

# A Guide to the agop 0.1-0 Package for R

## Aggregation Operators in R

Marek Gagolewski<sup>1,2</sup>, Anna Cena<sup>1,2</sup>

<sup>1</sup> Systems Research Institute, Polish Academy of Sciences  
ul. Newelska 6, 01-447 Warsaw, Poland

<sup>2</sup> Rexamine, Email: {gagolews,cena}@rexamine.com  
[www.rexamine.com/resources/agop/](http://www.rexamine.com/resources/agop/)

June 25, 2013

*The package, as well as this tutorial, is still in its early days – any suggestions are welcome!*

---

## Contents

<b>1</b>	<b>Getting started</b>	<b>1</b>
<b>2</b>	<b>Theoretical Background</b>	<b>2</b>
2.1	A Note on Representing Numeric Data and Applying Operations in R . . . . .	3
2.2	A Note on Storing Multiple Numeric Vectors in R . . . . .	4
2.3	Aggregation Operators and Their Basic Properties . . . . .	5
2.4	Impact Functions and The Producers Assessment Problem . . . . .	6
<b>3</b>	<b>Visualization</b>	<b>8</b>
3.1	Depicting producers . . . . .	8
<b>4</b>	<b>Pre-orders</b>	<b>8</b>
<b>5</b>	<b>Predefined Classes of Aggregation Operators in agop</b>	<b>12</b>
5.1	A Review of Notable Classes of Aggregation Operators . . . . .	12
5.2	Interesting Impact Functions . . . . .	15
<b>6</b>	<b>Aggregation Operators from the Probabilistic Perspective</b>	<b>17</b>
6.1	Some Notable Probability Distributions . . . . .	17
6.1.1	Pareto-Type II Distribution . . . . .	18
6.2	Stochastic Properties of Aggregation Operators . . . . .	20
<b>7</b>	<b>NEWS/CHANGELOG</b>	<b>21</b>
	<b>Bibliography</b>	<b>22</b>

## 1 Getting started

*“The process of combining several numerical values into a single representative one is called **aggregation**, and the numerical function performing this process is called **aggregation function**. This simple definition demonstrates the size of the field of application of aggregation: applied mathematics (e.g. probability, statistics, decision theory), computer science (e.g. artificial intelligence, operation research), as well as*

*many applied fields (economics and finance, pattern recognition and image processing, data fusion, multicriteria decision making, automated reasoning etc.). Although history of aggregation is probably as old as mathematics (think of the arithmetic mean), its existence has reminded underground till only recent (...).” [24, p. xiii]*

R [31] is a free, open source software environment for statistical computing and graphics, which includes an implementation of a very powerful and quite popular high-level language called S. It runs on all major operating systems, i.e. Windows, Linux, and MacOS X. To install R and/or find some information on the S language please visit R Project’s Homepage at [www.R-project.org](http://www.R-project.org). Perhaps you may also wish to install RStudio, a convenient development environment for R. It is available at [www.rstudio.org](http://www.rstudio.org).

*agop* is an open source (licensed under GNU LGPL 3) package for  $R \geq 2.12$  to which anyone can contribute. It started as a fork of the CITAN (*Citation Analysis Toolpack*, [13]) package.

To install latest “official” release of the package available on *CRAN* we type<sup>1</sup>:

```
install.packages('agop')
```

Alternatively, we may fetch its current development snapshot from *GitHub*:

```
install.packages('devtools')  
library('devtools')  
install_github('agop', 'Rexamine')
```

Note that in this case you will need a working C/C++ compiler.

Each session with *agop* should be preceded by a call to:

```
library('agop') # Load the package
```

To view the main page of the manual we type:

```
library(help='agop')
```

For more information please visit the package’s homepage [18]. In case of any problems, comments, or suggestions feel free to contact the authors. Good luck!

## 2 Theoretical Background

Let us establish some basic notation convention used throughout this tutorial. From now on let  $\mathbb{I} = [a, b]$ , possibly with  $a = -\infty$  or  $b = \infty$ . Note that in many practical situations we commonly choose  $\mathbb{I} = [-1, 1]$ ,  $\mathbb{I} = [0, 1]$  or  $\mathbb{I} = [0, \infty]$ . A set of all vectors of arbitrary length with elements in  $\mathbb{I}$  is denoted by  $\mathbb{I}^{1,2,\dots} = \bigcup_{n=1}^{\infty} \mathbb{I}^n$ .

For two equal-length vectors  $\mathbf{x}, \mathbf{y} \in \mathbb{I}^n$  we write  $\mathbf{x} \leq \mathbf{y}$  if and only if for all  $i = 1, \dots, n$  it holds  $x_i \leq y_i$ . Moreover, all binary arithmetic operations on vectors  $\mathbf{x}, \mathbf{y} \in \mathbb{I}^n$  will be performed element-wise, e.g.  $\mathbf{x} + \mathbf{y} = (x_1 + y_1, \dots, x_n + y_n) \in \mathbb{I}^n$ . Similar behavior is assumed for  $-$ ,  $\cdot$ ,  $/$ ,  $\wedge$  (min),  $\vee$  (max), etc. Additionally, each function of one variable  $f : \mathbb{I} \rightarrow \mathbb{I}$  can be extended to the vector space: we write  $f(\mathbf{x})$  to denote  $(f(x_1), \dots, f(x_n))$ .

---

<sup>1</sup>You are viewing the **development** version of the tutorial. Some of the features presented in this document may be missing in the CRAN release. Please, upgrade to the **latest** development version from *GitHub* if you need the new functionality. Note that you will need a working C/C++ compiler.

Let  $x_{(i)}$  denote the  $i$ th order statistic, i.e. the  $i$ th smallest value in  $\mathbf{x}$ . Moreover, for convenience, let  $x_{\{i\}} = x_{|\mathbf{x}|-i+1}$  denote the  $i$ th greatest value in  $\mathbf{x}$ .

For any  $n \in \mathbb{N}$  and  $c \in \mathbb{I}$ , we set  $(n * c) = (c, \dots, c) \in \mathbb{I}^n$ . Also,  $[n] := \{1, 2, \dots, n\}$  with  $[0] = \emptyset$ .

## 2.1 A Note on Representing Numeric Data and Applying Operations in R

Recall how we create numeric vectors in R:

```
(x1 <- c(5, 2, 3, 1, 0, 0))
## [1] 5 2 3 1 0 0
class(x1)
## [1] "numeric"
(x2 <- 10:1) # the same as seq(10, 1)
## [1] 10 9 8 7 6 5 4 3 2 1
(x3 <- seq(1, 5, length.out=6))
## [1] 1.0 1.8 2.6 3.4 4.2 5.0
(x4 <- seq(1, 5, by=1.25))
## [1] 1.00 2.25 3.50 4.75
```

To obtain  $(n * c)$ , e.g. for  $n = 10$  and  $c = 3$ , we call:

```
rep(10, 3)
## [1] 10 10 10
```

Note that in R all the arithmetic operations on vectors are performed element-wise, i.e. in a manner indicated above. This is called **vectorization**. The same holds for mathematical functions: they are extended to the vector space.

```
x <- c(1, 3, 3, 2)
y <- c(2, 3, -1, 0)
x+y
## [1] 3 6 2 2
x*y
## [1] 2 9 -3 0
pmin(x,y)
## [1] 1 3 -1 0
pmax(x,y)
## [1] 2 3 3 2
abs(y)
## [1] 2 3 1 0
```

Thus, we calculated  $\mathbf{x} + \mathbf{y}$ ,  $\mathbf{x} \cdot \mathbf{y}$ ,  $\mathbf{x} \wedge \mathbf{y}$ ,  $\mathbf{x} \vee \mathbf{y}$ , and  $|\mathbf{x}|$  (try to determine yourself what happens if we deal with vectors of unequal length in R).

Moreover, for the  $\leq$  relation we write:

```
all(x <= y)
## [1] FALSE
```

To get  $x_{\{i\}}$  we have to sort the given vector nonincreasingly.

```
(xs <- sort(x, decreasing=TRUE)) # `decreasing' may be misleading
## [1] 3 3 2 1
xs[3] # the third greatest value in x
## [1] 2
```

## 2.2 A Note on Storing Multiple Numeric Vectors in R

Vectors of the same length can be conveniently stored in a matrices. Please note that the `dimnames` attribute of a matrix may define its row and column labels. Its value may be set to `NULL` (no names given) or to a list with two character vectors (rows and columns names, respectively). Another simple way to set the labels is by using the `rownames()` and `colnames()` functions.

The `apply()` function may be called to evaluate a given method on each matrix row or column (parameter `MARGIN` set to 1 and 2, respectively).

```
expertopinions <- matrix(c(
  6,7,2,3,1, # this will be the first COLUMN
  8,3,2,1,9, # 2nd
  4,2,4,1,6  # 3rd
),
ncol=3,
dimnames=list(NULL, c("A", "B", "C"))) # only column names set
class(expertopinions)
## [1] "matrix"
print(expertopinions) # or print(authors)
##      A B C
## [1,] 6 8 4
## [2,] 7 3 2
## [3,] 2 2 4
## [4,] 3 1 1
## [5,] 1 9 6
apply(expertopinions, 2, mean) # apply the mean() function on each COLUMN
##      A      B      C
## 3.8 4.6 3.4
```

Vectors that are not of the same length may be store in a list (with possibly named elements). In that case, the functionality of `apply()` is provided by `lapply()` or `sapply()` functions.

```
authors <- list(
  "John S." = c(7,6,2,1,0),
  "Kate F." = c(9,8,7,6,4,1,1,0)
)
class(authors)
```

```
## [1] "list"

str(authors) # or print(authors)

## List of 2
## $ John S.: num [1:5] 7 6 2 1 0
## $ Kate F.: num [1:8] 9 8 7 6 4 1 1 0

index_h(authors[[1]]) # the h-index /see below/ for 1st author
## [1] 2

sapply(authors, index_h) # calculate the h-index for all vectors in a list
## John S. Kate F.
##      2      4

index_h(authors) # index_h() expects an numeric vector on input
## Error: argument 'x' should be a numeric vector (or an object coercible to)
```

## 2.3 Aggregation Operators and Their Basic Properties

Dealing with huge amounts of data faces us with the problem of constructing their synthetic descriptions. The aggregation theory, a relatively new research domain at the border of mathematics and computer science, is interested in the analysis of functions that may be used in this task. Thus, we should start with the formal definition of objects of our interest. Here is the most general setting:

**Definition 1.** A function  $F : \mathbb{I}^{1,2,\dots} \rightarrow \mathbb{I}$  is called an (*extended*<sup>2</sup>) **aggregation operator** if it is at least **nondecreasing** in each variable, i.e. for all  $n$  and  $\mathbf{x}, \mathbf{y} \in \mathbb{I}^n$  if  $\mathbf{x} \leq \mathbf{y}$ , then  $F(\mathbf{x}) \leq F(\mathbf{y})$ .

Note that each aggregation operator is a mapping into  $\mathbb{I}$ , thus for all  $n$  we have  $\inf_{\mathbf{x} \in \mathbb{I}^n} F(\mathbf{x}) \geq a$  and  $\sup_{\mathbf{x} \in \mathbb{I}^n} F(\mathbf{x}) \leq b$ . By nondecreasingness, however, these conditions reduce to  $F(n * a) \geq a$  and  $F(n * b) \leq b$ .

Also keep in mind that some authors assume (cf. [24]) that aggregation operators must fulfill the two following **strong boundary conditions**: for all  $n$  we have  $\inf_{\mathbf{x} \in \mathbb{I}^n} F(\mathbf{x}) = a$  and  $\sup_{\mathbf{x} \in \mathbb{I}^n} F(\mathbf{x}) = b$ . In our case, this does not necessarily hold – we want to be more general.

Here are some interesting properties of aggregation operators. Later on we will characterize the classes of aggregation operators that fulfill them.

**Definition 2.** We call  $F : \mathbb{I}^{1,2,\dots} \rightarrow \mathbb{I}$  **symmetric** if:

$$(\forall n \in \mathbb{N}) (\forall \mathbf{x}, \mathbf{y} \in \mathbb{I}^n) \mathbf{x} \cong \mathbf{y} \implies F(\mathbf{x}) = F(\mathbf{y}),$$

where  $\mathbf{x} \cong \mathbf{y}$  if and only if there exists a permutation  $\sigma$  of  $[n]$  such that  $\mathbf{x} = (y_{\sigma(1)}, \dots, y_{\sigma(n)})$ .

It may be shown, see [24, Thm. 2.34], that  $F : \mathbb{I}^n \rightarrow \mathbb{I}$  is symmetric if and only if there exists a function  $G : \mathbb{I}^{1,2,\dots} \rightarrow \mathbb{I}$  such that  $F(x_1, \dots, x_n) = G(x_{\{1\}}, \dots, x_{\{n\}})$ , i.e. it may be defined solely using order statistics. Intuitively,  $F$  is symmetric if its value is independent of the aggregated vector’s elements presentation.

Idempotence is well-known from algebra, where we say that element  $x$  is idempotent with respect to binary operator  $*$  if we have  $x * x = x$ . The following definition extends this property to  $n$ -ary aggregation functions, cf. [24].

---

<sup>2</sup>Extended to the space of vectors of arbitrary length, cf. e.g. [4, 24]; Classical approach considers only fixed-length vectors. In *agop* we are as much general as possible.

**Definition 3.** We call  $F : \mathbb{I}^{1,2,\dots} \rightarrow \mathbb{I}$  **idempotent** if:

$$(\forall n \in \mathbb{N}) (\forall x \in \mathbb{I}) F(n * x) = x.$$

Idempotent aggregation operators fulfilling the strong boundary conditions (see p. 5) are sometimes called **averaging functions**, cf. [24].

An example of such object is the arithmetic mean or median.

**Definition 4.** We call  $F : \mathbb{I}^{1,2,\dots} \rightarrow \mathbb{I}$  **additive** if:

$$F(\mathbf{x} + \mathbf{y}) = F(\mathbf{x}) + F(\mathbf{y}),$$

for all  $n \in \mathbb{N}, \mathbf{x}, \mathbf{y} \in \mathbb{I}^n$  such that  $\mathbf{x} + \mathbf{y} \in \mathbb{I}^n$ .

Please note that for  $a \leq 0$ , if  $F$  is additive, then necessarily it holds  $F(\mathbf{0}) = 0$ .

**Definition 5.** We call  $F$  **minitive** if:

$$(\forall n \in \mathbb{N}) (\forall \mathbf{x}, \mathbf{y} \in \mathbb{I}^n) F(\mathbf{x} \wedge \mathbf{y}) = F(\mathbf{x}) \wedge F(\mathbf{y}).$$

**Definition 6.** We call  $F$  **maxitive** if:

$$(\forall n \in \mathbb{N}) (\forall \mathbf{x}, \mathbf{y} \in \mathbb{I}^n) F(\mathbf{x} \vee \mathbf{y}) = F(\mathbf{x}) \vee F(\mathbf{y}).$$

**Definition 7.** We call  $F$  **modular** (cf. [3, 24, 28]) if:

$$(\forall n \in \mathbb{N}) (\forall \mathbf{x}, \mathbf{y} \in \mathbb{I}^n) F(\mathbf{x} \vee \mathbf{y}) + F(\mathbf{x} \wedge \mathbf{y}) = F(\mathbf{x}) + F(\mathbf{y})$$

It may easily be seen that each additive operator is also modular (i.e. modularity is more general than additivity), because for any additive aggregation operator  $F$ , since  $(\mathbf{x} \vee \mathbf{y}) + (\mathbf{x} \wedge \mathbf{y}) = \mathbf{x} + \mathbf{y}$ , we have  $F(\mathbf{x}) + F(\mathbf{y}) = F(\mathbf{x} + \mathbf{y}) = F((\mathbf{x} \vee \mathbf{y}) + (\mathbf{x} \wedge \mathbf{y})) = F(\mathbf{x} \vee \mathbf{y}) + F(\mathbf{x} \wedge \mathbf{y})$ .

Apart from the “ordinary” minitivity, maxitivity, and modularity we may introduce their symmetrized versions, using  $\mathbf{x} \overset{S}{+} \mathbf{y} = (x_{\{1\}} + y_{\{1\}}, \dots, x_{\{n\}} + y_{\{n\}})$ ,  $\mathbf{x} \overset{S}{\vee} \mathbf{y} = (x_{\{1\}} \vee y_{\{1\}}, \dots, x_{\{n\}} \vee y_{\{n\}})$  and  $\mathbf{x} \overset{S}{\wedge} \mathbf{y} = (x_{\{1\}} \wedge y_{\{1\}}, \dots, x_{\{n\}} \wedge y_{\{n\}})$ .

## 2.4 Impact Functions and The Producers Assessment Problem

We already noticed the important class of aggregation operators: the averaging functions. They may be used to represent the most “typical” value of a numeric vector. Here is another interesting class that represents solutions to some very interesting practical issue.

The **Producers Assessment Problem** (PAP, [22]) concerns evaluation of a set of **producers** (e.g. scientists, artists, writers, craftsman) according to some quality or popularity **ratings** of **products** (e.g. scientific articles, works, books, artifacts) that were outputted by an entity.

PAP instances may be found in many real-life situations, like those encountered for example in scientometrics, webometrics, marketing, manufacturing, or quality engineering, see Table 1 and e.g. [10]. Our main interest here is focused on constructing and analyzing aggregation operators which may be used in the producers’ rating task. Such functions should take into account the two following aspects of a producer’s quality:

- his/her ability to output highly-rated products,

**Tab. 1.** The Producer Assessment Problem – typical instances

	Producer	Products	Rating method	Discipline
A	Scientist	Scientific articles	Number of citations	Scientometrics
B	Scientific institute	Scientists	The $h$ -index	Scientometrics
C	Web server	Web pages	Number of in-links	Webometrics
D	Artist	Paintings	Auction price	Auctions
E	Billboard company	Advertisements	Sale results	Marketing
F	R package author	Packages	PageRank values on the citation graph	Software Engineering

- his/her overall productivity.

For the sake of illustration, we will consider PAP in the scientometric context, where scientists “produce” papers that are cited by peers.

Let  $\mathbb{I} = [0, \infty]$  represent the set of values that some a priori chosen paper quality measure may take. These may of course be non-integers, for example when we consider citations normalized with respect to the number of papers’ authors.

It is widely accepted, see e.g. [39, 38, 37, 32, 29, 30, 22, 12, 11], that each aggregation operator  $F : \mathbb{I}^{1,2,\dots} \rightarrow \mathbb{I}$  to be applied in PAP should at least be:

- nondecreasing in each variable (additional citations received by a paper or an improvement of its quality measure does not result in a decrease of the authors’ overall evaluation),
- arity-monotonic (by publishing a new paper we never decrease the overall valuation of the entity),
- symmetric (independent of the order of elements’ presentation, i.e. we may always assume that we aggregate vectors that are already sorted).

More formally, axiom (b) is fulfilled iff for any  $\mathbf{x} \in \mathbb{I}^{1,2,\dots}$  and  $y \in \mathbb{I}$  it holds  $F(\mathbf{x}) \leq F(x_1, \dots, x_n, y)$ . It may be seen that this property is **arity-dependent**, i.e. it takes into account the number of elements to be aggregated.

Moreover, (a) and (c) were defined in the previous section.

Here is a bunch of arity-dependent properties that can be useful while aggregating vectors of varying lengths.

**Definition 8.** We call  $F \in \mathcal{E}(\mathbb{I})$  a **zero-insensitive** aggregation operator if for each  $\mathbf{x} \in \mathbb{I}^{1,2,\dots}$  it holds  $F(\mathbf{x}, 0) = F(\mathbf{x})$ .

It may be seen that, under nondecreasingness, zero-insensitivity implies arity-monotonicity, see [20]. What is interesting, each zero-insensitive impact function  $F$  may be defined by means of  $G : \mathbb{I}^\infty \rightarrow \mathbb{I}$  such that  $F(\mathbf{x}) = G(\mathbf{x}, 0, 0, \dots)$ , i.e. of function which domain is the space of vectors of infinite length.

Zero-sensitivity may be strengthened as follows, cf. [20] and [39, Axiom A1].

**Definition 9.**  $F \in \mathcal{E}(\mathbb{I})$  is  **$F$ -insensitive** if

$$(\forall \mathbf{x} \in \mathbb{I}^{1,2,\dots}) (\forall y \in \mathbb{I}) y \leq F(\mathbf{x}) \implies F(\mathbf{x}, y) = F(\mathbf{x}).$$

Note that the above property was called R-stability in [2].

**Definition 10.**  $F \in \mathcal{E}(\mathbb{I})$  is  **$F$ -sensitive** if

$$(\forall \mathbf{x} \in \mathbb{I}^{1,2,\dots}) (\forall y \in \mathbb{I}) y > F(\mathbf{x}) \implies F(\mathbf{x}, y) > F(\mathbf{x}).$$

### 3 Visualization

#### 3.1 Depicting producers

The `plot_producer()` function may be used to draw a graphical representation of a given numeric vector, i.e. what is sometimes called a *citation function* in scientometrics.

A given vector  $\mathbf{x} = (x_1, \dots, x_n)$  can be represented by a step function defined for  $0 \leq y < n$  and given by:

$$\pi(y) = x_{(n-\lfloor y \rfloor + 1)}.$$

This function may be obtained by setting `type='right.continuous'` argument in `plot_producer()`. Recall that  $x_{(i)}$  denotes  $i$ -th smallest value in  $\mathbf{x}$ .

On the other hand, for `type='left.continuous'` (the default), we get

$$\pi(y) = x_{(n-\lfloor y \rfloor + 1)}$$

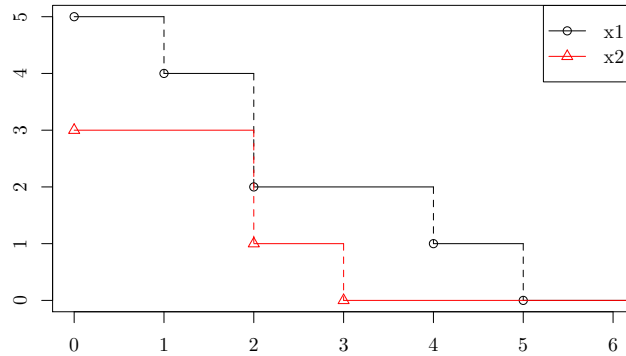
for  $0 < y \leq n$ .

Moreover, this function may depict the curve joining the sequence of points  $(0, x_{(n)}), (1, x_{(n)}), (1, x_{(n-1)}), (2, x_{(n-1)}), \dots, (n, x_{(1)})$ .

The `plot_producer()` function behaves much like the well-known R's `plot.default()` and allows for passing all its graphical parameters.

For example, let us depict the state of two given producers,  $\mathbf{x}^{(1)}$  and  $\mathbf{x}^{(2)}$ .

```
x1 <- c(5, 4, 2, 2, 1)
x2 <- c(3, 3, 1, 0, 0, 0, 0)
plot_producer(x1, extend=TRUE)
plot_producer(x2, add=TRUE, col=2, pch=2, extend=TRUE)
legend('topright', c('x1', 'x2'), col=c(1, 2), lty=1, pch=c(1, 2))
```



### 4 Pre-orders

Let us consider the following relation on  $\mathbb{I}^{1,2,\dots}$ . For any  $\mathbf{x} \in \mathbb{I}^n$  and  $\mathbf{y} \in \mathbb{I}^m$  we write  $\mathbf{x} \preceq \mathbf{y}$  if and only if  $n \leq m$  and  $x_{\{i\}} \leq y_{\{i\}}$  for all  $i \in \min\{n, m\}$ . Of course,  $\preceq$  is a pre-order – it would have been a partial order, if we had defined it on the set of *sorted* vectors.

Intuitively, we say that an author (scientometric context again)  $X$  is (weakly) dominated by an author  $Y$ , if  $X$  has no more papers than  $Y$  and each the  $i$ th most cited paper of  $X$  has no



more citations than the  $i$ th most cited paper of  $Y$ . Not that the  $(m - n)$  least cited  $Y$ ’s papers are not taken into account here. Most importantly, however, there exist pairs of vectors that are *incomparable* with respect to  $\preceq$  (see the illustration below).

This pre-order in *agop* may be determined using the `pord_weakdom()` function.

```
c(pord_weakdom(5:1, 10:1), pord_weakdom(10:1, 5:1)) # 5:1 <= 10:1
## [1] TRUE FALSE

c(pord_weakdom(3:1, 5:4), pord_weakdom(5:4, 3:1)) # 3:1 ?? 5:4
## [1] FALSE FALSE
```

We have the following result (Gagolewski, Grzegorzewski, [22]). Let  $F \in \mathcal{E}(\mathbb{I})$ . Then  $F$  is symmetric, nondecreasing in each variable and arity-monotonic if and only if for any  $\mathbf{x}, \mathbf{y}$  if  $\mathbf{x} \preceq \mathbf{y}$ , then  $F(\mathbf{x}) \leq F(\mathbf{y})$ . Therefore, the class of impact functions may be equivalently defined as all the aggregation operators that are nondecreasing with respect to this pre-order.

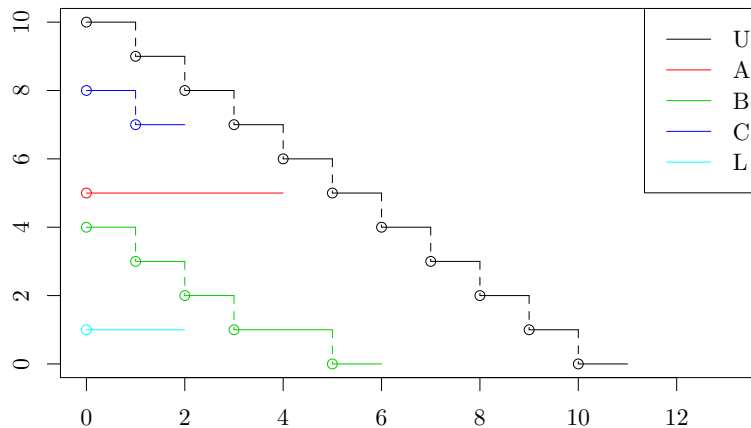
Additionally, we will write  $\mathbf{x} \triangleleft \mathbf{y}$  if  $\mathbf{x} \preceq \mathbf{y}$  and  $\mathbf{x} \neq \mathbf{y}$  (strict dominance).

**Example.** Let us consider the 5 following vectors.

```
ex1 <- list(
  U = 10:0,          # some upper bound
  A = c(5,5,5,5),    # moderate productivity & quality
  B = c(4,3,2,1,1,0), # high productivity
  C = c(8,7),         # high quality
  L = c(1,1)         # some lower bound
)
```

Plot of “citation” curves:

```
for (i in seq_along(ex1))
  plot_producer(ex1[[i]], add=(i>1), col=i)
legend("topright", legend=names(ex1), col=1:length(ex1), lty=1)
```



Here is the adjacency matrix for the pre-ordered set  $(\{A, B, C, L, U\}, \preceq)$ .

```
ord <- rel_graph(ex1, pord_weakdom)
print(ord)

## 5 x 5 sparse Matrix of class "dtCMatrix"
##   U A B C L
## U 1 . . . .
## A 1 1 . . .
## B 1 . 1 . .
## C 1 . . 1 .
## L 1 1 1 1 1

is_reflexive(ord) # is reflexive
## [1] TRUE

is_transitive(ord) # is transitive
## [1] TRUE

is_total(ord)      # not a total pre-order...
## [1] FALSE
```

We see that we have  $A \not\leq B$ ,  $A \not\leq C$ ,  $B \not\leq C$  (no pair from  $\{A, B, C\}$  is comparable w.r.t.  $\leq$ ):

```
incomp <- get_incomparable_pairs(ord)
incomp <- incomp[incomp[,1]<incomp[,2],] # remove permutations: ((1,2), (2,1)) -> (1,2)
incomp[,] <- rownames(ord)[incomp]
print(incomp) # all incomparable pairs

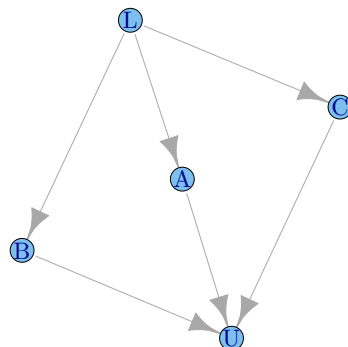
##      [,1] [,2]
## [1,] "A"  "B"
## [2,] "A"  "C"
## [3,] "B"  "C"

# the other way: generate maximal independent sets
lapply(get_independent_sets(ord), function(set) rownames(ord)[set])

## [[1]]
## [1] "A" "B" "C"
```

To draw the Hasse diagram, it will be good to de-transitivize the graph (for aesthetic reasons).

```
require(igraph)
hasse <- graph_adjacency(de_transitive(ord))
set.seed(1234567) # igraph's drawing facilities are far from perfect
plot(hasse, layout=layout_fruchterman_reingold(hasse, dim=2))
```



$(\{A, B, C, L, U\}, \leq)$  is not totally ordered, let's apply **fair totalization** (set  $x \leq'' y$  and  $y \leq'' x$  whenever  $\neg(x \leq y \text{ or } y \leq x)$ , see [16] for discussion), and then calculate transitive closure.

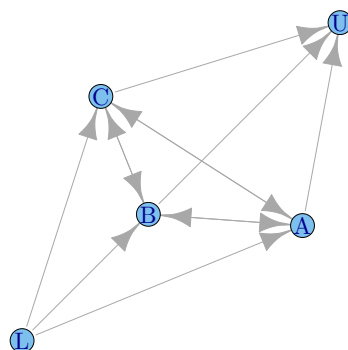
```

ord_total <- closure_transitive(closure_total_fair(ord)) # a total pre-order
print(ord_total)

## 5 x 5 sparse Matrix of class "dgCMatrix"
##   U A B C L
## U 1 . . . .
## A 1 1 1 1 .
## B 1 1 1 1 .
## C 1 1 1 1 .
## L 1 1 1 1 1

hasse <- graph_adjacency(de_transitive(ord_total))
set.seed(1234)
plot(hasse, layout=layout_fruchterman_reingold(hasse, dim=2))

```



Note that each total pre-order  $\preceq''$  induces an equivalence relation ( $x \simeq y$  iff  $x \preceq'' y$  and  $y \preceq'' x$ ; the equivalence classes may be ordered with  $\preceq''$ ). These may be explored with the `get_equivalence_classes()` function.

```
sapply(get_equivalence_classes(ord_total), function(set) rownames(ord)[set])
## [[1]]
## [1] "L"
##
## [[2]]
## [1] "A" "B" "C"
##
## [[3]]
## [1] "U"
```

Thus, we’ve obtained  $L \prec (A \simeq B \simeq C) \prec U$ .

## 5 Predefined Classes of Aggregation Operators in `agop`

### 5.1 A Review of Notable Classes of Aggregation Operators

Here are some well-known classes of aggregation operators. Originally, they were defined for fixed-length vector and for  $\mathbb{I} = [0, 1]$ .

**Definition 11.** Let  $\mathbf{w} = (w_1, \dots, w_n) \in [0, 1]^n$  be a **weighting vector** such that  $\sum_{i=1}^n w_i = 1$ . Then, for any  $\mathbf{x} \in \mathbb{I}^n$ :

1. The **weighted arithmetic mean** associated with  $\mathbf{w}$ ,  $\text{WAM}_{\mathbf{w}} : \mathbb{I}^n \rightarrow \mathbb{I}$ , is defined as

$$\text{WAM}_{\mathbf{w}}(\mathbf{x}) = \sum_{i=1}^n w_i x_i.$$

2. The **ordered weighted averaging operator** (cf. [40]) associated with  $\mathbf{w}$ ,  $\text{OWA}_{\mathbf{w}} : \mathbb{I}^n \rightarrow \mathbb{I}$ , is defined as

$$\text{OWA}_{\mathbf{w}}(\mathbf{x}) = \sum_{i=1}^n w_i x_{\{i\}}.$$

We see that both functions are idempotent, additive, and that OWA is the symmetrized version of WAM. Moreover, for  $\mathbf{w} = (n * \frac{1}{n})$ ,  $\text{WAM}_{\mathbf{w}}$  defines the arithmetic mean (`mean()` in R). **Truncated mean** is an interesting example of an OWA operator (see `mean(x, trim=...)`).

In `agop` the WAM and OWA operators are available as `wam()` and `owa()`.

```
wam(c(1,2,2,2), c(0.1,0.4,0.4,0.1))
## [1] 1.9
owa(c(1,3,5,2), rep(1,4)) # should be normalized
## Warning: elements of 'w' does not sum up to 1. correcting.
## [1] 2.75
```

**Definition 12.** Let  $\mathbf{w} = (w_1, \dots, w_n) \in \mathbb{I}^n$  be a vector such that  $\bigvee_{i=1}^n w_i = b = \sup \mathbb{I}$ . Then, for any  $\mathbf{x} \in \mathbb{I}^n$ :

1. The **weighted maximum** associated with  $\mathbf{w}$ ,  $\text{WMax}_{\mathbf{w}} : \mathbb{I}^n \rightarrow \mathbb{I}$ , is defined as

$$\text{WMax}_{\mathbf{w}}(\mathbf{x}) = \bigvee_{i=1}^n (w_i \wedge x_i).$$

2. The **ordered weighted maximum** (cf. [8]) associated with  $\mathbf{w}$ ,  $\text{OWMax}_{\mathbf{w}} : \mathbb{I}^n \rightarrow \mathbb{I}$ , is defined as

$$\text{OWMax}_{\mathbf{w}}(\mathbf{x}) = \bigvee_{i=1}^n (w_i \wedge x_{\{i\}}),$$

with  $x_{\{1\}} \geq \dots \geq x_{\{n\}}$  and  $w_1 \geq w_2 \geq \dots \geq w_n$ .

agop implementation: `wmax()` and `owmax()`.

```
wmax(c(1,3,5,2), Inf)
## Error: 'x' and 'w' should have the same length
wmax(c(1,3,5,2), rep(Inf, 4)) # greatest value /default behavior/
## [1] 5
owmax(1:10, 1:10)
## [1] 5
```

**Definition 13.** Let  $\mathbf{w} = c(w_1, \dots, w_n) \in [0, 1]^n$  be such that  $\bigwedge_{i=1}^n w_i = a = \inf \mathbb{I}$ . Then, for any  $\mathbf{x} \in \mathbb{I}^n$ :

1. The **weighted minimum**  $\text{WMin}_{\mathbf{w}} : \mathbb{I}^n \rightarrow \mathbb{I}$  associated with the weight vector  $\mathbf{w}$  is defined as

$$\text{WMin}_{\mathbf{w}}(\mathbf{x}) = \bigwedge_{i=1}^n ((1 - w_i) \vee x_i).$$

2. The **ordered weighted minimum**  $\text{OWMin}_{\mathbf{w}} : \mathbb{I}^n \rightarrow \mathbb{I}$  associated with the weight vector  $\mathbf{w}$  is defined as

$$\text{OWMin}_{\mathbf{w}}(\mathbf{x}) = \bigwedge_{i=1}^n ((1 - w_i) \vee x_{(i)}),$$

with  $x_{(1)} \leq \dots \leq x_{(n)}$  and  $w_1 \leq w_2 \leq \dots \leq w_n$ .

agop implementation: `wmin()` and `owmin()`.

It is clear to see that  $\text{OWMax}$  operators fulfill the maxitivity property and  $\text{OWMin}$  operators fulfill the minitivity property. Interestingly, it may be shown, cf. [24], that for each  $\text{OWMax}$  operator there exist an equivalent  $\text{OWMin}$  operator and inversely.

As stated above, “classical” aggregation operators were defined for vectors of fixed lengths. Let us present some notable generalizations of these operators.

Let  $\mathbb{I}^{\mathbb{I}}$  denote the set of functions from  $\mathbb{I}$  to  $\mathbb{I}$ . The following object will be needed for further considerations.

**Definition 14.** A **triangle of functions** is a sequence  $\Delta = (f_{i,n} \in \mathbb{I}^{\mathbb{I}} : i \in [n], n \in \mathbb{N})$ .

Here is a graphical interpretation of  $\Delta$ :

$$\begin{array}{cccc} f_{1,1} & & & \\ f_{1,2} & f_{2,2} & & \\ f_{1,3} & f_{2,3} & f_{3,3} & \\ \vdots & \vdots & \vdots & \ddots \end{array}$$

**Definition 15.** Let  $\Delta = (f_{i,n})_{i \in [n], n \in \mathbb{N}}$  be a triangle of functions such that  $(\forall n) \sum_{i=1}^n \inf f_{i,n} \geq a$  and  $(\forall n) \sum_{i=1}^n \sup f_{i,n} \leq b$ . Then the **quasi-L-statistic** generated by  $\Delta$  is a function  $\mathbf{qL}_\Delta : \mathbb{I}^{1,2,\dots} \rightarrow \mathbb{I}$  such that

$$\mathbf{qL}_\Delta(\mathbf{x}) = \sum_{i=1}^n f_{i,n}(x_{\{i\}}).$$

It is easily seen that quasi-L-statistics generalize OWA operators if we set  $f_{i,n}(x) = c_{i,n}x$ ,  $c_{i,n} \in [0, 1]$ , and  $(\forall n) \sum_{i=1}^n c_{i,n} = 1$ .

Assume that  $\mathbb{I} = [0, b]$ . Interestingly, it has been shown ([28], cf. also [14]) that an aggregation operator  $F : \mathbb{I}^{1,2,\dots} \rightarrow \mathbb{I}$  fulfills the symmetric modularity property if and only if  $F$  is a nondecreasing quasi-L-statistic. What is more, in [14] we may find that  $\mathbf{qL}_\Delta$  is nondecreasing if and only if there exists  $\nabla = (g_{i,n})_{i \in [n], n \in \mathbb{N}}$  such that  $(\forall n) (\forall i \in [n]) g_{i,n}$  is nondecreasing,  $(\forall n) \sum_{i=1}^n g_{i,n} \leq b$ ,  $(\forall n) (\forall i > 1) g_{i,n}(0) = 0$  and  $\mathbf{qL}_\Delta = \mathbf{qL}_{\nabla}$ .

**Definition 16.** The **quasi-S-statistic** for a given triangle of functions  $\Delta = (f_{i,n})_{i \in [n], n \in \mathbb{N}}$  is a function  $\mathbf{qS}_\Delta : \mathbb{I}^{1,2,\dots} \rightarrow \mathbb{I}$  such that

$$\mathbf{qS}_\Delta(\mathbf{x}) = \bigvee_{i=1}^n f_{i,n}(x_{\{i\}}),$$

for any  $\mathbf{x} \in \mathbb{I}^{1,2,\dots}$ .

Quasi-S-statistic generalize the OWMax operators, if  $f_{i,n}(x) = x \wedge c_{i,n}$ ,  $c_{i,n} \in \mathbb{I}$  and  $(\forall n) \bigvee_{i=1}^n c_{i,n} = b$ .

There is an equivalence between symmetric maxitive aggregation operators and nondecreasing quasi-S-statistics. Moreover, without loss of generality we may assume that a nondecreasing quasi-S-statistic is always generated by triangle of functions in which  $(\forall n) (\forall i \in [n]) f_{i,n}$  is nondecreasing,  $(\forall n) (\forall i \in [n]) f_{i,n}(a) = f_{n,n}(a)$  and  $(\forall n) f_{1,n} \preceq \dots \preceq f_{n,n}$ , see [14].

**Definition 17.** The **quasi-I-statistic** generated by  $\Delta = (f_{i,n})_{i \in [n], n \in \mathbb{N}}$  is a function  $\mathbf{qI}_\Delta : \mathbb{I}^{1,2,\dots} \rightarrow \mathbb{I}$  such that

$$\mathbf{qI}_\Delta(\mathbf{x}) = \bigwedge_{i=1}^n f_{i,n}(x_{\{i\}}),$$

for any  $\mathbf{x} \in \mathbb{I}^{1,2,\dots}$ .

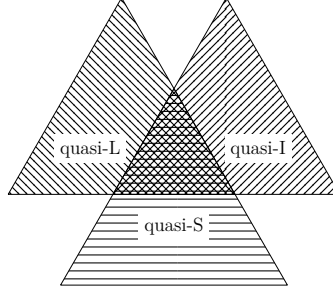
Quasi-I-statistics are generalizations of the OWMin operators, if  $f_{i,n}(x) = x \vee c_{i,n}$ ,  $c_{i,n} \in \mathbb{I}$  and  $(\forall n) \bigwedge_{i=1}^n c_{i,n} = a$ .

Like above, it has been shown that every symmetric minitive aggregation operator is a nondecreasing quasi-I-statistic, and conversely. Additionally, with no loss in generality we may assume that nondecreasing quasi-S-statistic is generated by triangle of functions in which  $(\forall n) (\forall i \in [n]) f_{i,n}$  is nondecreasing,  $(\forall n) (\forall i \in [n]) f_{i,n}(b) = f_{n,n}(b)$  and  $(\forall n) f_{1,n} \preceq \dots \preceq f_{n,n}$ , see [14].

Note: sometimes we also consider L-, S-, and I-statistics, i.e. special cases of the above-defined quasi--statistics, generated by triangles of coefficients (i.e. sequences  $\Delta = (c_{i,n} \in \mathbb{I} : i \in [n], n \in \mathbb{N})$ , cf. [4]). An **L-statistic** is a quasi-L-statistic for which we have  $f_{i,n}(x) = c_{i,n}x$ . Similarly, by setting  $f_{i,n}(x) = x \wedge c_{i,n}$  we obtain an **S-statistic** from the quasi-S-statistics class, and by setting  $f_{i,n}(x) = x \vee c_{i,n}$  we get an **I-statistic** from quasi-I-statistics.

Also note that L-statistics are known from the probability theory. However, sometimes under this name some authors understand sums of a function of order statistics.

Most interestingly, in [14] it has been shown that the intersection of any two of the three “quasi” classes is the same:



Basing on this result, the **OM3** class (symmetric maxitive, minitive, and also modular aggregation operators) was proposed in [5, 6].

**Definition 18.** A sequence of nondecreasing functions  $\mathbf{w} = (w_1, w_2, \dots)$ ,  $w_i : \mathbb{I} \rightarrow \mathbb{I}$ , and a triangle of coefficients  $\Delta = (c_{i,n})_{i \in [n], n \in \mathbb{N}}$ ,  $c_{i,n} \in \mathbb{I}$  such that  $(\forall n) \ c_{1,n} \leq c_{2,n} \leq \dots \leq c_{n,n}$ ,  $0 \leq w_n(0) \leq c_{1,n}$ , and  $w_n(b) = c_{n,n}$ , generates a nondecreasing **OM3 operator**  $M_{\Delta, \mathbf{w}} : \mathbb{I}^n \rightarrow \mathbb{I}$  such that for  $\mathbf{x} \in \mathbb{I}^n$  we have:

$$\begin{aligned} M_{\Delta, \mathbf{w}}(\mathbf{x}) &= \bigvee_{i=1}^n w_n(x_{(n-i+1)}) \wedge c_{i,n} = \bigwedge_{i=1}^n (w_n(x_{(n-i+1)}) \vee c_{i-1,n}) \wedge c_{n,n} \\ &= \sum_{i=1}^n \left( (w_n(x_{(n-i+1)}) \vee c_{i-1,n}) \wedge c_{i,n} - c_{i-1,n} \right). \end{aligned}$$

We see that the OM3 class contains i.a. all order statistics (whenever  $w_n(x) = x$ , and  $c_{i,n} = 0$ ,  $c_{j,n} = b$  for  $i < k$ ,  $j \geq k$ , and some  $k$ ), OWMax operators (for  $w_n(x) = x$ ), and the famous Hirsch  $h$ -index (see below).

## 5.2 Interesting Impact Functions

Let us go back to the Producers Assessment Problem. Below we assume that  $\mathbb{I} = [0, \infty]$ .

**The  $h$ -index.** Given a sequence  $\mathbf{x} = (x_1, \dots, x_n) \in \mathbb{I}^{1,2,\dots}$ , the **Hirsch index** [25] of  $\mathbf{x}$  is defined as  $H(\mathbf{x}) = \max\{i = 1, \dots, n : x_{\{i\}} \geq i\}$  if  $n \geq 1$  and  $x_{\{1\}} \geq 1$ , or  $H(\mathbf{x}) = 0$  otherwise. It may be shown that the  $h$ -index is a zero-insensitive OM3 aggeration operator, see [15], with:

$$H(\mathbf{x}) = \bigvee_{i=1, \dots, n} i \wedge \lfloor x_{\{i\}} \rfloor.$$

Interpretation: “an author has  $h$ -index of  $H$  if  $H$  of his/her  $n$  most cited papers have at least  $H$  citations each, and the other  $n - H$  papers are cited no more that  $H$  times each”. The  $h$ -index may also be expressed as a Sugeno integral [34] w.r.t. to a counting measure, cf. [35].

*agop* implementation: `index_h()`.

```
index_h(c(6,5,4,2,1,0,0,0,0,0,0))
```

```
## [1] 3
```

```
index_h(c(-1,3,4,2)) # only for x>=0
## Error: all elements in 'x' should be in [0,Inf]
```

Moreover, we have  $H(\mathbf{x}) \leq \min\{n, x_1\}$ .

Note that the  $h$ -index was defined in original context (aggregation of citation counts) for integer vectors. More generally, it is better to use the OM3 operator with  $w_i(x) = x = \text{Id}(x)$  and  $c_{i,n} = i$  (two identity “objects” = one of the simplest setting). Interestingly, such aggregation operator is then asymptotically idempotent, i.e. for all  $x \in \mathbb{I}$  we have  $\lim_{n \rightarrow \infty} M_{\Delta, \mathbf{w}}(n * x) = x$ .

**The  $g$ -index.** Egghe’s  $g$ -index [9] is defined as  $G(\mathbf{x}) = \max\{g = 1, \dots, n : \sum_{i=1}^g x_{\{g\}} \geq g^2\}$ , and is available in *agop* by calling `index_g()`. We have  $G(\mathbf{x}) \geq H(\mathbf{x})$  with  $G(n * n) = H(n * n) = n$ .

Note that this aggregation operator is not zero-insensitive, for example  $G(9, 0) = 2$  and  $G(9, 0, 0) = 3$ . Thus, we also provide the `index_g_zi()` function, which treats  $\mathbf{x}$  as it would be padded with 0s.

```
index_g(9)
## [1] 1
index_g(c(9,0,0))
## [1] 3
index_g_zi(9)
## [1] 3
```

The index is interesting from the computational point of view – it may be calculated on the nondecreasing vector of cumulative sums, `cumsum(sort(x, decreasing=TRUE))`, however, it cannot directly be expressed as a symmetric maxitive aggregation operator.

However, it might be shown (see [23] for the proof) that if  $\mathbf{x}$  is sorted nondecreasingly, then:

$$G(\mathbf{x}) = H(\mathbf{x})(0 \vee \text{cummin}(\text{cumsum}(x) - (1:n)^2 + (1:n))),$$

where  $1:n = (1, 2, 3, \dots, n)$ .

**The  $w$ -index.** The  $w$ -index [39] is defined as

$$W(\mathbf{x}) = \max \left\{ w = 0, 1, 2, \dots : x_{\{i\}} \geq w - i + 1, i = 1, \dots, w \right\} \quad (1)$$

and is available in *agop* by calling `index_w()`.

Interestingly, we have shown in [23] that if  $\mathbf{x}$  is sorted nondecreasingly, then:

$$W(\mathbf{x}) = H(\mathbf{x})(\text{cummin}(\mathbf{x} + (1:n) - 1)).$$

Thus, it is easily seen that this is a zero-insensitive impact function. What is more we have  $H(\mathbf{x}) \leq W(\mathbf{x}) \leq 2H(\mathbf{x})$  and  $W(\mathbf{x}) \leq \min\{n, x_1\}$ .

**The  $r_p$ -indices.** The  $r_p$ -index, for  $p \geq 1$  is expressed as

$$r_p(\mathbf{x}) = \sup \{ r > 0 : \mathbf{s}^{p,r} \leq \mathbf{x} \},$$

where  $\mathbf{s}^{p,r} = (\sqrt[p]{r^p - 0^p}, \sqrt[p]{r^p - 1^p}, \dots, \sqrt[p]{r^p - \lceil r \rceil^p})$ . For more details see [12, 19].

Please note that for integer vectors we have  $r_1 = W$  and  $r_\infty = H$  (cf. [19]). Hence it easily seen that, this is a zero-insensitive impact function.

*agop* implementation `index_rp()`.



**The  $l_p$ -indices.** The  $l_p$ -index (cf. [12, 19]) for  $p \in [1, \infty)$ ,  $u > 0$  and  $v > 0$  is a function  $l_p : \mathbb{I}^{1,2,\dots} \rightarrow \mathbb{I}^2$  given by the equation

$$l_p(\mathbf{x}) = \arg \sup_{(u,v)} \{uv : \mathbf{e}^{p,u,v} \preceq \mathbf{x}\},$$

where  $\mathbf{e}^{p,u,v} = \left( \sqrt[p]{v^p - (\frac{v}{u}0)^p}, \sqrt[p]{v^p - (\frac{v}{u}1)^p}, \dots, \sqrt[p]{v^p - (\frac{v}{u}\lfloor u \rfloor)^p} \right)$ .  
 agop implementation `index_lp()`.

**The MAXPROD-index.** The MAXPROD-index [27] is given by the equation

$$\text{MP}(\mathbf{x}) = \max \left\{ i \cdot x_{\{i\}} : i = 1, 2, \dots \right\} \quad (2)$$

is another example of zero-insensitive impact function. Interestingly, this index is a particular case of a projected  $l_\infty$ -index, see [19], and can be also expressed in terms of Shilkret integral [33], see [23] for discussion.

In *agop* the MAXPROD-index is implemented in the `index_maxprod()` function.

**Simple transformations of the  $h$ -index.** Bibliometricians in many papers considered very simple, direct modifications of the  $h$ -index. For example, the  $h(2)$ -index [26] is defined as:

$$\text{H2}(\mathbf{x}) = \max \left\{ h = 0, 1, 2, \dots : x_h \geq h^2 \right\}. \quad (3)$$

Some authors introduced other settings than “ $h^2$ ” on the right side of (3), e.g. “ $2h$ ”, “ $\alpha h$ ” for some  $\alpha > 0$ , or “ $h^\beta$ ”,  $\beta \geq 1$ , cf. [1].

It may easily be shown that these reduce to the  $h$ -index for properly transformed input vectors, e.g.  $\text{H2}(\mathbf{x}) = \text{H}(\sqrt{\mathbf{x}})$ .

## 6 Aggregation Operators from the Probabilistic Perspective

By default, theory of aggregation looks at the aggregation operators mainly from the algebraic perspective. Of course, we may also be interested in their probabilistic properties, e.g. in i.i.d. RVs models (the simplest and the most “natural” ones in statistics), cf. [12] for discussion.

Intuitively, a random variable is a method for “producing” input data. An aggregation operator applied on a random variable (possibly multidimensional) is classically called a **statistic**.

### 6.1 Some Notable Probability Distributions

Let  $(X_1, \dots, X_n)$  i.i.d.  $F$ , where  $\text{supp } F = \mathbb{I}$ . In social phenomena modeling, if  $F$  is continuous, we often assume that the underlying density  $f$  is decreasing and convex on  $\mathbb{I}$ , possibly with heavy-tails. E.g. in the bibliometric impact assessment problem, this assumption reflects the fact that higher paper valuations are more difficult to obtain than the lower ones, most of the papers have very small valuation (near 0), and the probability of attaining a high note decreases in at least linear pace.

Let us make a review of some useful statistical distributions, that are not available through “base” R (for other ones, e.g. exponential, normal, uniform, Weibull, etc. refer to the widely-available literature).

### 6.1.1 Pareto-Type II Distribution

Many generalizations of the Pareto distribution have been proposed (GPD, *Generalized Pareto Distributions*, cf. e.g. [36, 41]). Here we will introduce the so-called Pareto-Type II (Lomax) distribution, which has support  $\mathbb{I} = [0, \infty]$  and is defined with two parameters.

Formally,  $X$  follows the Pareto-II distribution with shape parameter  $k > 0$  and scale parameter  $s > 0$ , denoted  $X \sim \text{P2}(k, s)$ , if its density is of the form

$$f(x) = \frac{ks^k}{(s+x)^{k+1}} \quad (x \geq 0). \quad (4)$$

The cumulative distribution function of  $X$  is then:

$$F(x) = 1 - \frac{s^k}{(s+x)^k} \quad (x \geq 0). \quad (5)$$

The Pareto-Type II distribution is implemented in *agop*: `dpareto2()` gives the p.d.f. (4), `ppareto2()` gives the c.d.f. (5), `qpareto2()` calculates the quantile function,  $F^{-1}$ , and `rpareto2()` generates random deviates.

**Properties.** The expected value of  $X \sim \text{P2}(k, s)$  exists for  $k > 1$  and is equal to  $\mathbb{E}X = \frac{s}{k-1}$ . Variance exists for  $k > 2$  and is equal to  $\text{Var} X = \frac{ks^2}{(k-2)(k-1)^2}$ . More generally, the  $i$ -th raw moment for  $k > i$  is given by:  $\mathbb{E}X^i = \frac{\Gamma(i+1)\Gamma(k-i)}{\Gamma(k+1)} ks^i$ .

For a fixed  $s$ , if  $X \sim \text{P2}(k_x, s)$  and  $Y \sim \text{P2}(k_y, s)$ ,  $k_x < k_y$ , then  $X$  stochastically dominates  $Y$ , denoted  $X \succ Y$ . On the other hand, for a fixed  $k$ , if  $X \sim \text{P2}(k, s_x)$  and  $Y \sim \text{P2}(k, s_y)$ , then  $s_x > s_y$  implies  $X \succ Y$ .

Most importantly, if  $X \sim \text{P2}(k, s)$ , then the conditional distribution of  $X - t$  given  $X > t$ , is  $\text{P2}(k, s+t)$   $t \geq 0$ .

Additionally, it might be shown that if  $X \sim \text{P2}(k, s)$ , then  $\ln(s+X)$  has c.d.f.  $F(x) = 1 - s^k e^{-kx}$  and density  $f(x) = ks^k e^{-kx}$  for  $x \geq \ln s$ , i.e. has the same distribution as  $Z + \ln s$ , where  $Z \sim \text{Exp}(k) \equiv \Gamma(1, 1/k)$  (exponential distribution).

**Parameter estimation.** Let  $\mathbf{x} = (x_1, \dots, x_n)$  be a realization of the Pareto-Type II i.i.d. sample with known  $s > 0$ . The unbiased (corrected) maximum likelihood estimator for  $k$ :

$$\hat{k}(\mathbf{x}) = \frac{n-1}{\sum_{i=1}^n \ln\left(1 + \frac{1}{s}x_i\right)}.$$

It may be shown that for  $n > 2$  it holds  $\text{Var} \hat{k}(\mathbf{x}) = k^2 \frac{1}{n-2}$ .

*agop* implementation: `pareto2_estimate_mle()` with explicitly set argument `s`.

```
rowMeans(replicate(1000, {
  pareto2_estimate_mle(rpareto2(50, 2, 1.5), s=1.5)
}))

##          k          s
## 2.003142 1.500000
```

For both unknown  $k$  and  $s$  we have:

$$\begin{cases} \hat{k} = \frac{n}{\sum_{i=1}^n \ln(1+x_i/\hat{s})}, \\ 1 + \frac{1}{n} \sum_{i=1}^n \ln(1+x_i/\hat{s}) - \frac{n}{\sum_{i=1}^n (1+x_i/\hat{s})^{-1}} = 0. \end{cases} \quad (6)$$

Unfortunately, the second equation must be solved numerically. It is worth noting that the above system of equations may sometimes have no solution (as the local minimum of the likelihood function may not exist, see [7] for discussion). This estimator may be heavily biased and have a large mean squared error (of course, it is only asymptotically unbiased).

*agop* implementation: `pareto2_estimate_mle()` with explicitly set argument `s`.

```
rowMeans(replicate(1000, {
  pareto2_estimate_mle(rpareto2(50, 2, 1.5))
}), na.rm=TRUE)

##          k          s
## 2.794081 2.378223
```

We see that the estimator’s performance is weak.

A better (in general) estimation procedure was proposed in [42]. The Zhang-Stevens MMS (*minimum mean square error*) (Bayesian) estimator has relatively small bias (often positive) and mean squared error. In *agop* it is available as: `pareto2_estimate_mmse`.

```
rowMeans(replicate(1000, {
  pareto2_estimate_mmse(rpareto2(50, 2, 1.5))
}))

##          k          s
## 2.602839 2.089551
```

**Two-sample  $F$ -test.** The following simple test was introduced in [12]. Let  $(X_1, X_2, \dots, X_{n_1})$  *i.i.d.*  $P_2(k_1, s)$  and  $(Y_1, Y_2, \dots, Y_{n_2})$  *i.i.d.*  $P_2(k_2, s)$ , where  $s$  is an a-priori known scale parameter. We are going to verify the null hypothesis  $H_0 : k_1 = k_2$  against the two-sided alternative hypothesis  $K : k_1 \neq k_2$ .

It might be shown that  $\sum_{i=1}^n \ln(s + X_i) - n \ln s \sim \Gamma(n, 1/k)$ . This implies that under  $H_0$ , the following test statistic follows the Snedecor  $F$  distribution:

$$R(\mathbf{X}, \mathbf{Y}) = \frac{n_1 \sum_{i=1}^{n_2} \ln \left( 1 + \frac{Y_i}{s} \right)}{n_2 \sum_{i=1}^{n_1} \ln \left( 1 + \frac{X_i}{s} \right)} \stackrel{H_0}{\sim} F[2n_2, 2n_1]. \quad (7)$$

The null hypothesis is accepted iff

$$R(\mathbf{x}, \mathbf{y}) \in \left[ \mathbf{qf}\left(\frac{\alpha}{2}, 2n_2, 2n_1\right), \mathbf{qf}\left(1 - \frac{\alpha}{2}, 2n_2, 2n_1\right) \right],$$

where  $\mathbf{qf}(q, d_1, d_2)$  denotes the  $q$ -quantile of  $F^{[d_1, d_2]}$

The  $p$ -value may be determined as follows:

$$p = 2 \left( \frac{1}{2} - \left| \mathbf{pf}(R(\mathbf{x}, \mathbf{y}), 2n_2, 2n_1) - \frac{1}{2} \right| \right), \quad (8)$$

where  $\mathbf{pf}(x, d_1, d_2)$  is the c.d.f. of  $F^{[d_1, d_2]}$ .

*agop* implementation: `pareto2_test_f()`.

```
x <- rpareto2(35, 1.2, 1)
y <- rpareto2(25, 2.1, 1)
pareto2_test_f(x, y, s=1)

##
## Two-sample F-test for equality of shape parameters for Type
```

```
## II-Pareto distributions with known common scale parameter
##
## data:  x and y
## F = 0.3858, p-value = 0.000547
## alternative hypothesis: two-sided
```

## 6.2 Stochastic Properties of Aggregation Operators

Let us assume that  $(X_1, \dots, X_n)$  denotes a sequence of independent, identically distributed random variables and let c.d.f  $F$  of  $X_i$   $i = 1, \dots, n$  be continuous and strictly increasing in  $[a, b]$ , such that  $a = \inf \{x : F(x) > 0\}$  and  $b = \sup \{x : F(x) < 1\}$ ,  $a < b$ .

In this section we are going to recall some basic statistical properties of **L-** and **S-statistics**.

**Definition 19.** An **S-statistic** associated with function  $\kappa$  and a random sample  $(X_1, \dots, X_n)$  is a function

$$V_{n,k}(X_1, \dots, X_n) = \bigvee_{i=1}^n \kappa\left(\frac{i}{n}\right) \wedge X_{(n-i+1)},$$

where  $\kappa : [0, 1] \rightarrow [a, b]$ , called **control function**, is a strictly increasing continuous function such that  $\kappa(0) = a$  and  $\kappa(1) = b$ .

It is easily seen that  $V_{n,k}$  is equivalent with **S-statistic**  $S_{\Delta|_{\mathbb{I}^n}}$  (cf. Definition ??) associated with triangle of coefficients  $\Delta = (c_{i,n})_{n \in \mathbb{N}, i \in [n]}$  such that  $c_{i,n} = \kappa\left(\frac{i}{n}\right)$ .

Also please note that without loss of generality we may consider **S-statistic** form Definition 19 of a form

$$V_n(Y_1, \dots, Y_n) = \bigvee_{i=1}^n \frac{i}{n} \wedge Y_{(n-i+1)},$$

where  $(Y_1, \dots, Y_n) = (\kappa^{-1}(X_1), \dots, \kappa^{-1}(X_n))$  is a sequence of i.i.d. random variables given by the continuous c.d.f.  $G := F \circ \kappa$ . For more details see [21].

The exact distribution is given by the following Theorem.

**Theorem 20.** *The c.d.f. of  $V_n(Y_1, \dots, Y_n)$  is given by*

$$\begin{aligned} D_n(x) &= 1 - \sum_{i=1}^n \binom{n}{i} [1 - G(x)]^i [G(x)]^{n-i} \\ &= \mathcal{I}(G(x); n - \lfloor xn \rfloor, \lfloor xn \rfloor + 1), \end{aligned}$$

for  $x \in [0, 1)$ , where  $\mathcal{I}(p; a, b)$  is the regularized incomplete beta function.

Let us recall that Hirsch’s  $h$ -index, for  $x \in \mathbb{N}_0^n$  reduced itself to a **OWMax** operator, which is a special case of **S-statistic**. Therefore, the exact distribution of  $H$ , given by the following theorem, comes straightforward from the previous results (for more details see. [17]).

**Theorem 21.** *Let  $\mathbf{X} = (X_1, \dots, X_n)$  be a sequence of i.i.d. random variables with continuous c.d.f.  $F$  defined on  $\mathbb{R}_{0+}$ . Then the c.d.f. of  $H(\mathbf{X})$  for  $x \in [0, n)$  is given by*

$$D_n(x) = \mathcal{I}(F(\lfloor x + 1 \rfloor^{-0}); n - \lfloor x \rfloor, \lfloor x \rfloor + 1),$$

where  $\mathcal{I}(p; a, b)$  is the regularized incomplete beta function.

**Definition 22.** A  $\kappa$ -index of a random variable given by a c.d.f.  $F$  with respect to control function  $\kappa$  is a number  $\varrho_\kappa \in [0, 1]$  such that

$$\varrho_\kappa = 1 - F(\kappa(\varrho_\kappa)).$$

For survival function,  $S(x) = 1 - F(x)$ ,  $\varrho_\kappa$  satisfies  $\varrho_\kappa = S(\kappa(\varrho_\kappa)) = \mathbb{P}(X > \kappa(\varrho_\kappa))$ . It is clear to see that, intuitively,  $\varrho_\kappa$  is a number such that the probability of assuming a value greater than  $\kappa(\varrho_\kappa)$  is equal to  $\varrho_\kappa$ . What is more, in [21] we have following results.

**Theorem 23.**  $V_n$  is a strongly consistent estimator of  $\varrho = \varrho_{id}$ .

The asymptotic distribution of **S-statistic** is given by the following theorem.

**Theorem 24.** If  $G$  is a c.d.f. differentiable at  $\varrho$ , then

$$V_n \xrightarrow{D} \mathcal{N}\left(\varrho, \frac{1}{1 + G'(\varrho)} \sqrt{\frac{\varrho(1 - \varrho)}{n}}\right).$$

Let us now consider asymptotic distribution of **L-statistics**. Please note, that sometimes it is used to represent  $L_\Delta$  as

$$L_n = \frac{1}{n} \sum_{i=1}^n J\left(\frac{i}{n}\right) X_{(i)},$$

where function  $J(u)$  for  $u \in [0, 1]$  is such that  $J(i/n) = nc_{i,n}$ ,  $i = 1, \dots, n$ .

We present below the asymptotic normality result for normed  $L_n$ , cf ???. To this end, first define

$$\mu(J, F) = \int_0^1 J(u) F^{-1}(u) du$$

and

$$\sigma^2(J, F) = 2 \int \int_{0 < u < v < 1} J(u) J(v) u(1 - v) dF^{-1}(u) dF^{-1}(v).$$

**Theorem 25.** Assume that  $\text{var}(X_1) < \infty$  and that  $J(u)$  is bounded and continuous a.e.  $F^{-1}$ . Then

$$\lim_{n \rightarrow \infty} n\mathbb{E}(L_n) = \mu(J, F),$$

$$\lim_{n \rightarrow \infty} n\text{var}(L_n) = \sigma^2(J, F),$$

and if  $\sigma^2(J, F) > 0$ ,

$$L_n^* = \frac{L_n - \mathbb{E}(L_n)}{\sqrt{\text{var}(L_n)}} \xrightarrow{D} \mathcal{N}(0, 1).$$

What is more, it can be shown (cf. ??) that if  $\int_0^1 [u(1 - u)]^{1/2} dF^{-1}(u) < \infty$  and  $J$  is bounded and satisfies a Lipschitz condition of order  $\alpha > \frac{1}{2}$  except perhaps at finite number of continuity points of  $F^{-1}$ , then we have

$$\lim_{n \rightarrow \infty} n^{1/2} [\mathbb{E}(L_n) - \mu(J, F)] = 0.$$

## 7 NEWS/CHANGELOG

```

** agop package NEWS **

*****

0.1-0 /under development/

* initial release
  [the package started as a lightweight fork of the CITAN package]

*****

```

**Acknowledgments.** This document has been generated with L<sup>A</sup>T<sub>E</sub>X, knitr and the tikzDevice package for R. Their authors’ wonderful work is fully appreciated.

The contribution of Marek Gagolewski was partially supported by the European Union from resources of the European Social Fund, Project PO KL “Information technologies: Research and their interdisciplinary applications”, agreement UDA-POKL.04.01.01-00-051/10-00 (March-June 2013), and by FNP START Scholarship from the Foundation for Polish Science (2013).

## Bibliography

- [1] ALONSO, S., CABRERIZO, F. J., HERRERA-VIEDMA, E., AND HERRERA, F. *h-index: A review focused on its variants, computation and standardization for different scientific fields. Journal of Informetrics* 3 (2009), 273–289.
- [2] BELIAKOV, G., AND JAMES, S. Stability of weighted penalty-based aggregation functions. *Fuzzy Sets and Systems* (2013). to appear.
- [3] BELIAKOV, G., PRADERA, A., AND CALVO, T. *Aggregation Functions: A Guide for Practitioners*. New York: Springer-Verlag, 2007.
- [4] CALVO, T., AND MAYOR, G. Remarks of two types extended aggregation functions. *Tatra Mountains Mathematical Publications* 16 (1999), 235–253.
- [5] CENA, A., AND GAGOLEWSKI, M. OM3: ordered maxitive, minitive, and modular aggregation operators – Part I: Axiomatic analysis under arity-dependence. In *Aggregation Functions in Theory and in Practise (AISC 228)*, H. Bustince et al., Eds. Springer-Verlag, 2013, pp. 93–103.
- [6] CENA, A., AND GAGOLEWSKI, M. OM3: ordered maxitive, minitive, and modular aggregation operators – Part II: A simulation study. In *Aggregation Functions in Theory and in Practise (AISC 228)*, H. Bustince et al., Eds. Springer-Verlag, 2013, pp. 105–115.
- [7] DEL CASTILLO, J., AND DAOUDI, J. Estimation of the Generalized Pareto Distribution. *Statistics and Probability Letters* 79 (2009), 684–688.
- [8] DUBOIS, D., PRADÉ, H., AND TESTEMALE, C. Weighted fuzzy pattern matching. *Fuzzy Sets and Systems* 28 (1988), 313–331.
- [9] EGGHE, L. An improvement of the *H*-index: the *G*-index. *ISSI Newsletter* 2, 1 (2006), 8–9.

- [10] FRANCESCHINI, F., AND MAISANO, D. A. The Hirsch index in manufacturing and quality engineering. *Quality and Reliability Engineering International* 25 (2009), 987–995.
- [11] FRANCESCHINI, F., AND MAISANO, D. A. Structured evaluation of the scientific output of academic research groups by recent *h*-based indicators. *Journal of Informetrics* 5 (2011), 64–74.
- [12] GAGOLEWSKI, M. *Aggregation Operators and Their Application in a Formal Model for Quality Evaluation System of Scientific Research (Wybrane operatory agregacji i ich zastosowanie w modelu formalnym systemu jakości w nauce)*. PhD thesis, 2011. (In Polish).
- [13] GAGOLEWSKI, M. Bibliometric impact assessment with R and the CITAN package. *Journal of Informetrics* 5, 4 (2011), 678–692.
- [14] GAGOLEWSKI, M. On the Relationship Between Symmetric Maxitive, Minitive, and Modular Aggregation Operators. *Information Sciences* 221 (2013), 170–180.
- [15] GAGOLEWSKI, M. On the relationship between symmetric maxitive, minitive, and modular aggregation operators. *Information Sciences* 221 (2013), 170–180.
- [16] GAGOLEWSKI, M. Scientific assessment process cannot be fair?, 2013. submitted paper.
- [17] GAGOLEWSKI, M. Statistical Hypothesis Test for the Difference Between Hirsch Indices of Two Pareto-Distributed Random Samples. In *Synergies of Soft Computing and Statistics for Intelligent Data Analysis*, R. Kruse et al., Eds., vol. 190. Springer-Verlag, 2013, pp. 359–367.
- [18] GAGOLEWSKI, M., AND CENA, A. *agop: Aggregation Operators in R*, 2013. [www.rexamine.com/resources/agop/](http://www.rexamine.com/resources/agop/).
- [19] GAGOLEWSKI, M., AND GRZEGORZEWSKI, P. A geometric approach to the construction of scientific impact indices. *Scientometrics* 81, 3 (2009), 617–634.
- [20] GAGOLEWSKI, M., AND GRZEGORZEWSKI, P. Arity-monotonic extended aggregation operators. In *Information Processing and Management of Uncertainty in Knowledge-Based Systems*, E. Hüllermeier, R. Kruse, and F. Hoffmann, Eds., vol. 80. Springer-Verlag, 2010, pp. 693–702.
- [21] GAGOLEWSKI, M., AND GRZEGORZEWSKI, P. S-statistics and their basic properties. In *Combining Soft Computing and Statistical Methods in Data Analysis*, C. Borgelt et al., Eds. Springer-Verlag, 2010, pp. 281–288.
- [22] GAGOLEWSKI, M., AND GRZEGORZEWSKI, P. Possibilistic analysis of arity-monotonic aggregation operators and its relation to bibliometric impact assessment of individuals. *International Journal of Approximate Reasoning* 52, 9 (2011), 1312–1324.
- [23] GAGOLEWSKI, M., AND MESIAR, R. Monotone measures and universal integrals in a uniform framework for the scientific impact assessment problem, 2013. submitted paper.
- [24] GRABISCH, M., MARICHAL, J.-L., MESIAR, R., AND PAP, E. *Aggregation functions*. Cambridge University Press, 2009.
- [25] HIRSCH, J. E. An index to quantify individual’s scientific research output. *Proceedings of the National Academy of Sciences* 102, 46 (2005), 16569–16572.

- [26] KOSMULSKI, M. A new Hirsch-type index saves time and works equally well as the original  $h$ -index. *ISSI Newsletter* 2, 3 (2006), 4–6.
- [27] KOSMULSKI, M. MAXPROD — A new index for assessment of the scientific output of an individual, and a comparison with the  $h$ -index. *Cybermetrics* 11, 1 (2007).
- [28] MESIAR, R., AND MESIAROVÁ-ZEMÁNKOVÁ, A. The ordered modular averages. *IEEE Transactions on Fuzzy Systems* 19, 1 (2011), 42–50.
- [29] QUESADA, A. Monotonicity and the Hirsch index. *Journal of Informetrics* 3, 2 (2009), 158–160.
- [30] QUESADA, A. More axiomatics for the Hirsch index. *Scientometrics* 82 (2010), 413–418.
- [31] R DEVELOPMENT CORE TEAM. *R: A language and environment for statistical computing*. R Foundation for Statistical Computing, Vienna, Austria, 2013. URL: <http://www.R-project.org>.
- [32] ROUSSEAU, R. Woeginger’s axiomatisation of the  $h$ -index and its relation to the  $g$ -index, the  $h(2)$ -index and the  $R^2$ -index. *Journal of Informetrics* 2, 4 (2008), 335–340.
- [33] SHILKRET, N. Maxitive measure and integration. *Indag. Math.* 33 (1971), 109–116.
- [34] SUGENO, M. *Theory of fuzzy integrals and its applications*. PhD thesis, Tokyo Institute of Technology, 1974.
- [35] TORRA, V., AND NARUKAWA, Y. The  $h$ -index and the number of citations: Two fuzzy integrals. *IEEE Transactions on Fuzzy Systems* 16, 3 (2008), 795–797.
- [36] VILLASENOR-ALVA, J., AND GONZALEZ-ESTRADA, E. A bootstrap goodness of fit test for the Generalized Pareto Distribution. *Computational Statistics and Data Analysis* 53, 11 (2009), 3835–3841.
- [37] WOEGINGER, G. J. A symmetry axiom for scientific impact indices. *Journal of Informetrics* 2 (2008), 298–303.
- [38] WOEGINGER, G. J. An axiomatic analysis of Egghe’s  $g$ -index. *Journal of Informetrics* 2, 4 (2008), 364–368.
- [39] WOEGINGER, G. J. An axiomatic characterization of the Hirsch-index. *Mathematical Social Sciences* 56, 2 (2008), 224–232.
- [40] YAGER, R. On ordered weighted averaging aggregation operators in multicriteria decision making. *IEEE Transactions on Systems, Man, and Cybernetics* 18, 1 (1988), 183–190.
- [41] ZHANG, J. Improving on estimation for the Generalized Pareto Distribution. *Technometrics* 52, 3 (2010), 335–339.
- [42] ZHANG, J., AND STEPHENS, M. A. A new and efficient estimation method for the Generalized Pareto Distribution. *Technometrics* 51, 3 (2009), 316–325.