

## PATHWAY 2: DEVELOPING COMPOSITE INDICATORS

---

Christopher Gandrud

SG1022, City University London

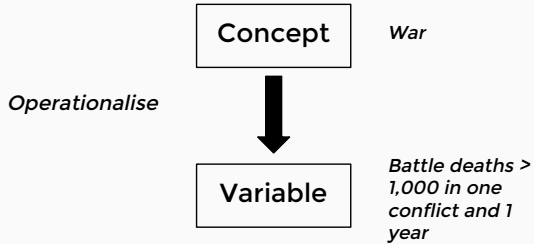
- What are composite indicators and who uses them?
- Making Composite Indicators
- Pros and cons

## WHAT ARE COMPOSITE INDICATORS

---

Last semester, we learned about concepts and variables.

- *Concept*: A phenomenon (e.g. poverty, democracy, human development, trust in institutions) we are interested in studying.
- *Variable*: Observable characteristic of a unit (e.g. person, city, country) that operationalises the concept.



Frequently in social science we are interested in complex concepts that cannot be operationalised with one variable.

Instead, they likely involve some combination of variables.

For example, what is democracy?

Conceptualisation: Dahl (1971) argued that democracy has two core attributes:

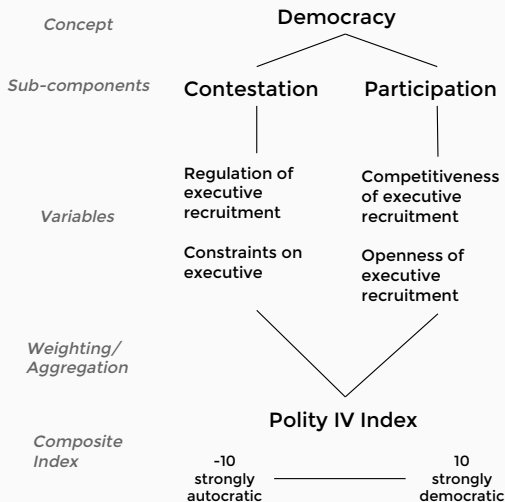
- *contestation*: competitive elections to choose leaders,
- *participation*: inclusive rules for and rates of participation.

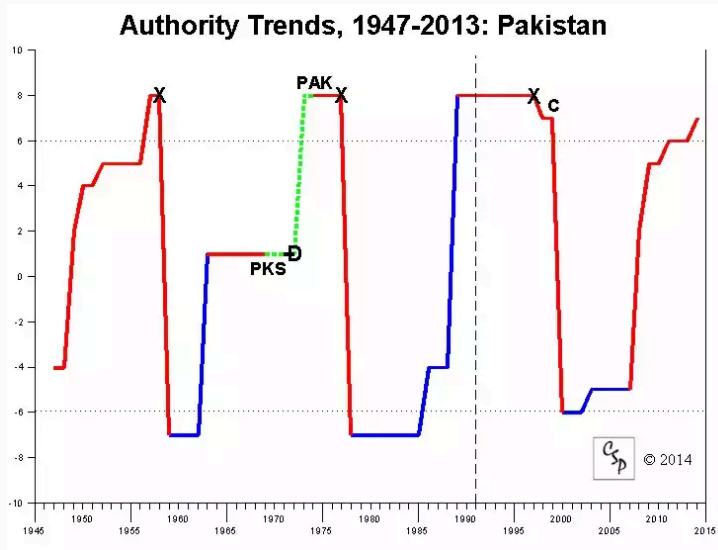
Note: not just elections/no elections or universal suffrage/limited suffrage. So ...



...we need a composite indicator to operationalise Dahl's democracy.

# CREATING THE POLITY IV DEMOCRACY INDEX



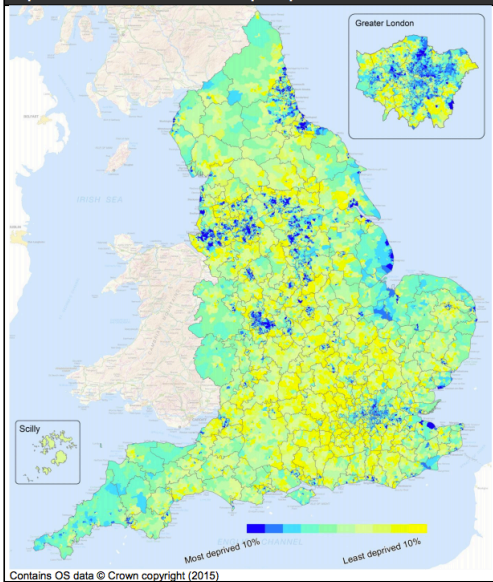


Why should we care about measuring concepts well?

If we don't have good measures of our concepts, we can't know how one thing effects another (relationships *between* concepts), how characteristics change over time, and how to improve the social world.

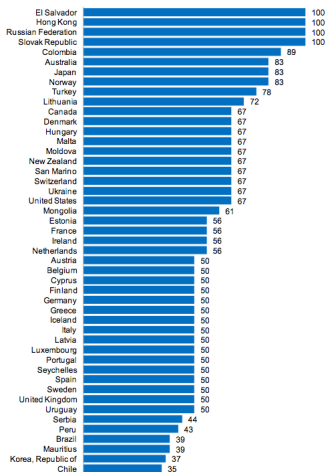
Composite indicators are a popular tool used by governments, international organisations, NGOs, and public policy advocacy groups to both understand the world and advocate change.

**Map 1: Distribution of the Index of Multiple Deprivation 2015**



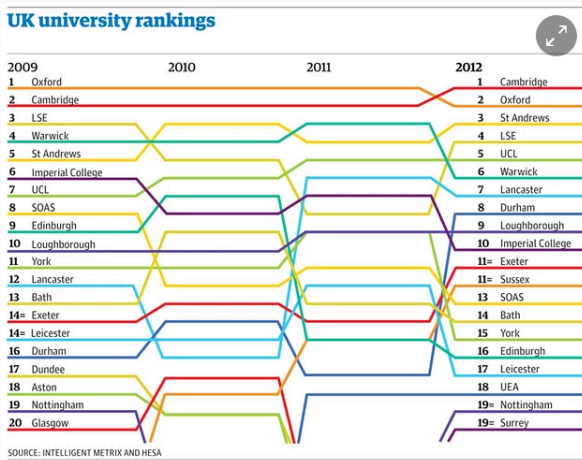
# IMF FISCAL TRANSPARENCY (GFS) INDEX

Figure 2. Average Comprehensiveness of GFS, 2011–13: Top Quartile





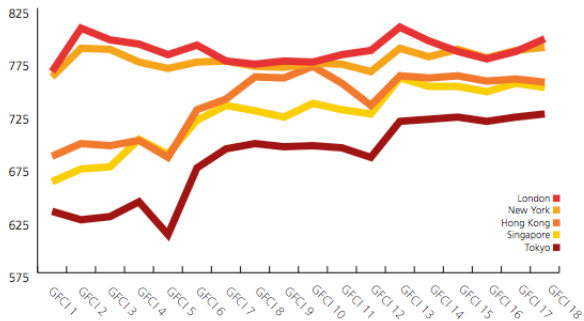
# UNIVERSITY RANKINGS



(<http://www.theguardian.com/news/datablog/2011/may/17/university-guide-2012-data-guardian>)

# GLOBAL FINANCIAL CENTRES INDEX (2015)

Chart 4 | Top Five Centres GFCI Ratings Over Time



## New York and London vie for crown of world's top financial centre

Financial centres face competition from Asia and increased regulation

Financial Times (1 Oct 2014)

## STEPS FOR CREATING COMPOSITE INDICATORS

---

1. Theoretical framework
2. Data selection
3. Address missing data
4. Multivariate analysis (note: we don't really much this in this course)
5. Normalisation
6. Weighting and aggregation
7. Validation

Modified from OECD (2008, 20-21)

All of these steps should be clearly documented in detail.

## THEORETICAL FRAMEWORK

---

“What is badly defined is likely to be badly measured.” (OECD 2008, 22)



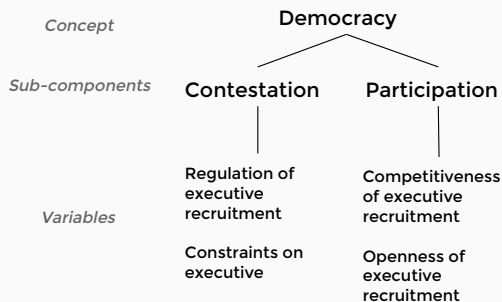
Concept definitions should:

- be clear and parsimonious,
- refer to a theoretical framework,
- select a relevant unit of analysis,
- link various sub-components and the underlying indicator.

### Sub-components:

- should be *statistically independent* of each other (want as few components as possible, don't want to double up. Though can compensate in the weighting stage),
- existing *linkages* between them should be *described theoretically and empirically* as much as possible.

# CONCEPTUALISING DEMOCRACY



## DATA SELECTION

---

A composite indicator is only as good as its parts.

Data on the indicators' component should be selected based on at least their:

- relevance to the theoretical framework,
- analytical soundness,
- timeliness,
- accessibility.

Gathering, cleaning, and merging the component variables is not trivial.

Has important substantive implications (e.g. error can bias your indicators) and can be time consuming.

So, needs to be well documented and reproducible.

## MISSING DATA

---

To aggregate your components into a composite indicator you need a complete data set, with no missing values.

For example, imagine we are going to make a component indicator with variable\_1, variable\_2, and variable\_3. A complete data set would look like:

country	variable_1	variable_2	variable_3
Algeria	1	10034	0
Cambodia	0	30020	0
Zambia	0	50302	1



However, we often (especially when working with country-level data) have missing data on at least some of our variables.

For example:

country	variable_1	variable_2	variable_3
Algeria	NA	10034	0
Cambodia	0	30020	0
Zambia	0	NA	1

Note: NA in R means missing value.

You should always assess how much missing data you have and consider why the data is missing.

Data can be:

- *Missing at random*: why the variable value is missing has nothing to do with the variable–random.
- *Not missing at random*: the variable value is missing because of some reason related to the variable.
  - E.g.: Low income countries don't have the money to pay for enough staff to gather GDP data.

Not missing at random is complex to deal with and requires explicit modeling of the missing process.

We do not cover this modeling process in this course.

However, if you suspect your data is not missing at random, you should discuss this, including the reasons why you believe the data is not missing at random.

If you have missing data, to get complete cases you can:

- drop incomplete cases,
- single impute data (e.g. replacing the missing value with the variable's mean/median/mode),
- multiple impute data (not covered in this course).

## ANALYSIS OF THE SUB-COMPONENT STRUCTURE

---

“Analysing the underlying structure of the data is still an art.”

(OECD 2008, 25)

Understanding the underlying structure is important for informing weighting and aggregation decisions, whether or not to include a variable.

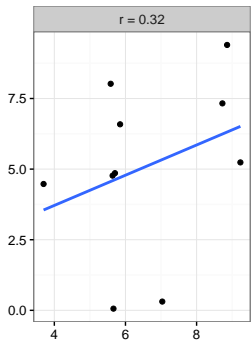
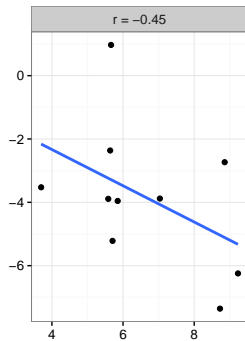
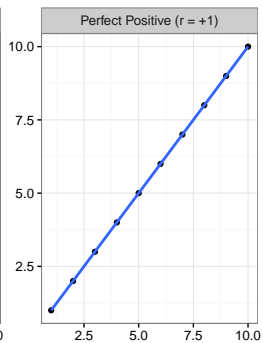
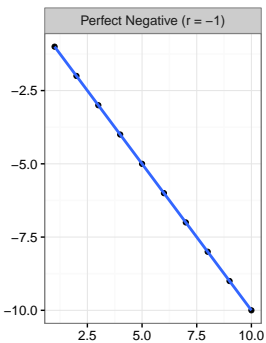
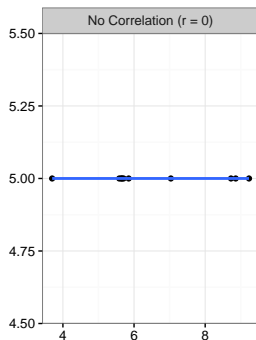
Note: important multivariate analysis tools–e.g. principal component analysis, cluster analysis–are beyond the scope of this class.

A simple way to examine the relationships between your sub-components is to create a correlation matrix.

Correlation: a mutual relationship between two variables.

Correlation coefficient: a number ranging from -1 to 1 calculated to represent the linear correlation between two variables. Usually denoted  $r$ .





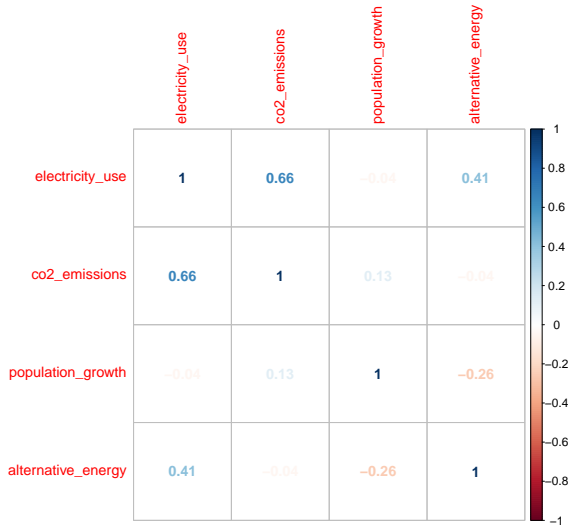
A simple way to examine the relationships between your sub-components is to create a correlation matrix.

Correlation: a mutual relationship between two variables.

Correlation coefficient: a number ranging from -1 to 1 calculated to represent the linear correlation between two variables. Usually denoted  $r$ .

Correlation matrix: A matrix showing correlations between many variables at once.

# CORRELATION MATRIX



The overall goal is to theoretically understand how the variables are related to each other theoretically and empirically.

Roughly balance:

- not wanting very highly correlated variables, as these indicate duplication
- and
- not wanting very lowly correlated variables, as this indicates they are related to different concepts.

Note: this is a very rough guide. More sophisticated tools await you in later courses.

## NORMALISE VARIABLES

---

Observable variables are often on different scales as they can have different measurement units.

For example, *life expectancy at birth* is in years ranging from 0 to > 100 and *GNI per capita* is in US dollars starting from > 300.

Obviously, adding these two variables together would weight GNI more than life expectancy.

Some notation:

$x_{u,t}$  observed variable value for unit  $u$  (e.g. country) at time  $t$  (e.g. year)

$l_{u,t}$  normalised value of the variable value for unit  $u$  at time  $t$

There are many approaches to normalising variables (see OECD 2008, 27-29).

Here are two for interval/ratio data ...



Min-Max rescales a variable to be between 0 and 1 based on the observed minimum and maximum values.

$$I_{u,t} = \frac{x_{u,t} - \min(X)}{\max(X) - \min(X)} \quad (1)$$

where  $X$  are all of the observed values of the variable. (For R students: it is a vector of the values.)

Note: extreme values can distort the measure and can widen the range of closely spaced values.

Z-Scores create a common scale with a mean of 0 and a standard deviation of 1.

$$l_{u,t} = \frac{x_{u,t} - \mu_X}{\sigma_X} \quad (2)$$

where  $\mu_X$  is the mean of the variable and  $\sigma_X$  is the standard deviation.

Note: extreme values can have a greater effect on the composite indicator.

All of the variables that you think have a “positive” influence on the underlying concept need to be positive.

All of the variables that you theoretically think have a “negative” influence on the underlying concept should be negative.

## NEED TO CONSIDER VARIABLE DIRECTION, EXAMPLE

For example, if we want to create a measure of environmental un-sustainability with the components: CO<sup>2</sup> Emissions, Population Growth, and Alternative Energy Use.

We would expect:

- More CO<sup>2</sup> Emissions → less sustainable.
- More Population Growth → less sustainable.
- More Alternative Energy Use → more sustainable

So, we need to reverse the scale of our Alternative Energy Use variable

To reverse the direction of a variable simply subtract each value by the maximum observed value, i.e.:

$$l_{u,t} = \max(\mathbf{X}) - x_{u,t} \quad (3)$$

## WEIGHTING AND AGGREGATION

---

Once you have your normalised variables ( $l_{c,t}$ ), then you need to consider how to aggregate them.

Specifically consider:

- **Weighting:** how important are each individual variables to the composite?

“No uniformly agreed methodology exists to weight individual indicators before aggregating them into a composite indicator.”

But generally, “what matters more ...weighs more.”

<https://composite-indicators.jrc.ec.europa.eu/?q=content/step-6-weighting>



If you simply add all of the variables together, you are implicitly assuming that they have an equal weight of 1. E.g.:

$$Cl_{u,t} = \sum I_{u,t} * 1 \quad (4)$$

Sometimes this makes sense: e.g. economic activity is often measured in the same currency, or if you have no prior knowledge about how important each component is.

There are many advanced techniques to determine variable weights (e.g. factor analysis, data envelopment analysis and unobserved components models).

In this course we are going to use “expert judgement” – you are the experts and will determine how much each component contributes to the composite indicator.

Imagine that you have four component variables  $X_1$ ,  $X_2$ ,  $X_3$  and  $X_4$ .

Based on your expert knowledge, you believe that  $X_1$  is much more important than the other three combined. So you give it a weight of 0.7 and the others 0.1:

$$Cl_{u,t} = 0.7x_{1u,t} + 0.1x_{2u,t} + 0.1x_{3u,t} + 0.1x_{4u,t} \quad (5)$$

Sometimes the concept we are measuring might be discrete. E.g. you are in a financial crisis or not in a financial crisis.

So, you might use your expert judgement to set a **threshold**, a point past which a unit goes from having the characteristic to not having the characteristics.

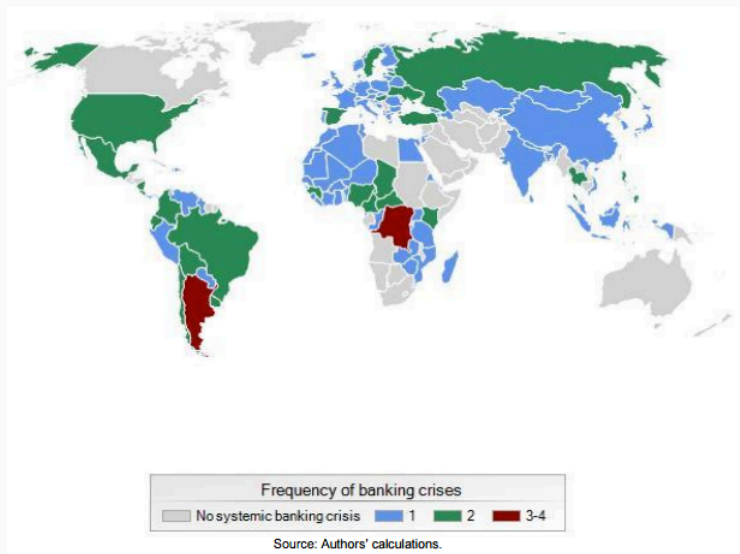
Laeven and Valencia (2013) determine a country has crossed a financial crisis threshold if:

- There is 'significant distress' in a country's financial system.

*and*

- At least three of six policy responses are used (e.g. bank holidays, bank nationalisations).

# LAEVEN & VALENCIA BANKING CRISES (1970-2011)



WHAT ARE WE MEASURING?

---

Once you have a composite indicator, your work is far from done.

You need to conduct numerous tests to determine if your indicator is a valid measure of the concept you are trying to measure.

We will discuss this more in Week 10.



There is likely **no perfect composite indicator** for any given concept.

**Transparency**—conceptualisation, data gathering/cleaning, construction, validation—is crucial for others to be able to understand and evaluate your indicators.

## PROS AND CONS

---

We need to remember that our indicators are **estimates** of what we want to measure.

We are **uncertain** about how well our indicators capture reality.

Uncertainty can be caused by at least:

- Error in our construct (i.e. including or omitting important variables)
- Measurement error in our raw data  $x_{c,t}$
- Error in our weighting/aggregation.

Composite indicators are popularly used to rank units (e.g. cities, countries).

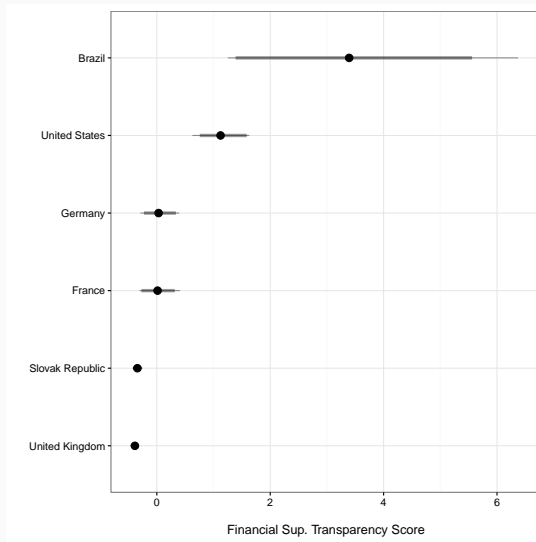
However, this can be highly misleading.

For example, Copelovitch, Gandrud, and Hallerberg (2015) created a measure of financial supervisory transparency.

**Table:** Financial Regulatory Transparency Index Ranking (2011)

Country	Rank
Brazil	1
United States	7
Germany	23
France	25
Slovak Republic	48
UK	61

# INDICATOR WITH ABSOLUTE SCORE & UNCERTAINTY



In this course, we won't be covering more advanced ways to quantify your uncertainty about your composite indicators.

However, you should be careful and honest when you compare units by ranking their composite scores.

Be honest about what you don't know.

# PROS AND CONS OF COMPOSITE INDICATORS

## Pros

Summarise complex multi-dimensional concepts so that we can understand relationships between them

Support policy-making

Facilitate communication with the public.

## Cons

May be misleading if poorly constructed & non-transparent

May lead to simplistic policy conclusions or missed to support particular predefined policy goal

May disguise problems in one dimension if construction is not transparent.

May exaggerate differences between units if uncertainty is not acknowledged.