

DEMO ARXIV TEMPLATE

A PREPRINT

H. Sherry Zhang 

Department of Econometrics and Business Statistics
Monash University
Melbourne, VIC
huize.zhang@monash.edu

Collaborators

Department of Econometrics and Business Statistics
Monash University
Melbourne, VIC

February 27, 2023

ABSTRACT

- indexes, useful, quantify severity, early monitoring,
- A huge number of indexes have been proposed by domain experts, however, a large majority of them are not being adopted, reused, and compared in research or in practice.
- One of the reasons for this is the plenty of indexes are quite complex and there is no obvious easy-to-use implementation to apply them to user's data.
- The paper describes a general pipeline framework to construct indexes from spatio-temporal data,
- This allows all the indexes to be constructed through a uniform data pipeline and different indexes to vary on the details of each step in the data pipeline and their orders.
- The pipeline proposed aim to smooth the workflow of index construction through breaking down the complicated steps proposed by various indexes into small building blocks shared by most of the indexes.
- The framework will be demonstrated with drought indexes as examples, but applicable in general to environmental indexes constructed from multivariate spatio-temporal data

Keywords indexes • data pipeline • software design

1 Introduction

Why index is useful, why people care about indexes

incorporate the following in why using index: multiple pieces of information (variables) that need to be taken into account

Many concepts relevant to decision making cannot be directly measured, however, they are crucial for resource allocation, early prevention, and other operational purpose. For example, fire authorities would be interested to quantify fire risk since bushfires can have a huge impact on monetary loss, health, and the local ecosystem. Climatologists would be interested in monitoring the change in global climate since variability in atmospheric and oceanic conditions has a direct impact on global weather and climate. Usually this concept of interest is associated with more than one variables and these variables need to be integrated to make decisions on the subject matter. A common approach to quantify concepts like these is to construct an index using these relevant variables. This allows researchers to compare the quantity of interest across entities (i.e. countries, regions) and also cross time.

Define what is an index, what is not

In this article, an index is defined as a tool to quantify a concept of interest that does not have a direct measure. The concept of interest doesn't have a direct measure can because it is impractical to measure at the population level. For example, it would be nearly impossible to include all the available stocks in the market to characterise stock market behavior, so indexes like Dow Jones Industrial Average, S&P 500, and Nasdaq Composite select a representative set of stocks to measure the overall market behavior. Also belonging to this category are the economic indexes like the Consumer Price Index, where price changes of a basket of items are weighted to measure inflation. The lack of direct measure could also because the concept itself is an unobservable human construction, rather than a physical quantity that can be measured. Many natural hazard and social concepts falls into this category. This includes drought indexes constructed from meteorological, agricultural, hydrological, and social-economic variables, e.g. Standardised Precipitation Index (SPI) (McKee et al. 1993) and Aggregated Drought Index (ADI) (Keyantash and Dracup 2004) among others. Social development indexes like Human Development Index (United Nations Development Programme 2022) and Global Liveability Index (Economist Intelligence Unit 2019) measure various aspects of the quality of human capital and urban life.

still need to tweak the tone a bit: "they are called index, they are not the index we will talk about"

Despite many quantity having the term *index* in their name, they cannot be technically classified as indexes according to the definition given above. The reason for these quantities to lose their index memberships is that they are variables can be accurately measured given the instrument precision. This includes quantities like precipitation of the driest month or percentage of days when maximum temperature is below 10th percentile. They are measures of precipitation and percentage of days under specific conditions (dries month, maximum temperature below 10th percentile). They are variables, or indicators, that can be used to construct indexes but are not indexes themselves. Similarly, a set of remote sensing indexes are not indexes, since they are measures of electromagnetic wave reflectance. This includes Normalized Difference Vegetation Index (NDVI) (Tucker 1979), derived from the ratio of difference over sum on two segments in the spectrum, also called band: near-infrared (NIR) and red. So are the "indexes" derived from NDVI, e.g. Vegetation Condition Index (Kogan 1995). Notice that this does not exclude all the construction derived from remotes sensor variables to be valid indexes. For example, Vegetation Drought Response Index (Brown et al. 2008) is a valid index since it integrates climate, satellite, and biophysical variables to quantify vegetation stress.

What is the challenges with current index construction

see if there is any paper describing this type of pains

useful to reference tidy data and tidy model that makes the workflow on modelling tidy somewhere in introduction

Currently, index construction lacks a standardised workflow. It is often up to researchers or research institutions to decide whether to provide open source code on the new indexes, what would be the best user interface for other researchers to use the new indexes, and how easily the new indexes can be compared with other existing indexes. This makes the computation lack transparency and indexes cumbersome to experiment with:

- Researchers who wish to validate the indexes calculated from large institutes need to reinvent the wheels themselves since the source code used for computing is often not available for public consumption;
- Open-source code provided by research groups has a narrow margin for exploring other options outside the provided;
- Similar steps used by different indexes are difficult to spot since the design of the user interface for indexes often includes all the steps under a single function call; and
- It is generally hard to inspect intermediate results during the index construction if users wish to check the output of a certain step.

what can be done if people adopt this pipeline/ why it is beneficial?

This paper proposes a data pipeline for index construction. By recognising the common steps shared by many indexes, we develop a pipeline that breaks down index construction into multiple modules and allow operations in various modules to be combined like building blocks to construct indexes. The pipeline approach is general while adaptable to most index construction. It allows indexes to be created, studied, and compared in a structured tidy form and enables statistical analysis of indexes to be performed easily: More specifically, it enables researchers to 1) validate the indexes calculated from external organisations, 2) unify various indexes under the same framework for computing, 3) swap or adjust individual steps in the index construction to study their contribution, 4) calculate uncertainty on indexes through bootstrap or others, 5) enhance existing indexes through comparing and studying their statistical properties, and finally, 6) propose new indexes from combining different steps in existing indexes.

who would benefit from this paper

This work is of interest to researchers actively developing new indexes since it encourages new indexes to be delivered in an easy-to-reproduce design. It would also provide analysts who wish to compute a range of indexes in their analysis a uniform interface to build relevant indexes from raw data. For statisticians and software developing engineers, this work frames the process of index construction in a more user-oriented workflow and could motivate similar research for other process in scientific computing.

The rest of the paper is structured as follows: Section 2 reviews the concept of data pipeline in R. The pipeline framework for index construction is presented in Section 3. Section 4 explains how to include a new building block in each pipeline module. Examples are given in Section 5 to demonstrate the index construction with the pipeline built.

2 Data pipeline

Think about if there is another word for data pipeline

Why you should care about pipeline

Data pipeline is not a new concept to computing. It refers to a set of data processing elements connected in series, where the output of one element is the input of the next one. Wickham et al. (2009) argues that whether made explicit or not, the pipeline has to be presented in every graphics program. The paper also argues that breaking down graphic rendering into steps is beneficial for understanding the implementation and comparing between different graphic systems. The discussion on pipeline construction is well documented in early interactive graphics software: Buja et al. (1988), Sutherland et al. (2000), and Xie, Hofmann, and Cheng (2014) and their pipeline steps include non-linear transformation, variable standardization, randomization and dimension reduction.

What is pipeline, its underlying software design philosophy, and how these are reflected in R

One of the most commonly known pipeline examples is perhaps the Unix pipeline where programs can be concatenated with | to flow the output from the last program into the next program, i.e.

```
command 1 | command 2 | command 3 | ...
```

To solve a complex problem, the Unix system builds simple programs that do one thing well and work well together. This design is also reflected in the tidyverse ecosystem in R. To solve a complicated data problem using tidyverse, analysts typically build the solution using a collection of tools from the tidyverse toolbox. The data object can flow smoothly from one command to the next, safeguarded by the tidy data format (Wickham 2014), which prescribes three rules on how to lay out tabular data. The tidyverse tools also embrace a strong human-centered design where function names are intuitive and easy to reference through autocomplete. With the tidyverse design principle in mind, the tidymodel suite enables analysts to build machine learning models through the data pipeline. It includes typical tasks required in machine learning like data resampling, feature engineering, model fitting, model tuning, and model evaluation. An advantage of tidymodel pipeline over separate software for individual models is that analysts no longer need to write model-specific syntax to work with each model, but pipeline-specific syntax that is applicable to all the models implemented in tidymodel. This allows users to easily experiment with a collection of machine learning models.

Constructing indexes would also benefit from pipeline and embracing the aforementioned design philosophy.

In index construction, data pipeline is often presented in a workflow diagram in the research paper to illustrate how the raw data is transformed into the final indexes. This agrees with Wickham's argument on the presence of the data pipeline, however, more often than not, the pipeline is not made explicit in the software. Often the time, all the steps are lumped into a single wrapper function, rather than being split into smaller, modulated functions. This increases the cost of maintaining and understanding the code base, gives analysts little freedom to customise the indexes for specific needs, and hinders reusing existing code for building new indexes. A pipeline approach unites a range of indexes under a single data pipeline and analysts can compose indexes from pipeline steps like building Legos from individual bricks. In this workflow, analysts are not limited by indexes that have been already proposed and can easily combine pipeline steps to compose novel indexes. Analysis of the indexes (i.e. calculation of uncertainty) is also feasible by adding external code into the pipeline.

3 A pipeline for building statistical indexes

3.1 How does the pipeline constructin of an index look like?

Consider a commonly used drought index: Standardized Precipitation-Evapotranspiration Index (SPEI) (Vicente-Serrano, Beguería, and López-Moreno 2010). Its construction involves:

- 1) transform the average temperature (TMED) into potential evapotranspiration (pet)

- 2) combine precipitation (prcp) and potential evapotranspiration (pet) into a single variable diff
- 3) aggregate the difference series with a sliding window (.scale)
- 4) fit a distribution to the aggregated series, and
- 5) derive the index value from the normal density values.

Conventionally approach may combine all these steps into a single function, with some level of modularity. However, these modules may only work for the selected index offered by the package.

Under the pipeline approach, analysts first need to identify which module each step belongs to.

Below shows the pseudocode for constructing SPEI with the pipeline:

```
DATA %>%
  var_trans(.method = thornthwaite, Tave = TMED, ..., .new_name = "pet") %>%
  dim_red(diff = prcp - pet) %>%
  aggregate(.var = diff, .scale = 12, .new_name = "agg") %>%
  dist_fit(.method = "lmoms", .var = agg, .dist = DIST) %>%
  augment(.var = agg)
```

The pipeline construct allows for multiple .scales and .dist to be evaluated in aggregate() and dist_fit() to compare index under different parameterisations. The result can then be passed into the ggplot2 to crease visualisation. Figure 1 compares the SPEI calculated with two distributions (log-logistic and pearson III).

Apart from evaluating multiples parameters, the pipeline approach allows the the steps written for one index can be directly extrapolate to another index building within the pipeline. The flexibility of the pipeline also integrate well with other existing packages, for examples, fitting distributions using L-moment is commonly used when constructing drought indexes. The package lmomco provides general L-moment fits to a wide range of distributions and users can easily access to all the distributions within the pipeline.

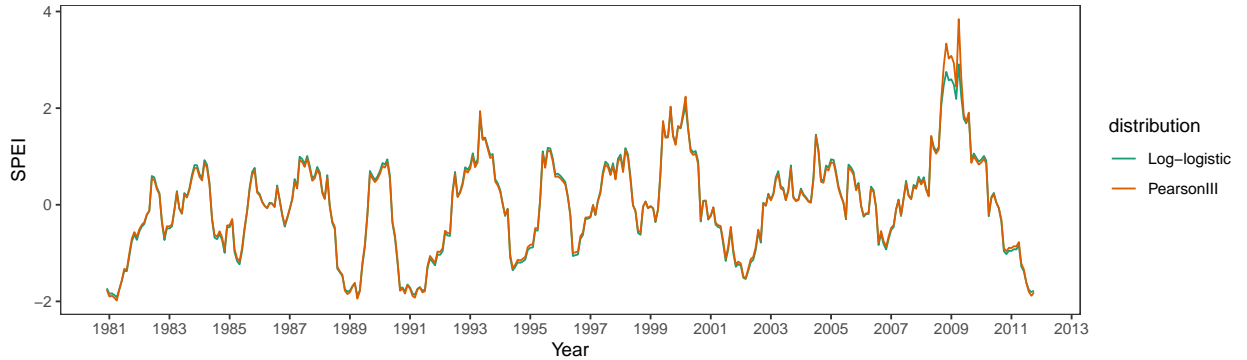


Figure 1: Standardised Precipitation Evapotranspiration Index (SPEI) calculated with two distribution fits in Step 3 described above: generalised logistic and pearson III distribution, using the wichita data from the package SPEI.

3.2 Pipeline steps for constructing indices

any index can be broken down into multiple steps and then we can do things with it: swap, change parameter, etc
variables == indicators

An overview of the pipeline is given in Figure 2 to illustrate the construction from raw data to the final indexes. The pipeline includes eight modules for operations in the spatial, temporal, and multivariate aspects of the data as well as modules for comparing and communicating indexes. Analysts are free to select the modules they need and arrange them in the order they see fit to construct indexes. While the starting point of the pipeline is raw data, there are steps prior to this that are crucial to the success of an index. For example, the defined index needs to be useful for measuring the concept of interest and variables need to be collected from reliable sources with proper quality control.

Before elaborating each of the eight pipeline modules as subsections, the data notation will be first introduced. Let $\mathbf{x}(\mathbf{s}; \mathbf{t})$ denote the raw data with spatial, temporal, and multivariate aspects: the spatial dimension $\mathbf{s} = (s_1, s_2, \dots, s_n)'$ is defined in the 2D space: $\mathbf{s} \in \mathcal{D}_s \subseteq \mathbb{R}^2$, the temporal dimension $\mathbf{t} = (t_1, t_2, \dots, t_J)'$ is defined in the 1D space: $\mathbf{t} \in \mathcal{D}_t \subseteq \mathbb{R}$. When more than one variable is involved, the multivariate data can also be written as: $\mathbf{x}(\mathbf{s}; \mathbf{t}) = (x_1(\mathbf{s}; \mathbf{t}), x_2(\mathbf{s}; \mathbf{t}), \dots, x_P(\mathbf{s}; \mathbf{t}))'$.

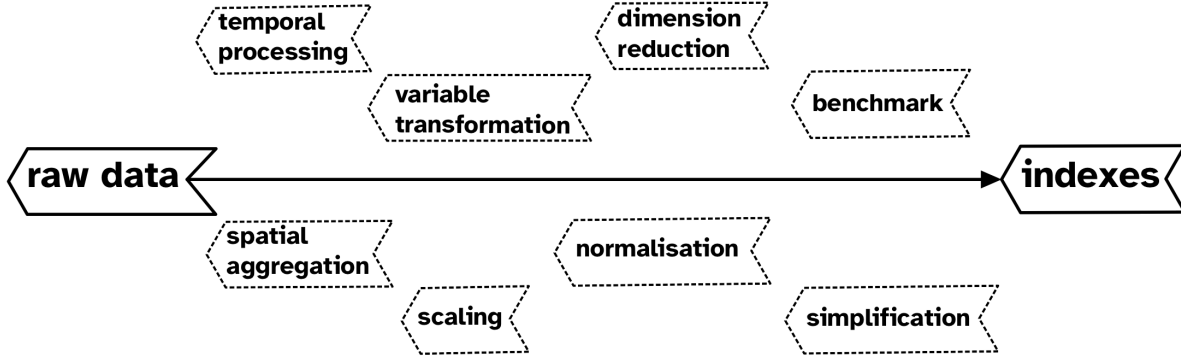


Figure 2: Diagram of pipeline steps for index construction. will need to be updated with better design and the distribution fitting step.

3.2.1 Temporal processing

The construction of an index sometimes needs to consider information from neighbouring time periods. The temporal processing is a general operator on the time dimension of the data in the form of

$$f_{\psi}(x(\mathbf{s}; \mathbf{t})), \quad (1)$$

where $\psi \in \Psi \subseteq \mathbb{R}^{d_{\psi}}$ is the parameters associated with the temporal operation and d_{ψ} is the number of parameter of ψ . A typical example of temporal processing is aggregation, which is used in the drought index SPI to measure the lack of precipitation for meteorological drought. In SPI, monthly precipitation is aggregated by a time scale parameter k : $x(s_i; t_{j'}) = \sum_{j=j'-k+1}^{j'} x(s_i; t_j)$, where j' is the new time index after the aggregation. In this notation, each spatial location is separately aggregated and precipitation is summed from k month back, $j' - k + 1$, to the current period, j' , to create the aggregated series, indexed by j' .

more explicit on k will influence 1) long term vs. short term, 2) uncertainty

The choice of time scales parameter k can result in variation in the calculated index values: a small k of 3 or 6 months produces the index more sensitive to individual months, while a large k of 24 or 36, an equivalent to a 2- or 3-year aggregation, gives dryness information relative to the long term condition. As will be shown in section [SECTION EXAMPLE], this variation may even lead to conflicting conclusions on the dry/wet condition of the area, highlighting the importance to account for index uncertainty when interpreting index values for decision-making.

Effective drought index

3.2.2 Spatial processing

Spatial processing may be needed when indexes are not calculated independently on each collected location or when variables collected from multiple sources need to be fused before further processing. The process can be written as a general operation in the form of

$$x(\mathbf{s}'; \mathbf{t}) = g_\theta(x(\mathbf{s}; \mathbf{t})), \quad (2)$$

where $\theta \in \Theta \subseteq \mathbb{R}^{d_\theta}$ is the associated parameters in the process and d_θ is the number of parameter of θ . An example of spatial processing is to align variables collected in different resolutions. When variables are collected at different resolutions, analysts may choose to down-sample those in a finer resolution, i , to match those in a coarser resolution, i' . This is a spatial aggregation and if aggregate using the mean, it can be written as

$$g(x) = \frac{\sum_{i \in i'} x}{n_{i'}}, \quad (3)$$

where $i \in i'$ includes all the cells from the finer resolution in the coarser grid and $n_{i'}$ is the number of observations falls into the coarser grid. Other examples of spatial processing include 1) borrowing information from neighbouring spatial locations to interpolate unobserved locations and 2) fusing variables from ground measures with satellite imageries.

3.2.3 Variable transformation

The purpose of variable transformation is to create variables that fits assumptions for further computing. These assumptions include a stable variance, normal distribution, or a certain scale required by some algorithms down the pipeline. Variable transformation is a general notion of a functional transformation on the variable:

$$h_\tau(x(\mathbf{s}; \mathbf{t})), \quad (4)$$

where $\tau \in T \subseteq \mathbb{R}^{d_\tau}$ is the parameter in the transformation if any, and d_τ is the number of parameter of τ . Transformation is needed for data that are highly skewed and some common transformations include log, quadratic, and square root transformation.

3.2.4 Scaling

While scaling can be seen as a specific type of variable transformation, it is separated into its own step to make the step explicit in the pipeline. The key difference between the two steps is that variable transformation typically changes the shape of the data while scaling only changes the data scale and can usually be written in the form of

$$[x(s_i; t_j) - \alpha] / \gamma. \quad (5)$$

For example, a z-score standardisation can be written in the above form with $\alpha = \bar{x}(s; t)$ and $\gamma = \sigma(s; t)$, a min-max standardisation uses $\alpha = \min[x(s_i, t_j)]$ and $\gamma = \max[x(s_i, t_j)] - \min[x(s_i, t_j)]$. Figure 3 shows a collection of variable pre-processing operations and uses color to differentiate whether the operation is a variable transformation or a scaling step. While both variable transformation and scaling are pre-processing steps, the scaling operations in green show the same distribution as the original data.

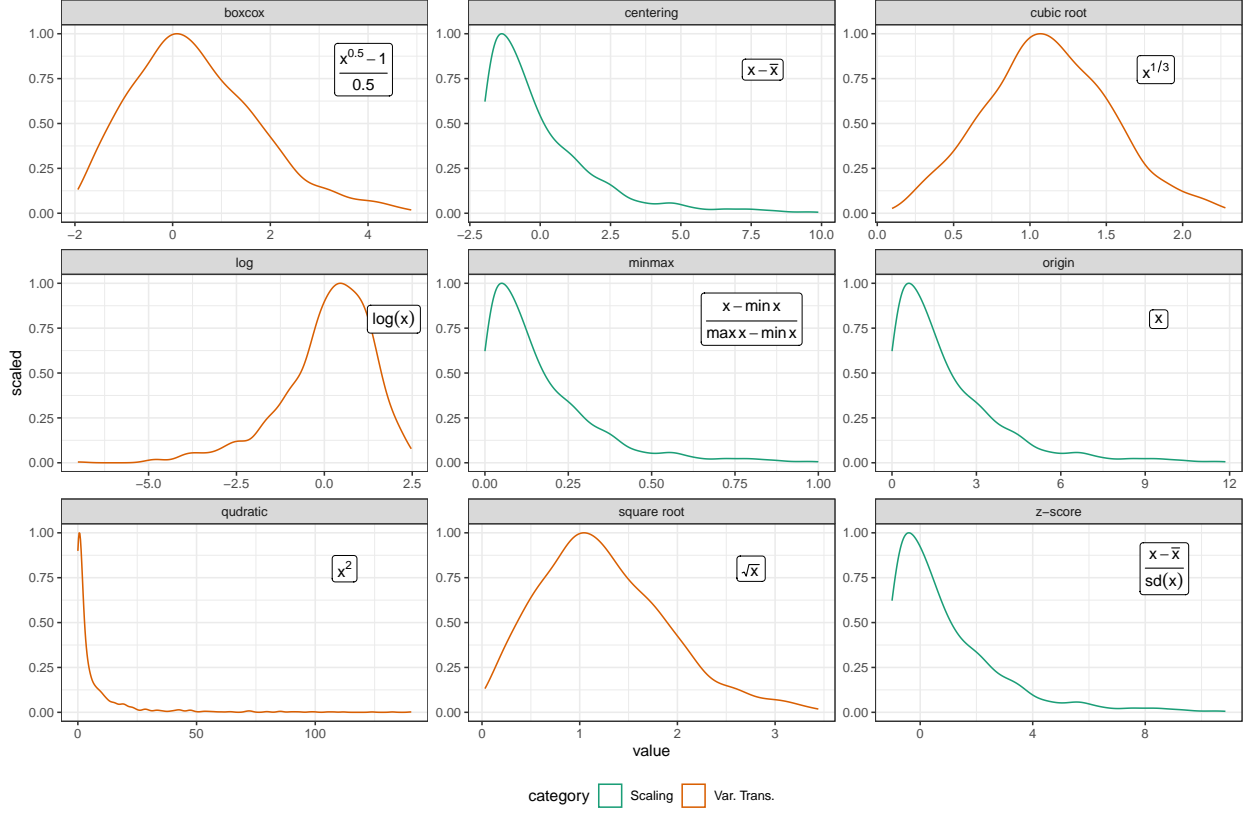


Figure 3: Comparison of operations in scaling (green) and variable transformation (orange) steps in free scale. Variables after the scaling operations have the same distribution as the origin, while the distribution changes after variable transformation.

3.2.5 Dimension reduction

When indexes are constructed from multivariate information, dimension reduction methods combine that information into a univariate series. In the pipeline, dimension reduction includes methods that take multivariate inputs and output the data in a lower dimension (often univariate):

$$x_{p^*}(\mathbf{s}; \mathbf{t}) \rightarrow x_p(\mathbf{s}; \mathbf{t}), \quad (6)$$

where $p^* = 1, 2, \dots, P^*$ and $p = 1, 2, \dots, P$ reduce the variable dimension from P to P^* . The most commonly used dimension reduction technique is Principal Component Analysis (PCA), also called Empirical Orthogonal Function (EOF) in earth science. It can be seen as a special case of weighting, where variables are summed up in a linear combination:

$$x_{p^*}(\mathbf{s}; \mathbf{t}) = \sum_{p=1}^P \lambda_p x_p(\mathbf{s}; \mathbf{t}),$$

with restrictions imposed on the weight coefficient: $\sum_{p=1}^P \lambda_p^2 = 1$. In other cases of weighting, the coefficients can be as simple as giving equal weight to each variables.

Some dimension reduction can also be formulated from domain-specific knowledge. This can be theories that describe the physics of the phenomenon being indexed or practical formulations used to extract useful features from the raw variables. For example, in the index SPEI, a difference series is calculated between precipitation and potential evapotranspiration (PET) and the validity of this formulation is backed up by climate water balance model [Thornthwaite, 1948], which describes [...]. Add another example of remote sensing variables i.e. $NDVI = (NIR - Red) / (NIR + Red)$?

While suggested weights and formulas can indicate norms adored by practitioners, analysts should be given the flexibility to experiment with different combinations when constructing indexes. This could help understand index behavior from its sensitivity to the variables and suggest alternative weights that better suit the specific tasks.

3.2.6 Distribution fit

model fit?

Distribution fit can be seen as the model fitting in its simplest term. It can be represented by

$$F_{\eta}(x(s; \mathbf{t})), \quad (7)$$

where $\eta \in H \subseteq \mathbb{R}^{d_{\eta}}$ is the distribution parameter and d_{η} is the number of parameter of η . A distribution fit typically aims at finding the distribution that best fits the data. Analysts may start from a pool of candidate distributions with a chosen fitting method and goodness of fit measure. While it is useful to find the ultimate best distribution to fits the data, from a probabilistic perspective, the fitting procedure itself has an uncertainty associated with the data fed and the parameter chosen. A reasonable alternative is to understand how much the index values can vary given different distributions, fitting methods, and goodness of fit tests, and whether these variations are negligible in a given application.

3.2.7 Normalising

This step maps the univariate series into a different scale, typically for ease of comparison across regions. For example, a normal scale, $[0, 1]$, or $[0, 100]$ may be favored for reporting certain indexes. In drought indexes, i.e. SPI or SPEI, the quantiles from the fitted distribution are converted into the normal scale via the normal reverse CDF function: $\Phi^{-1}(\cdot)$. Normalising is usually used at the end of the pipeline and its main difference from the scaling step is that here the change of scale also changes the distribution of the variable. While being commonly used, this step can get criticism from analysts for forcing the data into the decided scale, which can be either unnecessary or inaccurately exaggerate or downplay the outliers. Also, the use of a normal scale needs to be interpreted with caution. Figure 4 illustrates the normal density not being directly proportional to its probability of occurrence. This is concerning, especially at the extreme values, since a small difference in the tail density can have magnitudes of difference in its probability of occurrence.

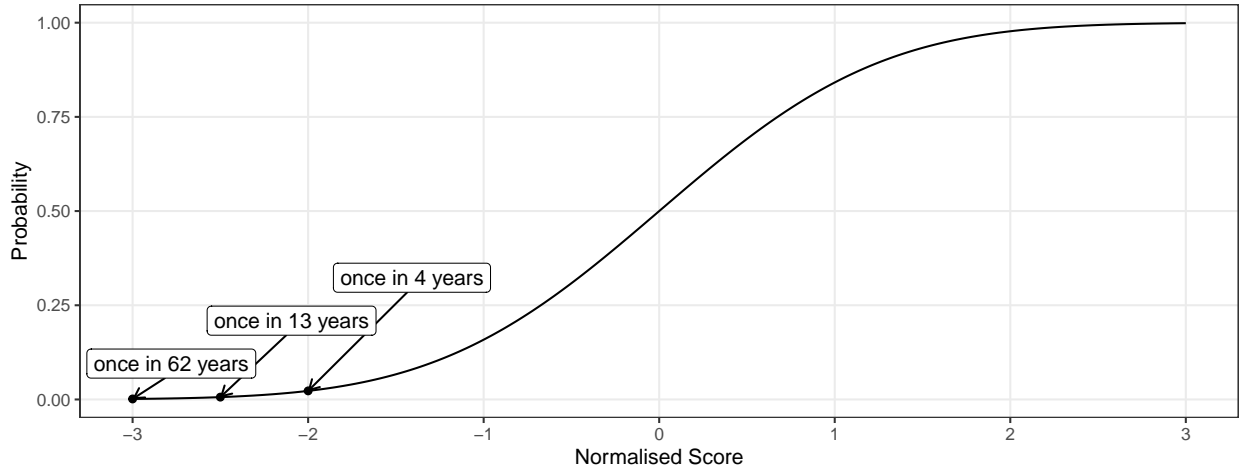


Figure 4: Scatterplot of normal quantiles against their density values. Three tail density values are highlighted with its probability of occurrence labelled. Probability is calculated assuming monthly data: with a density of -2, the probability of occurrence is $1/\text{pnorm}(-2)/12 = 4$ years. The non-linear relationship between the two quantities suggests normalised indexes need to be interpreted with caution since a slight change in the tail distribution can result in magnitudes of difference in its probability of occurrence.

3.2.8 Benchmarking

Benchmarking sets a constant value to allow the constructed index to be compared across time. Here we denote it with $u[x(s_i, t_j)]$ where u is a scalar of interest in the index constructed. A benchmark value could be a constant or a function of the data, i.e. mean.

3.2.9 Simplification

In public communication, the index values are usually accompanied by a categorical grade. The categorised grades are an ordered set of descriptive words or colors to communicate the severity or guide the comprehension of the indexes. The mapping from continuous index values to the discrete grades is called simplification in the pipeline and it can be written as a piece-wise function:

$$\begin{cases} C_0 & c_1 \leq (s_i; t_j) < c_0 \\ C_1 & c_2 \leq x(s_i; t_j) < c_1 \\ C_2 & c_3 \leq x(s_i; t_j) < c_2 \\ \dots & \\ C_z & c_z \leq x(s_i; t_j) \end{cases} \quad (8)$$

where C_0, C_1, \dots, C_z are the categories and c_0, c_1, \dots, c_z are the thresholds for each category. In SPI, droughts are sorted into four categories: mild drought: $[-0.99, 0]$; moderate drought: $[-1.49, -1]$; severe drought: $[-1.99, -1.5]$, and extreme drought: $[-\infty, -2]$. In this case, C_0, C_1, C_2, C_3 are the drought categories: mild, moderate, severe, and extreme drought ($z = 3$) and $c_0 = 0, c_1 = -1, c_2 = -1.5, c_3 = -2$ are the cutoff value for each class.

4 Incorporating alternative methods into the pipeline components

5 Examples

5.1 A drought index example

A common task for drought researchers is to compute indexes at different parameter combinations. This can be used to identify the spatial and temporal extent of drought events, recommend the best parameter choice, or compare the effectiveness of indexes for monitoring drought. The example below computes two indexes: SPI and SPEI, at various time scales and fitted distributions, for stations in the state of Queensland in Australia. The purpose of the example is to demonstrate the interfaces the tidyindex package built to allow easy computing at different parameter combinations.

The state of Queensland in Australia is frequently affected by natural disaster events such as flood and drought, which can have significant impacts on its agricultural industry. This study uses daily data from Global Historical Climatology Network Daily (GHCND), accessed via the package rnoaa to examine drought/flood condition in Queensland. Daily data is average into monthly and stations are excluded if monthly data contains missings, which is required for calculating both SPI and SPEI. This gives 29 stations with complete records from 1990 January to 2022 April.

The function `compute_indexes()` can be used to collectively compute multiple indexes. The tidyindex offers wrapper functions, with the prefix `idx_`, that simplify the calculation of commonly used indexes by combining a set of pipeline steps into a single function. For example, the function `idx_spei()` includes the five steps previously described in Section 3.1 (variable transformation, dimension reduction, temporal aggregation, distribution fit, and normalise). Each `idx_xxx()` function specifies the relevant parameters relevant to the index: the thornthwaite method is used to calculate PET in SPEI, with the average temperature (`tavg`) and latitude (`lat`) used as inputs. The SPEI is computed at four time scales (6, 12, 24, and 36 months) and fitted with two distributions (Log-logistic and General Extreme Value (GEV)). The SPI is also computed at the same four time scales and uses the default gamma distribution to fit the aggregated series.

```
.scale <- c(6, 12, 24, 36)
(idx <- queensland %>%
  init(id = id, time = ym) %>%
  compute_indexes(
    spei = idx_spei(
      .pet_method = "thornthwaite", .tavg = tavg, .lat = lat,
      .scale = .scale, .dist = c(gev(), loglogistic()),
      spi = idx_spi(.scale = .scale)
    ))
```

A tibble: 128,586 x 19

	.idx	.period	id	ym	prcp	tmax	tmin	tavg	long	lat	name	pet	diff
	<chr>	<dbl>	<chr>	<mth>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<chr>	<dbl>	<dbl>
1	spei	6	ASN00029~	1990 Jun	170	29.7	16.2	23.0	142.	-15.5	KOWA~	67.8	102.
2	spei	6	ASN00029~	1990 Jun	170	29.7	16.2	23.0	142.	-15.5	KOWA~	67.8	102.
3	spei	6	ASN00029~	1990 Jun	0	23.0	11.8	17.4	139.	-20.7	MOUN~	30.8	-30.8
4	spei	6	ASN00029~	1990 Jun	0	23.0	11.8	17.4	139.	-20.7	MOUN~	30.8	-30.8
5	spei	6	ASN00031~	1990 Jun	794	25.8	18.1	21.9	146.	-16.9	CAIR~	66.3	728.
6	spei	6	ASN00031~	1990 Jun	794	25.8	18.1	21.9	146.	-16.9	CAIR~	66.3	728.
7	spei	6	ASN00031~	1990 Jun	504	23.0	13.8	18.4	145.	-17.1	WALK~	48.4	456.
8	spei	6	ASN00031~	1990 Jun	504	23.0	13.8	18.4	145.	-17.1	WALK~	48.4	456.
9	spei	6	ASN00032~	1990 Jun	1970	23.9	16.4	20.2	146.	-17.6	SOUT~	54.5	1915.
10	spei	6	ASN00032~	1990 Jun	1970	23.9	16.4	20.2	146.	-17.6	SOUT~	54.5	1915.

... with 128,576 more rows, and 6 more variables: .scale <dbl>, .agg <dbl>,
.method <chr>, .fitted <dbl>, .dist <chr>, .index <dbl>

The output from `compute_indexes()` contains index values and associated parameter in a long tibble. It includes the original variables (`id`, `ym`, `prcp`, `tmax`, `tmin`, `tavg`, `long`, `lat`, and `name`), index parameters (`.idx`, `.scale`, `.method`, and `.dist`), intermediate variables (`pet`, `.agg`, and `.fitted`), and the final index (`.index`). This data can be visualised across space or time, or simultaneously, to explore the wet/dry condition in Queensland. Figure 5 visualises the spatial distribution of SPI at two periods (2010 October - 2011 March and 2019 October - 2020 March) with significant natural disaster events: 2010/11 Queensland flood and 2019 Australia drought, which contributes to the notorious 2019/20 bushfire. Figure 6 displays the sensitivity of the SPEI series for one particular station, Texas post office, at different time scales and fitted distributions. These two plots demonstrate some possibilities to explore the indexes after they are computed from `compute_indexes()`.

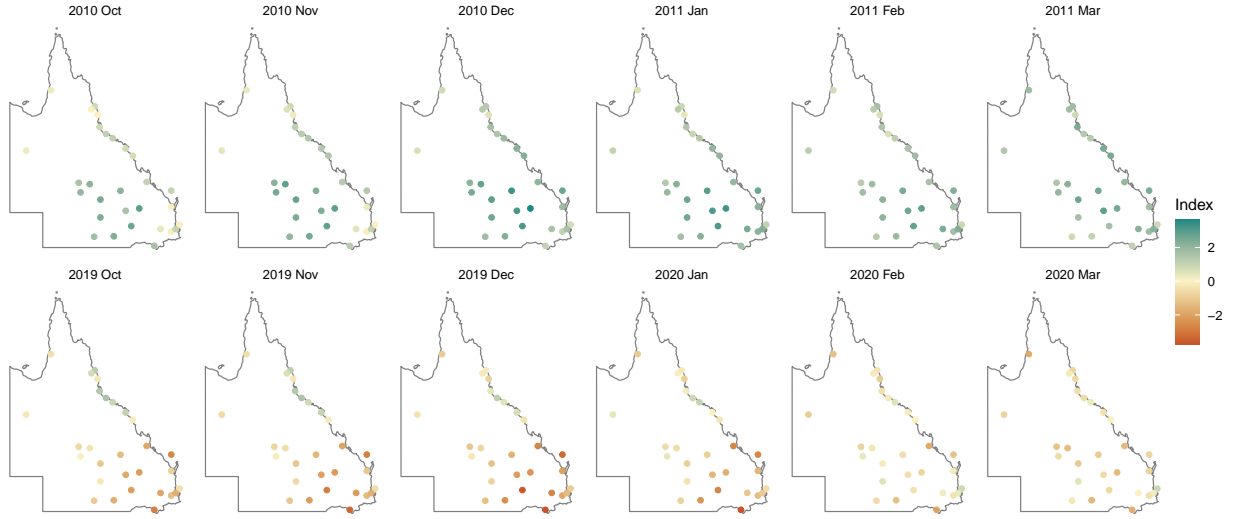


Figure 5: Spatial distribution of Standardized Precipitation Index (SPI-12) in Queensland, Australia during two major flood and drought events: 2010/11 and 2019/20. The map shows a continuous wet period during the 2010/11 flood period and a mitigated drought situation, after its worst in 2019 December and 2020 January, likely due to the increased rainfall in February from the meteorological record.

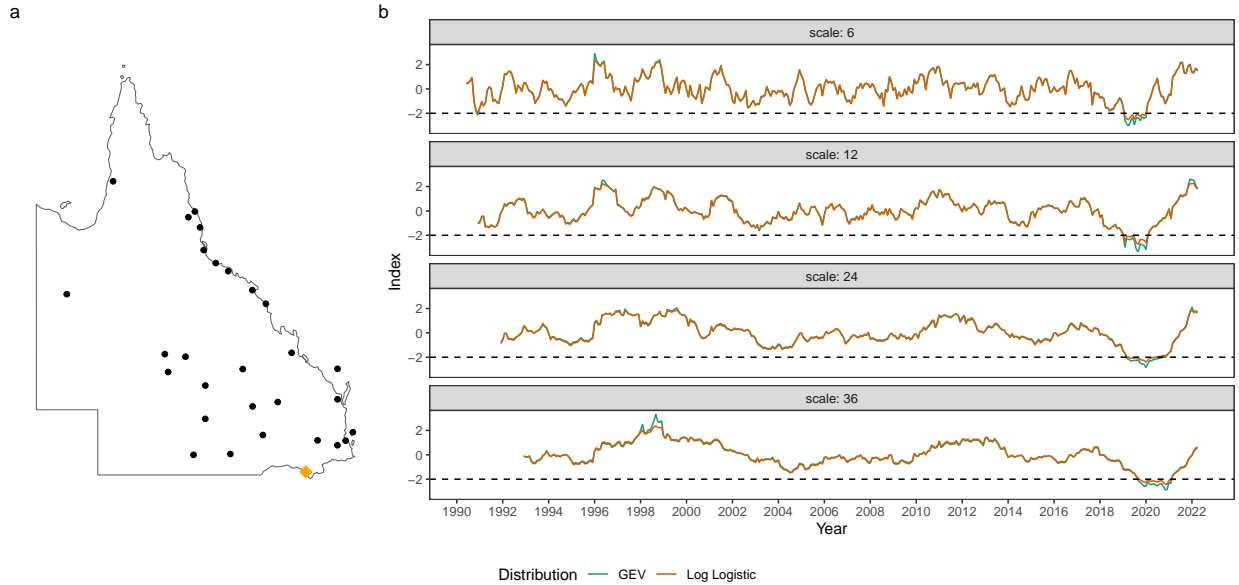


Figure 6: Time series plot of Standardized Precipitation-Evapotranspiration Index (SPEI) at the Texas post office station (highlighted by a diamond shape in panel a). The SPEI is calculated at four time scales (6, 12, 24, and 36 months) and fitted with two distributions (Log Logistic and GEV). The dashed line at -2 represents the class “extreme drought” by the SPEI. A larger time scale gives a smoother index series, while also takes longer to recover from an extreme situation as seen in the 2019/20 drought period. The SPEI values from two distribution fits mostly agree, while GEV can results in more extreme values, i.e. in 1998 and 2020.

5.2 An example on sustainable development indicators

In the following example, the Human Development Index (HDI) will first be constructed using the pipeline syntax. The section will then demonstrate how expressions or parameters in the pipeline can be swapped to alternatives to study the index property.

Human Development Index (HDI) measures the development of countries through more than just economic growth, but also people's life expectancy and opportunity to receive education. These three identified dimensions are measured using four indicators: life expectancy at birth (health), expected years of schooling (education), mean years of schooling (education), and Gross National Income (GNI) *per capita* (standard of living). The technical notes (United Nations Development Programme 2021) have documented the procedures to calculate HDI and they are summarised below:

1. take log on GNI *per capita*,
2. rescale the four indicators into [0, 1] using mini-max,
3. aggregate the two education indicators using arithmetic mean, and
4. aggregate the three dimensions into the index using geometric mean.

The values used in mini-max rescaling are summarised in Table 1 and the justification of these numbers can be found in the technical notes mentioned above.

Table 1: Maximum and minimum values used to rescale the four HDI indicators into [0, 1] range. The maximum of GNI per capita is taken as the common log of 75,000, approximately 4.875.

Dimension	Indicator	Variable	Minimum	Maximum
Health	Life expectancy at birth (years)	life_exp	20	85.000
Education	Expected years of schooling (years)	exp_sch	0	18.000
Education	Mean years of schooling (years)	avg_sch	0	15.000
Standard of living	GNI per capita (2017 PPP\$)	gni_pc	2	4.875

Among the four steps listed in the calculation, the first two are variable transformation. The next two can be considered as the dimension reduction step in the data pipeline given its intention to combine multivariate data into univariate. With the index pipeline proposed in the paper, HDI can be calculated as:

```
dt <- raw %>% init(id = country, indicators = life_exp:gni_pc)
res <- dt %>%
  var_trans(gni_pc = log(gni_pc, base = 10)) %>%
  var_trans(.method = rescale_minmax, .vars = life_exp:gni_pc,
            min = c(20, 0, 0, 2), max = c(85, 18, 15, log10(75000))) %>%
  dim_red(sch = (exp_sch + avg_sch) / 2) %>%
  dim_red(index = (life_exp * sch * gni_pc)^(1/3))
```

Index analysts can be interested in studying the property of the index and the countries through the “what-if” questions. In the context of HDI, this can be what if an alternative dimension reduction expression is used to reduce the three dimensions into the HDI? Apart from the geometrical mean, linear combination is a common method to reduce the data dimension. With different weighting schemes, the combination can be compared to study the variable importance of each dimension or to reveal the underlying structures of countries. The following code tests five expressions (arithmetic mean, principal component analysis weight, and three weighting schemes emphasizing life expectancy, education, and GNI per capita respectively) using the `swap_exprs()` function. The function first locates the operation to be tested (`.var = index`), then accepts a list of expressions for testing `.expr`, before requesting the raw data:

```
(res2 <- res %>%
  swap_exprs(
    .var = index,
    .expr = list(
      index1 = (life_exp + sch + gni_pc)/3,
      # PCA recommended weight
      index2 = 0.569 * life_exp + 0.576 * sch + 0.586 * gni_pc,
      index3 = 0.8 * life_exp + 0.1 * sch + 0.1 * gni_pc,
```

```

index4 = 0.1 * life_exp + 0.8 * sch + 0.1 * gni_pc,
index5 = 0.1 * life_exp + 0.1 * sch + 0.8 * gni_pc),
.raw_data = dt))

```

Index pipeline:

Summary:

```

var_trans: `log(gni_pc, base = 10)` -> gni_pc
var_trans: `rescale_minmax(var = life_exp, min = 20, max = 85)` -> life_exp
var_trans: `rescale_minmax(var = exp_sch, min = 0, max = 18)` -> exp_sch
var_trans: `rescale_minmax(var = avg_sch, min = 0, max = 15)` -> avg_sch
var_trans: `rescale_minmax(var = gni_pc, min = 2, max = 4.875061)` -> gni_pc
dim_red: `(exp_sch + avg_sch)/2` -> sch
dim_red: `(life_exp * sch * gni_pc)^(1/3)` -> index0
dim_red: `(life_exp + sch + gni_pc)/3` -> index1
dim_red: `0.569 * life_exp + 0.576 * sch + 0.586 * gni_pc` -> index2
dim_red: `0.8 * life_exp + 0.1 * sch + 0.1 * gni_pc` -> index3
dim_red: `0.1 * life_exp + 0.8 * sch + 0.1 * gni_pc` -> index4
dim_red: `0.1 * life_exp + 0.1 * sch + 0.8 * gni_pc` -> index5

```

Data:

```

# A tibble: 191 x 13
   id country    life_~1 exp_sch avg_sch gni_pc  sch index0 index1 index2 index3 index4
  <dbl> <chr>      <dbl>   <dbl>   <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
1     1  Switzerland 0.984   0.917   0.924 0.983 0.920 0.962 0.963 1.67 0.978 0.933
2     2    Norway 0.973     1   0.867 0.978 0.933 0.961 0.961 1.66 0.969 0.942
3     3   Iceland 0.964     1   0.918 0.955 0.959 0.959 0.959 1.66 0.963 0.959
4     4 Hong Kong 1     0.960 0.815 0.973 0.887 0.952 0.953 1.65 0.986 0.907
5     5 Australia 0.993     1   0.848 0.936 0.924 0.951 0.951 1.65 0.980 0.932
6     6  Denmark 0.944     1   0.864 0.967 0.932 0.948 0.948 1.64 0.945 0.937
7     7  Sweden 0.969     1   0.841 0.952 0.920 0.947 0.947 1.64 0.962 0.928
8     8  Ireland 0.954     1   0.772 1     0.886 0.945 0.947 1.64 0.952 0.904
9     9  Germany 0.933 0.945 0.939 0.952 0.942 0.942 0.942 1.63 0.936 0.942
10    10 Netherla~ 0.949     1   0.839 0.956 0.919 0.941 0.941 1.63 0.947 0.926
# ... with 181 more rows, 1 more variable: index5 <dbl>, and abbreviated variable name
#   1: life_exp

```

The resulting data includes five more columns corresponding to the five alternative indexes. The country ranking can be further computed on the data and plotted in a scatterplot matrix as in Figure 7. The plot shows almost similar ranks suggested by the indexes calculated using geometric mean (rank0), the arithmetic mean (rank1), and the PCA recommended weights (rank2). The agreement of the three methods supports the use of geometric mean as the dimension reduction method in the original index. The scatterplot matrix also reveals the difference in ranking when the three dimensions are given different weights in rank3, rank4, and rank5. The panel rank4 vs rank5 shows a cluster of highlighted points with relatively high ranks according to rank5 but low in rank4. These countries are printed in Table 2 prints along with their three dimensions, calculated indexes, and rankings. While having a high life expectancy and GNI *per capita*, these countries (Andorra, Qatar, San Marino, Kuwait, and Brunei Darussalam) score relatively low in education, which is weighted heavily in index4 ($0.1 * \text{life_exp} + 0.1 * \text{sch} + 0.8 * \text{gni_pc}$).

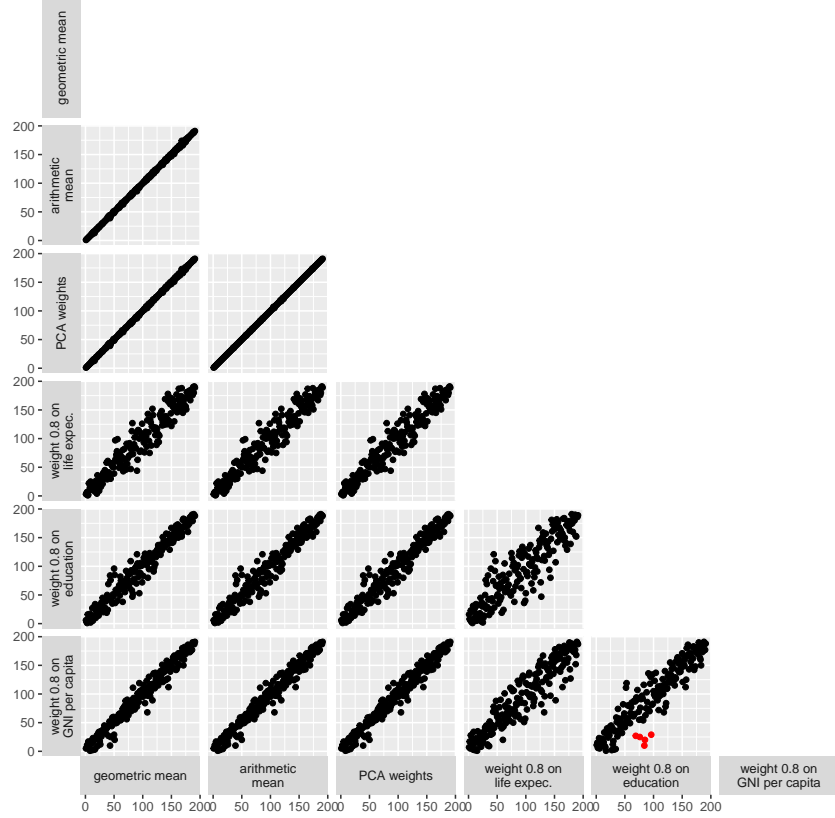


Figure 7: Comparing country ranks from six different index constructions: geometric mean (rank0), arithmetic mean (rank1), principal component analysis recommended weights (rank2), and heavier weight on life expectancy, education, and GNI per capita (rank3 - rank5). Geometric mean, arithmetic mean, and PCA produce similar rankings where countries are aligned in a diagonal line, while country rankings varies when compare among rank3 to rank5. Panel rank5 vs. rank4 highlights a cluster of countries that rank high by rank5 but low by rank4. The plot demonstrates the necessity to test on various alternative expressions to verify the robustness of the method and to capture distinctive characteristics of different countries.

Table 2: A selected number of countries with low rankings when education is given a heavy weight (index4) while ranks high when a heavier weight is given on GNI per capital (index5).

country	life expectancy	education	GNI per capita	index3	index4	rank3	rank4
Andorra	0.929	0.721	0.942	0.909	0.764	32	69
Qatar	0.912	0.684	1.000	0.898	0.739	34	84
San Marino	0.937	0.701	0.947	0.914	0.749	30	76
Kuwait	0.903	0.670	0.947	0.884	0.721	39	96
Brunei Darussalam	0.841	0.694	0.977	0.840	0.737	60	85

6 Conclusion

The paper presents a data pipeline with nine modules for constructing and analysing indexes. The pipeline increases transparency in the practice for index analysts to experiment with different index design and parameter choices to better design and apply their indexes. The significance of this work is its ability to provide a universal framework for index construction, which can be applied across different domains.

Examples have been given in the drought indexes and human development index to demonstrate computing of indexes with different parameters combinations and how alternative index design can provide insights to understand distinctive country characteristics that could sometimes be overlooked. The accompanied package, *tidyindex*, is not meant to provide comprehensive implementation for all indexes across all domains. Instead, it demonstrates implementing individual pipeline steps that are versatile to multiple indexes and composing new indexes from existing steps. Domain experts are welcomed to adopt the pipeline approach to develop specialised packages for specific-domains indexes.

Reference

- Brown, Jesslyn F., Brian D. Wardlow, Tsegaye Tadesse, Michael J. Hayes, and Bradley C. Reed. 2008. “The Vegetation Drought Response Index (VegDRI): A New Integrated Approach for Monitoring Drought Stress in Vegetation.” *GIScience & Remote Sensing* 45 (1): 16–46. <https://doi.org/10.2747/1548-1603.45.1.16>.
- Buja, A, D Asimov, C Hurley, and JA McDonald. 1988. “Elements of a Viewing Pipeline for Data Analysis.” In *Dynamic Graphics for Statistics*, 277–308. Wadsworth, Belmont.
- Economist Intelligence Unit. 2019. “The Global Liveability Index 2019.” *The Economist*. <https://www.cbeinternational.ca/pdf/Liveability-Free-report-2019.pdf>.
- Keyantash, John A., and John A. Dracup. 2004. “An Aggregate Drought Index: Assessing Drought Severity Based on Fluctuations in the Hydrologic Cycle and Surface Water Storage.” *Water Resources Research* 40 (9). <https://doi.org/10.1029/2003WR002610>.
- Kogan, F. N. 1995. “Application of Vegetation Index and Brightness Temperature for Drought Detection.” *Advances in Space Research*, Natural Hazards: Monitoring and Assessment Using Remote Sensing Technique, 15 (11): 91–100. [https://doi.org/10.1016/0273-1177\(95\)00079-T](https://doi.org/10.1016/0273-1177(95)00079-T).
- McKee, Thomas B, Nolan J Doesken, John Kleist, et al. 1993. “The Relationship of Drought Frequency and Duration to Time Scales.” In *Proceedings of the 8th Conference on Applied Climatology*, 17:179–83. 22. Boston, MA, USA.
- Sutherland, Peter, Anthony Rossini, Thomas Lumley, Nicholas Lewin-Koh, Julie Dickerson, Zach Cox, and Dianne Cook. 2000. “Orca: A Visualization Toolkit for High-Dimensional Data.” *Journal of Computational and Graphical Statistics* 9 (3): 509–29. <https://www.jstor.org/stable/1390943>.
- Tucker, Compton J. 1979. “Red and Photographic Infrared Linear Combinations for Monitoring Vegetation.” *Remote Sensing of Environment* 8 (2): 127–50. [https://doi.org/10.1016/0034-4257\(79\)90013-0](https://doi.org/10.1016/0034-4257(79)90013-0).
- United Nations Development Programme. 2021. “Human Development Report 2021-2022: The Next Frontier: Human Development and the Anthropocene.” https://hdr.undp.org/sites/default/files/2021-22_HDR/hdr2021-22_technical_notes.pdf.
- . 2022. “Human Development Report 2021-22.” New York. <http://report.hdr.undp.org>.
- Vicente-Serrano, Sergio M., Santiago Beguería, and Juan I. López-Moreno. 2010. “A Multiscalar Drought Index Sensitive to Global Warming: The Standardized Precipitation Evapotranspiration Index.” *Journal of Climate* 23 (7): 1696–1718. <https://journals.ametsoc.org/view/journals/clim/23/7/2009jcli2909.1.xml>.
- Wickham, Hadley. 2014. “Tidy Data.” *Journal of Statistical Software* 59 (September): 1–23. <https://doi.org/10.18637/jss.v059.i10>.
- Wickham, Hadley, Michael Lawrence, Dianne Cook, Andreas Buja, Heike Hofmann, and Deborah F. Swayne. 2009. “The Plumbing of Interactive Graphics.” *Computational Statistics* 24 (2): 207–15. <https://doi.org/10.1007/s00180-008-0116-x>.
- Xie, Yihui, Heike Hofmann, and Xiaoyue Cheng. 2014. “Reactive Programming for Interactive Graphics.” *Statistical Science* 29 (2): 201–13. <https://www.jstor.org/stable/43288470?seq=1>.