countries with their corresponding regions are listed in **SI Table 2**). We identified observations
5 after the selection above for four categories, namely countries where we have net income data for
6 at least one year, countries where we had both net-income and consumption distribution data for
7 at least one year (in case of these countries we selected the net income distribution value for
8 deciles), countries where we had only consumption data, and countries where there were no data
9 (these countries only had data on aggregate measures of inequality such as the GINI coefficient
10 but no data on income deciles). Table 3 below summarizes the number of observations (country
11 years) by category of data.

12

13

Data availability (for at least 1 year) by income concept	Number of countries	Notes on use
Net income only	33	Use net income share data.
Both net income and consumption	54	Use net income share data.
Consumption only	83	Imputed income shares to be calculated (See section 2.3)
No decile data available but GINI is available	14	Impute deciles based on GINI coefficient (See section 2.4)
No data available	39	Drop from data set (section 5)
Total	229	

14

15 *Table 23: Summary of data availability by income concept.*

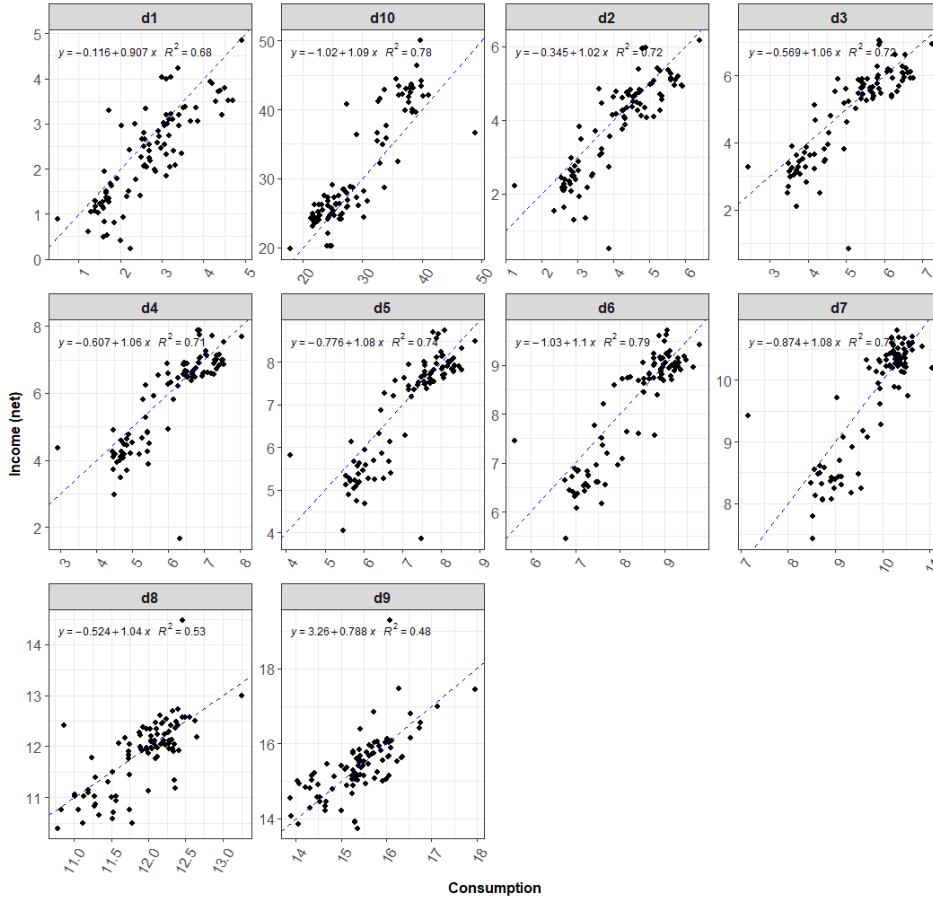
16 **2.3 Imputing net income shares using consumption shares**

17 Using data for countries which had both income and consumption distribution observations for
18 the same years (n=257, across 54 countries where each of which have data for ten deciles of
19 consumption and the ten deciles of net income), we constructed linear regression equations based
20 on a training dataset (n=148) for each decile to impute the net income shares using the
21 consumption shares of the income distribution (Figure 3). The highest R squared value was
22 observed for the fifth, sixth, seventh and tenth deciles d10 of 0.74 and the lowest R squared value

- 1 was observed for d9 of 0.37. We calculate values for 9 deciles d1-d8 and d10 and the re-calculate
- 2 d9 as the residual. We have verified that all imputed decile values add up to 1.

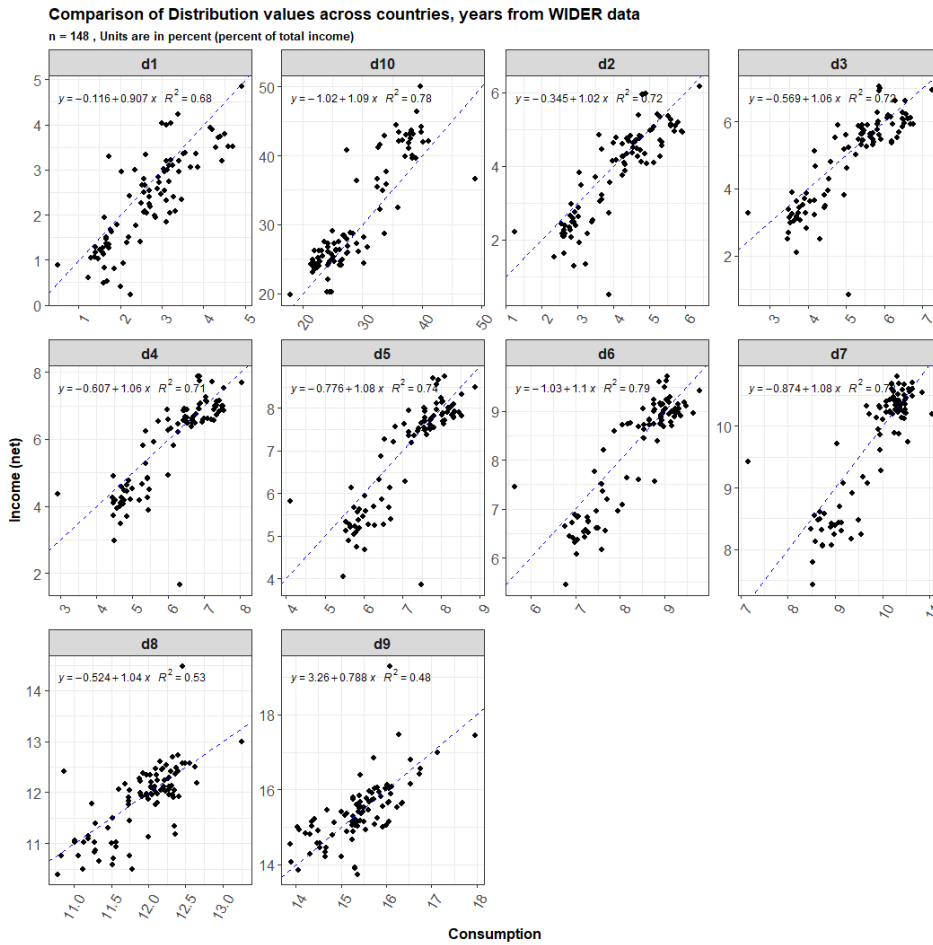
1

Comparison of Distribution values across countries, years from WIDER data
n = 148, Units are in percent (percent of total income)



2

1



2

Figure 3: Consumption distribution deciles (x axis) compared to Net income distribution deciles (y axis) across all country-year observations. Dashed lines show the 1:1 linear relationship. Solid line is the used regression line. Only observations for half the dataset are selected (Pre 2004) for the plot

3 Consumption distribution deciles are converted into net income deciles using the equation (1)
 4 below,

$$5 \quad D_{netincome_{n,r,t}} = Coeff_n * D_{consumption_{n,r,t}} + Intercept_n \quad (1)$$

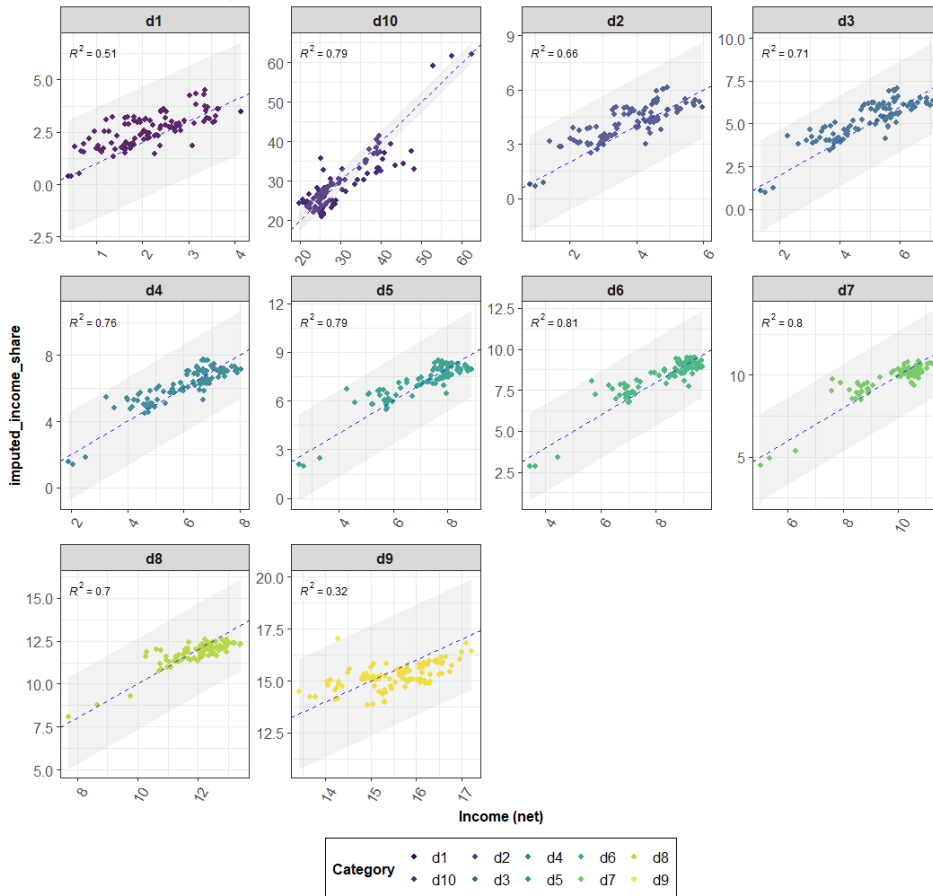
6 where,

1 D is the share of consumption or income in a particular decile between 0 and 100,
2 Coeff is the coefficient applied to each decile parameterized using a linear regression,
3 documented in Table 4,
4 Intercept is derived from linear regressions run for each decile, documented in Table 4,
5 n is the decile ranging from 1 to 10, and
6 r, t are the region and the time step respectively.

7 **Validation of our approach-** We then verified the performance of our regression on a testing
8 dataset (Figure 4). We note that the R squared values in our testing dataset is similar to our
9 training dataset and we also noted that the imputed values are within a 5 percent confidence
10 interval of actual values. To validate our imputation method we calculated errors (Imputed shares
11 - actual shares) for our testing dataset (n=109). We compared the error by decile for the dataset
12 (See SI Figure 1). The mean error across deciles is generally within half a percent across all
13 years. There are larger differences for the year 2011, where we had very few observations. We
14 have also verified that all imputed decile values add up to 1.

Formatted: Font: Bold, Italic

Comparison of Distribution values (using imputed income shares) across countries, years from the testing dataset
 n = 109 , Units are in percent (percent of total income)



- 1
- 2 [Figure 4: Comparison of actual vs imputed values on our testing dataset. Different deciles are shown as facets and we also show](#)
- 3 [the confidence interval. All imputed values are found to be within a 5 % CI of the original values except d10 where a few](#)
- 4 [observations are outside the range.](#)
- 5 [We note that this imputation method is applied to a small subset of observations \(494\) out of the](#)
- 6 [total observations in our dataset 8522. We also acknowledge that this method is simple and](#)
- 7 [should be improved upon in future updates/analysis.](#)

Decile	Intercept	Coefficient	Adjusted R ²
1	-0.02	0.81	0.5
2	-0.39	1.00	0.64

3	-0.65	1.06	0.69
4	-0.76	1.08	0.72
5	-0.91	1.10	0.75
6	-1.12	1.12	0.78
7	-1.10	1.10	0.78
8	-0.74	1.06	0.66
9	4.81	0.69	0.29
10	-1.39	1.11	0.75

1 *Table 34: Summary of coefficients and intercepts by decile used by Equation 1. These are fit*
2 *based on 257 data points.*

3 ~~The final dataset therefore includes 8422 observations based on distributions of consumption or~~
4 ~~net income across 170 countries spanning the time period 1958-2018.~~

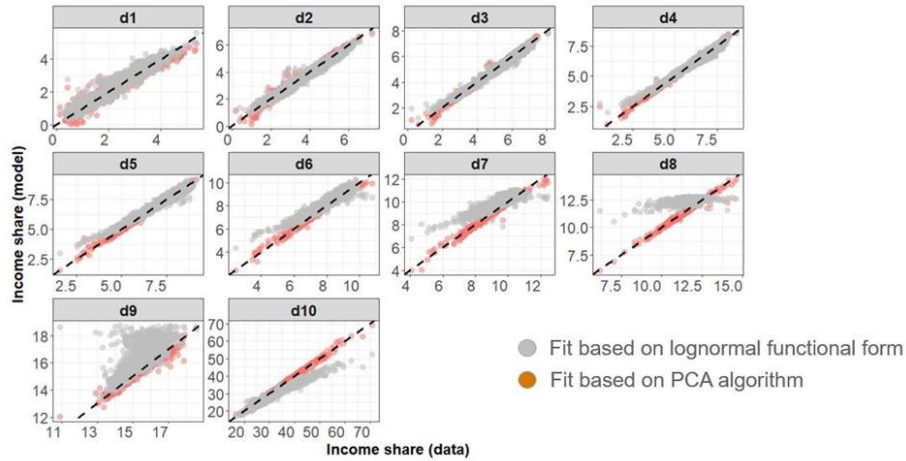
5 **2.4 Imputing net income deciles based on summary measures of the GINI coefficient.**

6 As observed in Table 1, the majority of observations in our dataset are those from the imputation
7 from the GINI coefficient. In this section we will explain this imputation approach, why a new
8 imputation approach was necessary and why this approach is an improvement upon existing
9 methods.

10 For many countries, years, no data are available for the income or consumption deciles based on
11 household survey data. However, World Development Indicators (WDI) dataset (Reid, 2012) do
12 provide aggregate measures of the income distribution such as the GINI coefficient for some
13 country-year observations¹. Many studies have utilized the GINI coefficient in combination with
14 different functional forms to estimate the underlying income distribution (Shorrocks & Wan,
15 2008; Soergel et al., 2021). Most prominent amongst these methods is the usage of the lognormal
16 functional form along with the GINI coefficient to derive the underlying distribution.

17 However, methods such as the lognormal functional form have documented limitations. For
18 example, the observations are known to deviate from the lognormal in the tails of the
19 distribution (Badel et al., 2020; Chotikapanich, 2008). Moreover, the lognormal functional form
20 is assumed for every country for every year. Recently, a non-parametric approach was developed
21 which uses the GINI coefficient in combination with a two-component model based on a
22 principal components analysis (PCA) to produce a more accurate estimate of income deciles
23 (Narayan et al., 2023). This method addresses some of the limitations of the lognormal
24 functional form. The performance of the non-parametric PCA based approach compared to the
25 lognormal functional form is described in more detail in ~~SI 2 Figure 4~~ Figure 5 below. We found
26 that the PCA based approach improves the fit across several deciles compared to the lognormal
27 functional form. The paper by Narayan et al. (2023) contains a more extensive discussion on the
28 model fit and comparisons of fit across countries-, years and individual deciles.

¹ The WDI dataset has observations of the GINI coefficient from various research studies. However, the underlying income concept of the GINI coefficient is not always available.



1

2 *Figure 5: Comparison of fit of lognormal functional form (grey dots) with PCA based fit (orange dots) with data for each decile*
 3 *(facet). Lines represent 1 to 1 fit between x and y axis. Income shares are expressed as a percent of total income.*

4 For country-years where we could not find data on net income or consumption, we used this
 5 PCA based approach along with observed values of the GINI coefficient from the World
 6 Development Indicators Database (Reid, 2012) to impute the underlying net income distribution.

7 The PCA based approach can be described as follows.

8 The income deciles are calculated as

$$9 \quad D_{r,t} = a_{r,t}PC1 + b_{r,t}PC2 \quad (2)$$

10 Where,

11 D is a 10-dimensional vector of income shares for all population deciles in region r at time t.

12 PC1 and PC2 are the two principal components, also vectors of length 10 (Values of PC1, PC2
 13 are provided in **SI 2 Figure 2, SI 2 Table 3**)

14 a and b are coefficients of the two principal components specific to each region and time

15 The coefficient a is derived from the GINI coefficient using a regression equation estimated on
 16 **1659** observations of national net income distribution

$$17 \quad a_{r,t} = -11.4815 + 29.71708 * GINI_{r,t} \quad (3)$$

1 And the coefficient b is estimated using lagged values of the Palma Ratio ($d10/(d1+d2+d3+d4)$)
2 and income share in the ninth decile and the current period labor share of GDP

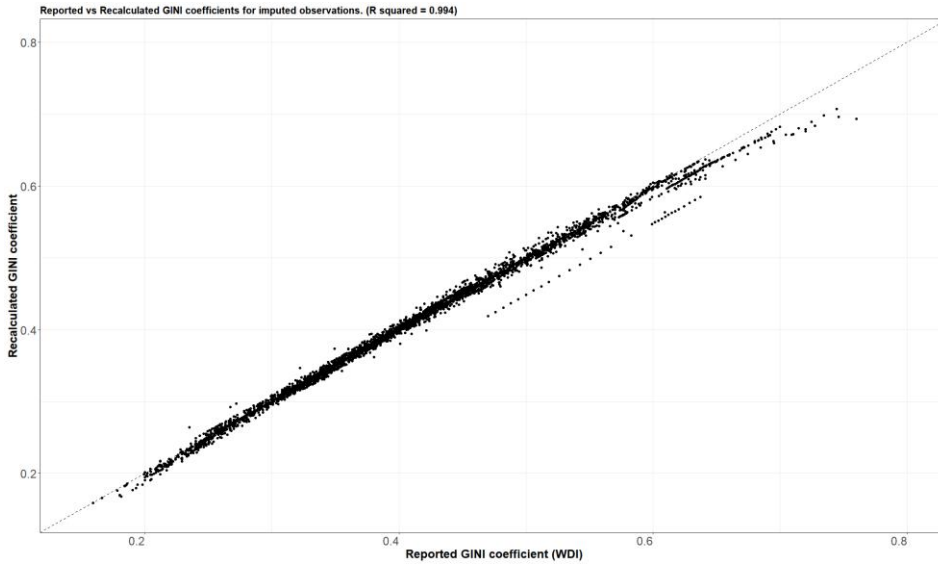
$$3 \quad b_{r,t} = -17.18222 + (1.07957 * LabShareGDP_{r,t}) + (113.10810 * Ninth Decile_{t-1}) \\ 4 \quad \quad \quad + (-0.36392 * PalmaRatio_{r,t-1}) \quad (4)$$

5
6 Using this approach, we were able to fill in values for various country-years. The observations in
7 our dataset are now summarized in [Table 1 above](#).

8 [As mentioned and discussed above](#),

9 Since we used a summary measure (GINI coefficient) to derive the underlying distribution, we
10 also validated our imputation approach by recalculating the GINI coefficient from the imputed
11 distribution and comparing it with the reported GINI coefficient (Figure 6). We observe that our
12 re-calculated values largely have a one-to-one correlation with the input GINI values suggesting
13 that the imputation did not introduce many errors (overall R squared value of the comparison is
14 0.99). However, the relationship does start to weaken for countries with very high GINI
15 coefficients such as South Africa where the recalculated GINI coefficient is different from the
16 observed GINI coefficient by as much as 0.07 points. This is a result of the parameters of the
17 PCA algorithm which do not reproduce well values for outlier countries with extreme GINI
18 coefficients. We also observe that the re-calculated GINI coefficients for some countries are
19 different in different years. For example, in Malawi, there are large year to year jumps in the
20 reported GINI coefficients from year to year (**SI 2 Figure 3**).

1



2

3 *Figure 6: Comparison of the reported GINI coefficients from the WDI (x axis) with the recalculated GINI coefficients from the*
 4 *imputed distribution (y axis). Each dot is a country-year observation. The dashed line represents a one-to-one relationship.*

5 We also evaluated temporal trends in the complete dataset which now include values from direct
 6 observations and also imputed values. The top two panels in Figure 5 below shows trends in the
 7 income shares for the 10th decile for India and China across time from all data sources.

8 This approach helps us generate better coverage in our dataset and the PCA model provides a
 9 statistically valid method to generate the data from GINI coefficients. This approach does have
 10 some limitations, however. The GINI coefficients from the WDI can represent multiple income
 11 concepts. For example, in the US, the GINI for the US from the World Development Indicators
 12 database is based on gross income and the income distribution based on surveys (From LIS) is
 13 for net income, i.e., it includes adjustments for direct taxation². As a good first step in this this
 14 direction would be in addressing this, we used data from the “All the GINIs” dataset which
 15 clearly specifies the income concept of the derived GINI coefficient (G. Ferreira et al., 2015;
 16 Smeeding & Latner, 2015), to identify the income concepts of the GINIs used for interpolation.
 17 Based on that, we identified that roughly 4200 observations of the GINIs used for imputation are
 18 net income GINIs while the remaining are consumption/expenditure GINIs or Gross income
 19 GINIs (Table 45). Therefore, data points when derived from imputation of a
 20 consumption/expenditure/gross income GINI have been marked as such in our final dataset.
 21 Users can choose to use all data points together or filter data depending upon their needs.

² Note that the examination of the metadata for the LIS values for the US shows that the values are computed as
 the gross income distribution minus an imputed tax adjustment.

1
2
3

<u>Type of imputation</u>	<u>values</u>
<u>Imputed from Net income GINI</u>	<u>4201</u>
<u>Imputed from Expenditure and Consumption GINI</u>	<u>1303</u>
<u>Imputed from Gross income GINI</u>	<u>1333</u>
<u>Total</u>	<u>6837</u>

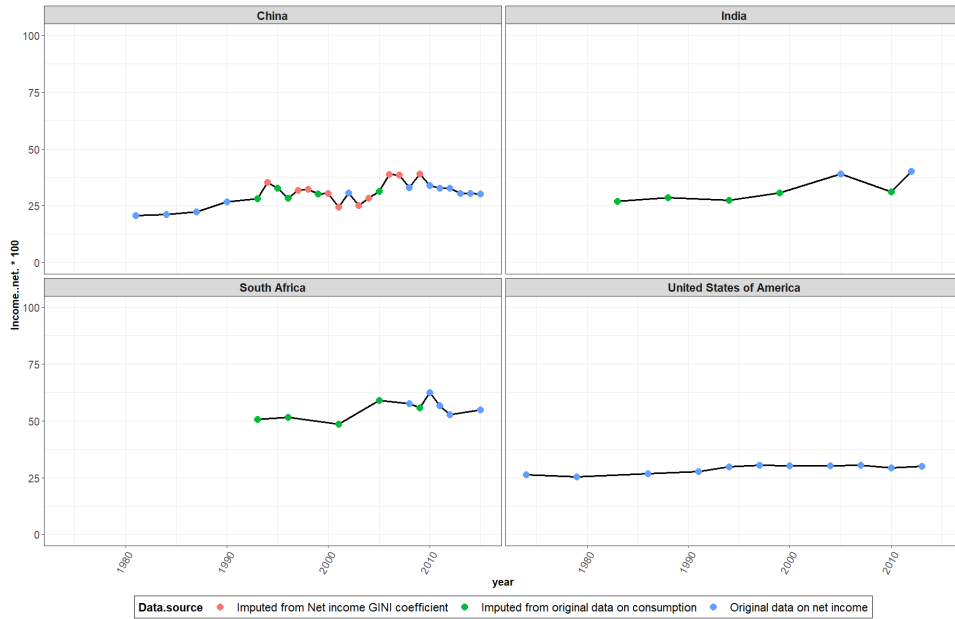
4 *Table 45: Description of source of GINI used for imputation.*

Formatted: Font: Italic

5 Given that the “All the GINIs” dataset still offers only a limited time series, ~~In the United States~~
6 ~~for example, we observe that the imputed income distribution values are consistently higher than~~
7 ~~observed values in all years (with income shares in upper deciles being approximately 5% higher~~
8 ~~when imputed compared to the actual data). This is likely because the GINI for the US from the~~
9 ~~World Development Indicators database is based on gross income and the income distribution~~
10 ~~based on surveys (From LIS) is for net income, i.e., it includes adjustments for direct taxation².~~
11 ~~[This still suggests a limitation in our imputation approach and one possible next step would be~~
12 ~~to only use net income GINIs for the imputation of the decile level income distribution. Figure 7~~
13 ~~below shows the full time series of our dataset based on different types of imputation performed.~~
14 ~~To implement this next step, we would require a dataset that clearly defines the income concept~~
15 ~~for the GINI coefficient provided. A good first step in this this direction would be to use data~~
16 ~~from the “All the GINIs” dataset which clearly specifies the income concept of the derived GINI~~
17 ~~coefficient (G. Ferreira et al., 2015; Smeeding & Latner, 2015).~~

²Note that the examination of the metadata for the LIS values for the US shows that the values are computed as the gross income distribution minus an imputed tax adjustment.

1



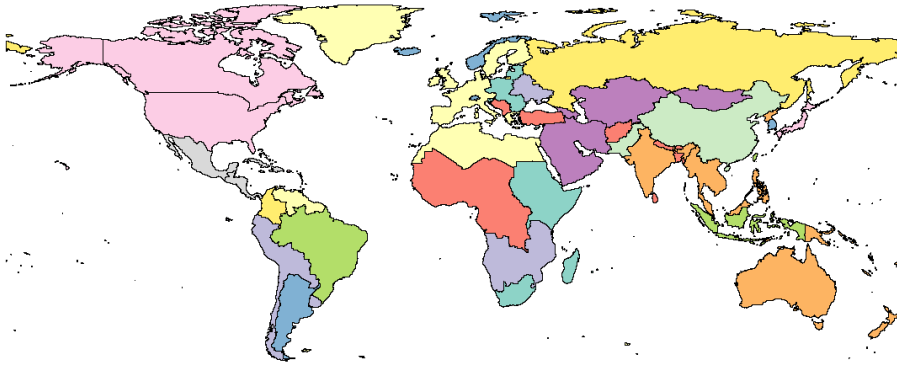
2

3 *Figure 7: Temporal trends in the 10th decile (presented as a % of total income) for the complete dataset. Colors represent*
 4 *different data sources.*

5

6 **3. Aggregating income distributions to the regional level**

7 The motivation for developing this country-level dataset was to initialize decile level income
 8 distribution values for the Global Change Analysis Model (GCAM). Models like GCAM operate
 9 on regional boundaries and therefore would require the income distributions to be aggregated to
 10 their respective regional boundary conditions. We aggregated the income
 11 distributions from the country level to 32 geographic regions represented by GCAM. The 32
 12 regions are shown as a map in Figure 8.



1

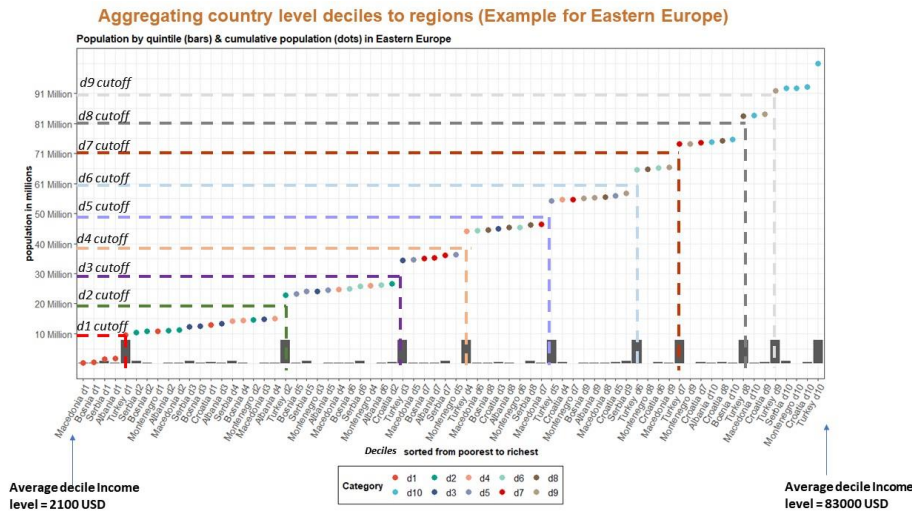
2 *Figure 8: Map of the 32 GCAM regions. These 32 GCAM regions are based on 229 country boundaries.*

3 Aggregating income distributions to the regional (where a region is made up of multiple
 4 countries) level is not straightforward because countries within regions differ in population size,
 5 average income level, and level of inequality in the income distribution. For example, an
 6 individual who belongs to the 10th decile in Romania would not necessarily be counted amongst
 7 the 10th decile of Europe as a whole, given the difference in the overall income level of Romania
 8 relative to higher income level of other European countries such as Germany and France.
 9 Similarly, even countries with similar average income levels may differ significantly in how
 10 income is distributed across deciles.

11 The aggregation of the country level income distributions to the regional income distributions
 12 involved the following steps:

- 13 1. First, we sorted all country net-income deciles in the region by the average decile income
 14 level, from lowest to highest income (The net income distribution shares are applied to
 15 this GDP per capita, measured in at PPP (2011 USD) to arrive at the income level). We
 16 use GDP per capita here, since that variable is the income proxy in GCAM.
- 17 2. Next, we calculated the cumulative population for each of these country income groups.
 18 The cumulative population over all country income groups matches the regional total
 19 population.
- 20 3. We then calculated cumulative population cutoffs that would create regional population
 21 deciles by dividing the regional population by 10.
- 22 4. Based on these cutoffs, we calculated the regional decile shares of income by assuming a
 23 uniform distribution of income within each country-decile. Thus, wherever a country
 24 decile spanned a regional cutoff, its income was split between regional deciles in
 25 proportion to the country population falling in each regional decile.

1 Figure 9 below illustrates our aggregation approach for GCAM region 14, Europe Non-EU,
 2 which is made up of Albania, Bosnia, Croatia, Macedonia, Montenegro, Serbia and Turkey. The
 3 figure demonstrates that a given regional decile can contain a mix of deciles at the country level.
 4 For example, the regional d2 consists of d3 and d4 values of some low-income countries such as
 5 Serbia and Albania. The regional d10 contains both the d9 and d10 values from Turkey.

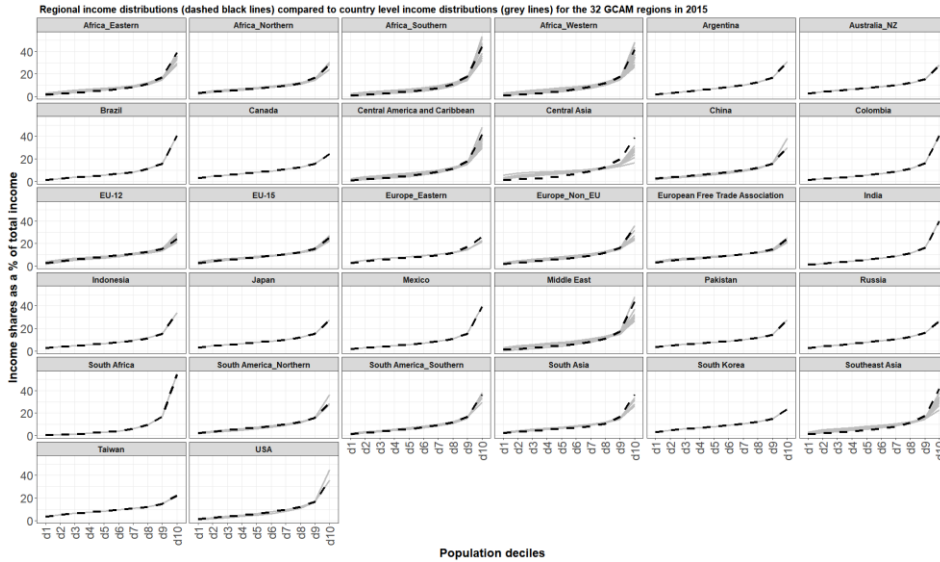


6
 7 Figure 9: Explanation of our aggregation approach. On the x axis all deciles within the region are sorted from low income to high
 8 income. Bars track the population. The dots show the cumulative population compared to the decile level income. Dashed lines
 9 show the new regional cutoffs for the deciles.

10 We also compared the aggregated income distribution to the country level income distributions
 11 for 2015 (Figure 10). We find that the aggregated income distributions are mostly driven by
 12 trends in the income distribution of the most populous countries in the region, as expected. In the
 13 example above, the income distribution for GCAM region 14 (Europe Non-EU) is largely driven
 14 by the income distribution of Turkey, which is the most populous, and most unequal, country in
 15 that region (e.g., Turkey represents approximately 75% of the regional population in 2015).
 16 There are certain cases where the regional distribution is significantly different than the country-
 17 level distributions. In Central Asia for example, the regional income distribution is much more
 18 unequal (regional GINI is 0.53) compared to the country level GINIs (Highest GINI is 0.39).
 19 This is because there is considerable variation in the income levels across countries. The
 20 country-level average incomes range from USD 2011 in Tajikistan to USD 23485 in Uzbekistan.
 21 This further illustrates why a specific aggregation method was necessary to construct these
 22 regional income distributions (Simple aggregation methods would miss such intra-regional
 23 dynamics).

24

1



2

3 *Figure 10: Regional income distributions (Dashed black line) compared to the national income distributions (grey lines) in each of*
 4 *the 32 regions in 2015.*

5

6 **4. Quantifying coverage and assessing regional bias in the data**

7 As mentioned earlier, we intended to develop a dataset for net income distribution for the 229
 8 countries aggregated to 32 regions used in GCAM. As shown in Table 5, we were unable to
 9 find any data on net income or consumption for 39 of those 229 countries. Previous models that
 10 have been developed for projecting income distributions have been based largely on data for high
 11 income countries (Rao et al., 2019; Sauer et al., 2020).

12 In order to evaluate whether the lack of data for the 56 countries introduces a bias, we assessed
 13 the data coverage in terms of percent of global population (total population of 229 countries) and
 14 percent of global GDP (total GDP at MER for 229 countries) for our dataset. We found that our
 15 dataset covers 98% of the global population and 93% of the global GDP in any given year.

16 Similarly, we also compared the average population and GDP of countries with and without data
 17 for five years (Table 6) and found that the average population of countries with data in the last
 18 historical year, i.e., 2015, is significantly higher (19 times) than the average population of
 19 countries without data. Similarly, the average GDP of countries with data is roughly 4.5 times
 20 the average GDP of countries without data.

Average national population (in thousands)	Average national GDP at MER (Billion 2010 USD)
--	--

Year	Data available	Data not available	Data available	Data not available
2010	37988	2835	370	90
2011	38881	2777	385	90
2012	39351	2808	394	90
2013	39822	2838	404	91
2014	40066	2915	414	91
2015	40610	2063	423	93

1 *Table 56: Comparison of national average population and national average GDP (at MER) for*
2 *countries with and without data for five historical years.*

3 Since this data will be used to initialize income distributions in the GCAM model, we also
4 evaluated whether the data would introduce a bias for any GCAM region (e.g., is there no
5 coverage or poor data coverage for any given GCAM region).

6 To evaluate this, we divided the countries in our dataset into the 32 geographical regions
7 modelled by GCAM. We then assessed the data coverage in terms of a percent of population (SI
8 3 Table 4) and GDP (SI 3 Table 5) for each of these regions. While these regions are specific to a
9 particular model, they also well represent heterogeneity across countries in terms of regional
10 economic and demographic conditions.

11 An example of a result of this assessment is that in the region of Africa Eastern we found data
12 that covers 64% of the region's population in 2010 and 40% of the region's GDP for the same
13 year. We performed this assessment for 5 years from 2010 to 2015. The purpose of this
14 assessment is to verify whether we have some coverage of data for all regions of the world
15 within those 5 years which would increase our confidence that our models are not biased towards
16 high income countries. The lowest coverage in our dataset is found for the Middle East region
17 where our data covers roughly 60% of the region's population and 40% of the region's GDP.

18 **6.5 Discussion**

19 In this paper we present a new consistent dataset on the net income distribution across 190
20 countries from 1958-2015. This dataset is also available for 32 aggregated regions ~~and the world~~
21 ~~as a whole~~. To our knowledge there is no other dataset that presents consistent data at multiple
22 geographical scales that has been documented in a peer-reviewed article. This complete and
23 harmonized dataset may be useful for efforts related to modelling of the net income distribution.

24 The aggregation method presented in this paper (section 43) takes into account both within-
25 country and across-country inequality when aggregating income distributions to regional
26 ~~boundaries of the world~~. This is important to regions where there is significant diversity in the
27 income distribution across countries such as Central Asia, where the aggregated income
28 distribution is significantly more unequal than any of the member countries (Figure 108).

29 There are a number of areas of improvement that we have noted that can be explored as next
30 steps or in future updates to this dataset. First, we have used a simple linear regression approach
31 when converting the consumption distributions to net income distribution. This can be improved
32 upon if more data becomes available related to the savings rate across countries or if the income

Formatted: Font: (Default) Times New Roman, 12 pt, Bold

Formatted: Numbered + Level: 1 + Numbering Style: 1, 2, 3, ... + Start at: 1 + Alignment: Left + Aligned at: 0.25" + Indent at: 0.5"

1 within countries can be broken down into the various incomes and expenditures similar to a
2 Computable General Equilibrium (CGE) framework.

3 Similarly, while our imputation approach greatly increased spatio-temporal coverage in our
4 dataset

5 , we noticed

6 that the GINI values from the WDI can represent multiple income concepts.

7 In the future, these gross income or consumption GINIs should also be converted to net income

8 In the future, these gross income or consumption GINIs should also be converted

9 to net income GINIs before the imputation. This would require more detailed data on the input

10 GINI coefficients. One possible next step would be to construct a method for such a

11 conversion using GINI values from datasets such as the “All the GINIs”

12 dataset which tracks the type of the GINI coefficient (G. Ferreira et al., 2015; Smeeding &

13 Latner, 2015). Another option would be to explicitly generate a tax adjustment to convert gross

14 income values to net income.

15 We further found that the PCA based imputation approach generates some error when imputing
16 the income distributions of highly unequal regions such as South Africa. As more data on income
17 distributions becomes available, the PCA algorithm can be re-parameterized to newer data.

18 When this happens, the imputation should be re-performed.

19 Finally, the data generation described above is documented as an open-source workflow of a
20 software package called *pridr* which can be used to generate and re-aggregate these data. The
21 software package is available on GitHub and the dataset itself is available as a version-
22 controlled release on Zenodo (See data availability statement below).

23

24 **7.6.Data availability**

25 The main dataset is available here on Zenodo- <https://zenodo.org/record/7093997> (Narayan et al.
26 2022) There are 2 main datasets available –

27 1. 32 region income deciles from 1958 to 2015

28 3.2.ISO level income distributions from 1958-2015

29

30 **Competing interests**

31 The authors declare that none of the authors have any competing interests.

32 **Acknowledgements**

33 This research was supported by the U.S. Department of Energy, Office of Science, as part of
34 research in Multi Sector Dynamics, Earth and Environmental System Modeling Program. The
35 Pacific Northwest National Laboratory is operated for DOE by Battelle Memorial Institute under
36 contract DE-AC05-76RL01830

37 **References**

- 1 Babones, S. J., & Alvarez-Rivadulla, M. J. (2007). Standardized income inequality data for use in cross-
2 national research. *Sociological Inquiry*, 77(1), 3-22.
- 3 Badel, A., Huggett, M., & Luo, W. (2020). Taxing top earners: a human capital perspective. *The Economic*
4 *Journal*, 130(629), 1200-1225.
- 5 Bank, W. (2015). PovcalNet. In.
- 6 Calvin, K., Patel, P., Clarke, L., Asrar, G., Bond-Lamberty, B., Cui, R. Y., Di Vittorio, A., Dorheim, K.,
7 Edmonds, J., & Hartin, C. (2019). GCAM v5. 1: representing the linkages between energy, water,
8 land, climate, and economic systems. *Geoscientific Model Development*, 12(2), 677-698.
- 9 Chotikapanich, D. (2008). *Modeling income distributions and Lorenz curves* (Vol. 5). Springer Science &
10 Business Media.
- 11 Deaton, A., & Zaidi, S. (2002). *Guidelines for constructing consumption aggregates for welfare analysis*
12 (Vol. 135). World Bank Publications.
- 13 Frank, M. W. (2009). Inequality and growth in the United States: Evidence from a new state-level panel
14 of income inequality measures. *Economic Inquiry*, 47(1), 55-68.
- 15 Fujimori, S., Hasegawa, T., & Oshiro, K. (2020). An assessment of the potential of using carbon tax
16 revenue to tackle poverty. *Environmental Research Letters*, 15(11), 114063.
- 17 G. Ferreira, F. H., Lustig, N., & Teles, D. (2015). Appraising cross-national income inequality databases:
18 An introduction. *The Journal of Economic Inequality*, 13, 497-526.
- 19 Hallegatte, S., & Rozenberg, J. (2017). Climate change through a poverty lens. *Nature Climate Change*,
20 7(4), 250-256. <https://doi.org/10.1038/nclimate3253>
- 21 Hughes, B. B. (2019). *International futures: Building and using global models*. Academic Press.
- 22 Jafino, B. A., Walsh, B., Rozenberg, J., & Hallegatte, S. (2020). Revised estimates of the impact of climate
23 change on extreme poverty by 2030.
- 24 [Lakner, Christoph, Branko Milanovic, and Branko Milanovic. "World Panel Income Distribution \(LM-
25 WIPD\)." Washington, DC: The World Bank \(2013\).](#)
- 26 Narayan, K. B., O'Neill, B. C., Waldhoff, S. T., & Tebaldi, C. (2023). Non-parametric projections of national
27 income distribution consistent with the Shared Socioeconomic Pathways. *Environmental*
28 *Research Letters*, 18(4), 044013.
- 29 Narayan, K. B., O'Neill, B. C., Waldhoff, S., and Tebaldi, C.: A consistent dataset for net income deciles
30 for 190 countries, aggregated to 32 geographical regions and the world from 1958-2015 (1.0.0),
31 <https://doi.org/10.5281/zenodo.7093997>, 2022.
- 32 Piketty, T., & Saez, E. (2003). Income inequality in the United States, 1913–1998. *The Quarterly journal*
33 *of economics*, 118(1), 1-41.
- 34 Rao, N. D., & Min, J. (2018). Less global inequality can improve climate outcomes. *Wiley Interdisciplinary*
35 *Reviews: Climate Change*, 9(2), e513.
- 36 Rao, N. D., Sauer, P., Gidden, M., & Riahi, K. (2019). Income inequality projections for the Shared
37 Socioeconomic Pathways (SSPs). *Futures*, 105, 27-39.
38 <https://doi.org/https://doi.org/10.1016/j.futures.2018.07.001>
- 39 Ravallion, M. (2015). The Luxembourg Income Study. *The Journal of Economic Inequality*, 13(4), 527-547.
40 <https://doi.org/10.1007/s10888-015-9298-y>
- 41 Reid, C. D. (2012). World development indicators 2011. *Reference Reviews*, 26(8), 26-27.
- 42 Sauer, P., Rao, N. D., & Pachauri, S. (2020). WIDER Working Paper 2020/65-Explaining income inequality
43 trends: an integrated approach.
- 44 Shorrocks, A., & Wan, G. (2008). *Ungrouping income distributions: Synthesising samples for inequality*
45 *and poverty analysis* (929230058X).
- 46 Smeeding, T., & Latner, J. P. (2015). PovcalNet, WDI and 'All the Ginis': a critical review. *The Journal of*
47 *Economic Inequality*, 13(4), 603-628.

- 1 Smeeding, T. M., & Grodner, A. (2000). Changing Income Inequality in OECD Countries: Updated Results
- 2 from the Luxembourg Income Study (LIS). In (pp. 205-224). Springer Berlin Heidelberg.
- 3 https://doi.org/10.1007/978-3-642-57232-6_10
- 4 Soergel, B., Kriegler, E., Bodirsky, B. L., Bauer, N., Leimbach, M., & Popp, A. (2021). Combining ambitious
- 5 climate policies with efforts to eradicate poverty. *Nature Communications*, 12(1).
- 6 <https://doi.org/10.1038/s41467-021-22315-9>
- 7 Van der Mensbrugge, D. (2015). Shared socio-economic pathways and global income distribution.
- 8