

OM3: ordered maxitive, minitive, and modular aggregation operators — Part II: A simulation study

Anna Cena and Marek Gagolewski

Abstract This article is a second part of the contribution on the analysis of the recently-proposed class of symmetric maxitive, minitive and modular aggregation operators. Recent results (Gagolewski, Mesiar, 2012) indicated some unstable behavior of the generalized h -index, which is a particular instance of OM3, in case of input data transformation. The study was performed on a small, carefully selected real-world data set. Here we conduct some experiments to examine this phenomena more extensively.

This is a revised version of the paper:

Cena A., Gagolewski M., *OM3: ordered maxitive, minitive, and modular aggregation operators — Part II: A simulation study*, In: Bustince H. et al, *Aggregation Functions in Theory and in Practise* (AISC 228), 2013, pp. 105–115, doi:10.1007/978-3-642-39165-1_14.

1 Introduction

In the first part of our contribution on OM3 aggregation operators, see [4], we carried out their axiomatic analysis under arity-dependence. Our motivation was that in many applications the “classical” assumption about fixed length of input vectors being aggregated, cf. [3, 14], is too restrictive. For example, in the Producer Assessment Problem (PAP), cf. [10], we wish to evaluate a set of producers according to their productivity and – simultaneously – the quality of the items they create.

Anna Cena

Systems Research Institute, Polish Academy of Sciences, ul. Newelska 6, 01-447 Warsaw, Poland
e-mail: Anna.Cena@ibspan.waw.pl

Marek Gagolewski

Systems Research Institute, Polish Academy of Sciences, ul. Newelska 6, 01-447 Warsaw, Poland
e-mail: gagolews@ibspan.waw.pl

Faculty of Mathematics and Information Science, Warsaw University of Technology,
ul. Koszykowa 75, 00-662 Warsaw, Poland

In Table 1 we list some typical instances of such situation, see also e.g. [6, 11]. It is easily seen that the number of artifacts varies from producer to producer. Thus, our main aim was to determine conditions required for the OM3 operators to poses some desirable properties such as zero- and F -insensitivity, or F +sensitivity.

Table 1 Typical instances of the Producer Assessment Problem (PAP).

Producer	Products	Rating method
Scientist	Scientific articles	Number of citations
Scientific institute	Scientists	The h -index
Web server	Web pages	Number of targeting web-links
R package author	R packages	Number of dependencies
Artist	Paintings	Auction price

The mentioned class of aggregation operators was of our interest, because these are the only functions which are symmetric modular, minitive, and – at the same time – maxitive, see [7]. To recall, given a closed interval of the extended real line $\mathbb{I} = [0, b]$ (possibly with $b = \infty$), the OM3 operators are defined as follows. Note that we assume that the reader is familiar with notation convention introduced in [4].

Definition 1. 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 $M_{\Delta, \mathbf{w}} \in \mathcal{P}_{(\text{nd})}$ such that for $\mathbf{x} \in \mathbb{I}^n$ we have:

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

Please note that this class includes i.a. the well-known h -index [15], all order statistics, and OWM_{ax}/OWM_{in} operators [5].

In the second part of our contribution we perform a simulation study of OM3 operators. Recently, it was noted in [11] that the generalized h -index (which is also an OM3 operator) exhibits a very unstable behavior upon some simple input elements' transformations. The study was performed on a small-sized, but carefully selected bibliometric data set. We therefore pose a question: does this undesirable behavior is also observed in a large-scale study?

The paper is organized as follows. In Sec. 2 we present some theoretical results connecting the issue of ranking of vectors using OM3 operators. The simulation results, concerning both fixed- and variable-length scenarios, are discussed in Sec. 3. Finally, Sec. 4 concludes the paper.

2 Theoretical results

We are going to analyze the correlation/association between rankings naturally created by aggregation with OM3 operators to assess their “global” change caused by vector “calibration”. This is because precise values of OM3 operators applied to variously transformed input vectors are rather meaningless. Such approach is often encountered in many domains in which aggregation operators are applied. For example, in scientometrics, we sometimes wish to order a set of authors according to the value of some citation-based quality measure, just to indicate a potential group of prominent scientists.

Keeping this in mind, let us present some theoretical results that may be useful when it comes to *comparing* OM3 operators’ values. From now on we assume that $\mathbb{I} = [0, \infty]$.

First of all, it turns out that – as far as the ranking problem is concerned – we may assume with no loss in generality that Δ is of the following, very simple form.

Proposition 1 *Let $M_{\Delta, w} \in \mathcal{P}_{(nd)} \cap \mathcal{P}_{(a0)}$ (see [4]) such that $M_{\Delta, w}(x_1, \dots, x_n) = \bigvee_{i=1}^n w(x_{(n-i+1)}) \wedge c_i$, where w is strictly increasing and $c_1 < c_2 < \dots$. Then there exist increasing functions $f, w' : \mathbb{I} \rightarrow \mathbb{I}$ for which for all $\mathbf{x} \in \mathbb{I}^{1,2,\dots}$ it holds $M_{\Delta, w}(\mathbf{x}) = f(M_{\nabla, w'}(\mathbf{x})) = f(\bigvee_{i=1}^n (w'(x_{(n-i+1)}) \wedge i))$.*

Proof. Let f be a piecewise linear continuous function such that for $i = 1, 2, \dots$ we have $f(i) = c_i$. It is obvious that f is a strictly increasing function, since the sequence $(c_i)_{i \in \mathbb{N}}$ is strictly increasing, and onto \mathbb{I} . Hence, there exists its (also strictly increasing) inverse, f^{-1} , for which we have $f^{-1}(c_i) = i$. Thus, $f^{-1}(M_{\Delta, w}(\mathbf{x})) = \bigvee_{i=1}^n (f^{-1}(w(x_{(n-i+1)})) \wedge f^{-1}(c_i)) = \bigvee_{i=1}^n ((f^{-1} \circ w)(x_{(n-i+1)}) \wedge i)$ for any $\mathbf{x} \in \mathbb{I}^{1,2,\dots}$. We may therefore set $w' = f^{-1} \circ w$, which completes the proof. \square

Moreover, please note that for $M_{cw}(\mathbf{x}) = \bigvee_{i=1}^n cw(x_{(n-i+1)}) \wedge i$, where $w : \mathbb{I} \rightarrow \mathbb{I}$ is increasing, $w(\infty) < \infty$, we may easily show that the following results hold.

Remark 1 *For any $\mathbf{x} \in \mathbb{I}^n$, $x_{(n)} < \infty$, we have $\lim_{c \rightarrow 0^+} M_{cw}(\mathbf{x}) \sim \text{MAX}(w(\mathbf{x}))$.*

Remark 2 *For any $\mathbf{x} \in \mathbb{I}^n$, $x_{(1)} > 0$, it holds $\lim_{c \rightarrow \infty} M_{cw}(\mathbf{x}) = n$.*

Therefore, we see that, intuitively, the rankings generated by some zero-insensitive OM3 operators “fall somewhere between” those generated by two very simple functions, one concerning only the producer’s ability to output artifacts of high quality, and the other reflecting solely his/her productivity.

3 Simulation study

We conducted simulation studies to assess the impact of input vector calibration on the output values of OM3 operators. We considered the following classes of functions:

- $M_c(\mathbf{x}) = \bigvee_{i=1}^n cx_{(n-i+1)} \wedge i$,
- $M_{c \log}(\mathbf{x}) = \bigvee_{i=1}^n c \log(1 + x_{(n-i+1)}) \wedge i$,
- $M_{\log c}(\mathbf{x}) = \bigvee_{i=1}^n \log(1 + cx_{(n-i+1)}) \wedge i$,

where $c \in \mathbb{R}_+$ is a scaling parameter. Note that the scaling operation is often performed on real-world data. For example, in scientometrics one may be interested in “normalizing” citations so that they reflect various characteristics of different fields (e.g. a citation in mathematics may be “worth” more than in biology), cf. [1]. The use of the logarithm is motivated by the fact that in most instances of the Producers Assessment Problem we encounter heavily-tailed and skewed data distributions.

Additionally, we considered four reference aggregation operators:

- $\text{MAX}(\mathbf{x}) = x_{(n)}$,
- generalized h -index given by $\text{HIRSCH}(\mathbf{x}) = \bigvee_{i=1}^n x_{(n-i+1)} \wedge i$, cf. [9],
- $\text{MED}(\mathbf{x})$ (sample median),
- $\Sigma \log(\mathbf{x}) = \sum_{i=1}^n (\log(1 + x_{(n-i+1)}))$.

Note that the first and the second function belongs to OM3. Moreover, as it was mentioned in [4], MAX is the only OM3 operator satisfying properties $\mathcal{P}_{(F+)}$ and $\mathcal{P}_{(F0)}$ (and $\mathcal{P}_{(a0)}$). HIRSCH, on the other hand, fulfills $\mathcal{P}_{(F0)}$ (which implies $\mathcal{P}_{(a0)}$). For $c \in \mathbb{R}_+$ such that $c \leq 1$, M_c fulfills $\mathcal{P}_{(F0)}$, but when $c > 1$, then it belongs only to $\mathcal{P}_{(a0)}$. Operators $M_{c \log}$ and $M_{\log c}$ satisfy $\mathcal{P}_{(a0)}$.

Spearman’s rank correlation coefficient. The effect of input vector calibration was evaluated by measuring the correlation between rankings created by OM3 values calculated for different scaling parameters. To assess the strength of correlation, we used Spearman’s correlation coefficient, which is a rank-based measure of association between two vectors. Technically, it is defined as the Pearson correlation coefficient between the ranks of elements. However, unlike Pearson’s r , which gives good results only when there is linear dependency, Spearman’s ρ gives sensible results when \mathbf{y} is a monotonic transformation of \mathbf{x} . What is more, since it is a non-parametric measure, it releases us from assumptions about variables’ distribution. In this paragraph we recall the definition of Spearman’s ρ and its basic properties.

Definition 2. Let $((x_1, y_1), \dots, (x_n, y_n))$ be a two-dimensional sample and let $R_i = r(x_i)$ and $S_i = r(y_i)$ denote the ranks of x_i and y_i , respectively, i.e. $x_i = x_{(R_i)}$ and $y_i = y_{(S_i)}$. Then Spearman’s rank correlation coefficient is given by

$$\rho(\mathbf{x}, \mathbf{y}) = \frac{\sum_{i=1}^n (R_i - \frac{n+1}{2})(S_i - \frac{n+1}{2})}{\sqrt{\sum_{i=1}^n (R_i - \frac{n+1}{2})^2 \sum_{i=1}^n (S_i - \frac{n+1}{2})^2}}.$$

Spearman’s ρ takes its values in $[-1, 1]$ and represents the degree of correlation between \mathbf{x} and \mathbf{y} . In particular, the closer Spearman’s ρ is either to 1 or -1 , the stronger the correlation between \mathbf{x} and \mathbf{y} is. The sign of the Spearman correlation indicates the direction of association between \mathbf{x} and \mathbf{y} . Moreover, in the context of probability, when variables are independent, the distribution of ρ not only does not depend on the joint probability distribution of (\mathbf{x}, \mathbf{y}) , but also it holds $\mathbb{E}\rho(\mathbf{x}, \mathbf{y}) = 0$.

Experimental data. Input vectors were generated from type II Pareto (Lomax) distribution family, $P2(k, s)$ (where $s > 0$ and $k > 0$), given by density function $f(x) = \frac{ks^k}{(s+x)^{k+1}}$, for $x \in \mathbb{I} = [0, \infty]$. This class of heavy-tailed, right-skewed distributions is often used in e.g. scientometrical modeling (where sometimes $k \in [1, 2]$ and $s = 1$ is assumed), see e.g. [2, 12, 13]. In this setting, a Pareto distribution describes a producer's ability to produce artifacts of various quality measures. Therefore, our knowledge of the producer's skills are given solely by k and s here.

For the sake of simplicity, we assumed that $s = 1$. The shape parameter k was randomly generated for every vector from the uniform distribution on interval $(1, 2)$, i.e. $k \sim U(1, 2)$, or the $P2(1, 1)$ distribution shifted by one, i.e. $k \sim P2(1, 1) + 1$. These model a population of producers of different abilities. Additionally, we considered the cases of producers of equal skills, with k equal to 1, 1.5, or 2. The calibration parameters c were taken from $[0.001, 10000]$.

We considered three simulation scenarios. In the first one, input vectors' length n was the same for all vectors. In the second one, we will examine the correlation for vectors of equal lengths and their expanded versions (cf. the arity-monotonicity property from [4]). In the last scenario, for each vector we generated their lengths randomly. In each step, $MC = 100000$ Monte Carlo samples were generated. The computations were performed with the `agop` package [8] for R.

3.1 Vectors of fixed lengths

First we analyzed fixed input vectors' lengths which were set to $n = 25, 100, 250$, and 1000 elements. Please note that this may be interpreted as an evaluation of producers of the same productivity.

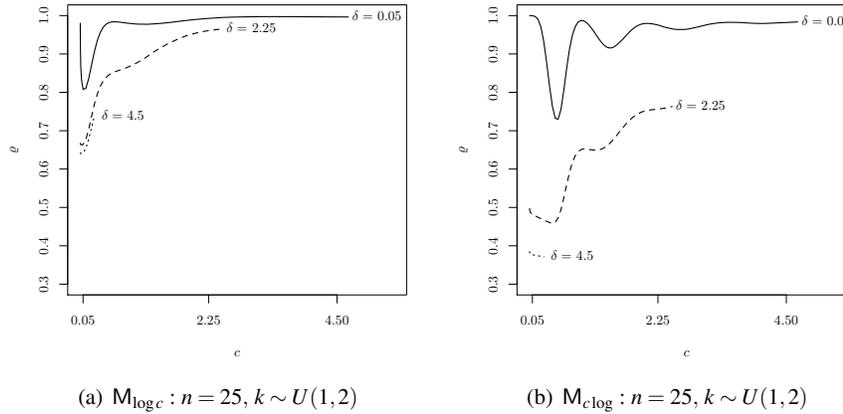


Fig. 1 The effect of adding small values to the calibration parameter.

Let us examine the sensitivity of OM3-generated rankings to vectors calibration. We calculated Spearman's rank correlation for $(M_c(\mathbf{x}_1), M_{c+\delta}(\mathbf{x}_1)), \dots, (M_c(\mathbf{x}_{MC}), M_{c+\delta}(\mathbf{x}_{MC}))$, and the same for the other operators. Two plots in Fig. 1 depict some exemplary, but representative results concerning, respectively, the functions $M_{\log c}$ and $M_{c \log}$ for $n = 25, k \sim U(1, 2)$.

We note that for small δ , the value of $\rho(M_c, M_{c+\delta})$ is relatively high (≥ 0.9 for $n = 25$). However, in most of the analyzed cases we observe a decrease in correlation strength for $c \simeq 0.4$, which may indicate some sort of ranking instability. Therefore, as far as applications are concerned, the scaling parameters should be chosen with care.

Let us now consider the correlation between $M_c, M_{\log c}$, and $M_{c \log}$, and the reference rankings, i.e. those generated by MAX, MED, HIRSCH, and $\Sigma \log$. Two exemplary cases are depicted in Fig. 2. Please note the log scale for c on the x axis.

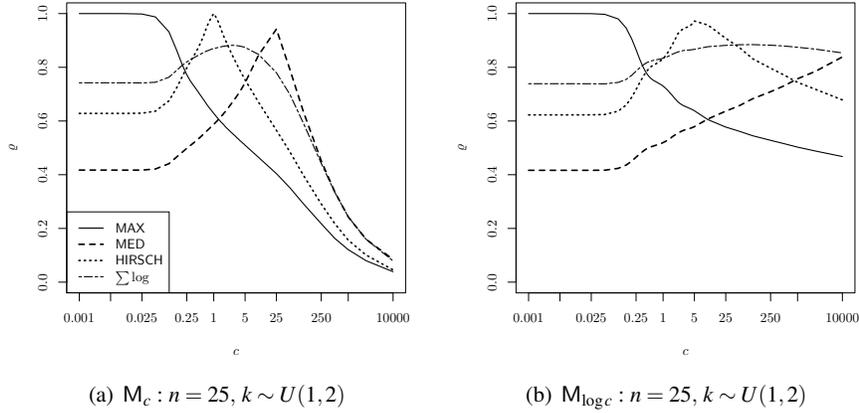


Fig. 2 Spearman's rank correlation coefficient between OM3 and the reference rankings.

Obviously, in each case for small c we get the same ranking as for the MAX function (see Remark 1). On the other hand, as c approaches ∞ , the OM3 rankings are uncorrelated with the reference ones (see Remark 2, e.g. MAX and n are independent random variables). Also note that $M_1 = \text{HIRSCH}$.

In all the analyzed cases we observed quite similar behavior of the four functions. Interestingly, with each of the three OM3 classes we can obtain, with a good accuracy, the reference, MAX-, HIRSCH-, MED-, and $\Sigma \log$ -based rankings. This may indicate, of course as far as the Paretian model and fixed n is concerned, that the OM3 aggregation operators may be sufficiently comprehensive in some applications. What is more, we observed that for $k \sim U(1, 2)$ or $k \sim P2(1, 1) + 1$ the correlations are higher than for fixed k . Likewise, when vectors' lengths increase, the correlations also increase. Let us now investigate the influence of shape parameter's and vectors' lengths n selection deeper.

How does k affect the rankings? As we can see in Fig. 3, the correlation between OM3- and HIRSCH-based rankings is greater in a case of $k \sim U(1,2)$ and $k \sim P2(1,1) + 1$, i.e. when k was generated randomly for each vector (producers of diverse characteristics), than in case of fixed k (producers of uniform abilities). What is more, the results obtained for $k = 1, 1.5$, and 2 are quite similar.

We see that in the HIRSCH case we observe that small change in the calibration parameter in the neighborhood of 1 causes noticeable decrease of the degree of correlation – cf. also [11].

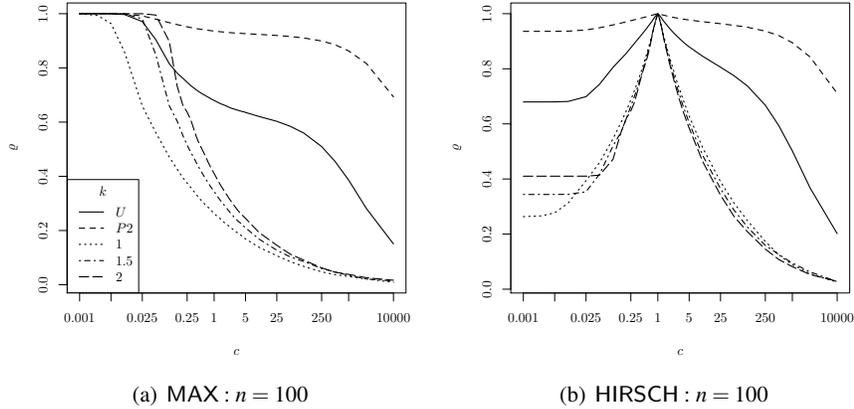


Fig. 3 Spearman's rank correlation coefficient between OM3- and, respectively, MAX- and HIRSCH-based rankings for different k generation methods.

How does n affect the rankings? In Fig. 4 we depict the case of producers of different productivity. For HIRSCH we observe that for randomly generated k , the bigger n is, the larger correlations we get. However, for fixed k the behavior is more complicated. For small c we notice larger ρ for smaller n .

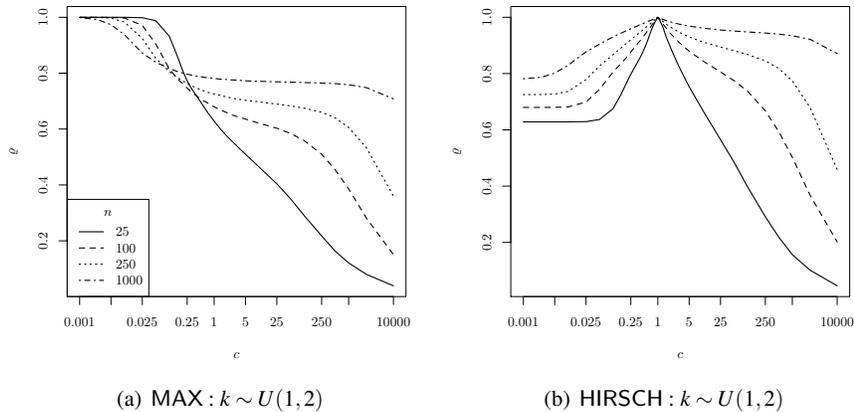


Fig. 4 Spearman's ρ between OM3- and, respectively, MAX- and HIRSCH-based rankings for different n .

3.2 Vector expansion

In the next scenario we represent the case in which we have a set of producers with $n_0 = 25$ artifacts. They are assessed with different OM3 operators. Then, to each vector describing the producer, we add new elements. Of course, according to the arity-monotonicity property, their valuation does not decrease (cf. [4]). The number of added elements, Δn , was independently generated for each producer from the heavy-tailed $[P2(1,1) + 1]$ distribution and the shifted Poisson distribution $Pois(5) + 1$ ($\text{Var}\Delta n = 5$). Moreover, $\Delta n = 25$ was also considered.

In Fig. 5 we presented a typical output. First of all, there is no substantial influence of the Δn distribution in the analyzed cases. Here, of course, as $c \rightarrow \infty$, $\rho \rightarrow 0$ (n and Δn are independent). For small and moderate values of c the correlation between original and extended vectors' valuations are high, but yet not perfect. Thus, the productivity of a producer indeed affects also his/her valuation with OM3 operators.

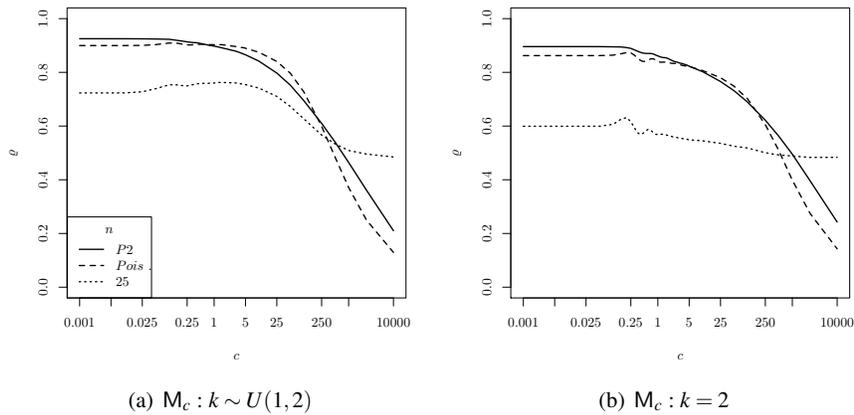


Fig. 5 Spearman's rank correlation coefficient between OM3-ranked original vectors and their expanded versions.

For $\Delta n = 25$ the correlation is lower, but much more insensitive to the value of the calibration coefficient. Note that the default ranking method for tied (equal) observations in R's `cor()` function uses averaging, therefore for $c \rightarrow \infty$ we get $\rho \rightarrow 0.5$ for fixed Δn .

3.3 Vectors of random lengths

In the last scenario let us examine a set of producers of random productivity. We considered $n \sim [P2(1,1) + 1]$, $n \sim Pois(5) + 1$, and $n \sim [U[1,500]]$.

Fig. 6ab depicts the correlation between OM3 and reference-based rankings in the first two cases. Note that the density functions of these distributions are decreasing. Therefore, there is a relatively large probability of obtaining small values of n : for $n \sim [P2(1, 1) + 1]$ we have $\text{Med}n = 1$ and for $n \sim \text{Pois}(5) + 1$ we get $\text{Med}n = 6$.

Fig. 6c depicts the $n \sim [U[1, 500]]$ and the fixed k case. It may be shown that if (X_1, \dots, X_n) i.i.d. $P2(k, 1)$ then $Y_n := \sum_{i=1}^n \log(X_i + 1) \sim \Gamma(n, 1/k)$ and $\mathbb{E}Y_n = n/k$. This explains a high degree of correlation between $\Sigma \log$ and M_c for $c \rightarrow \infty$.

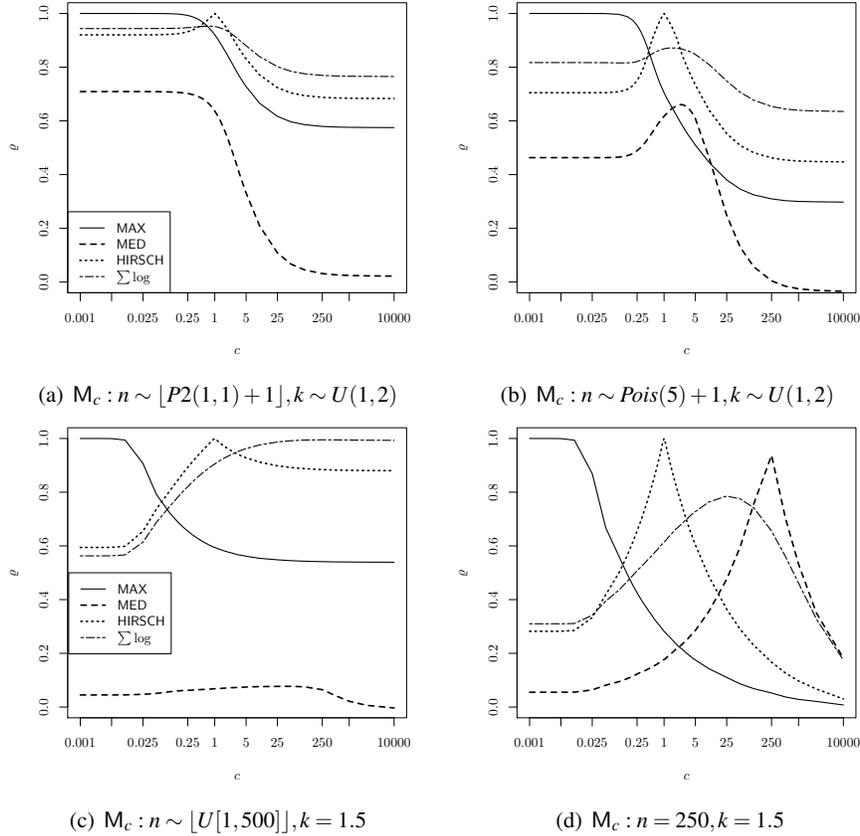


Fig. 6 Spearman's rank correlation coefficient between OM3 and references operators.

While describing the results instantiated in Fig. 2 (fixed n) we noted that OM3 class is quite flexible in terms of approximating the two reference aggregation operators, not mentioned in Remark 1 and 2. From Fig. 6abc we may deduce that for variable n such a nice property does not hold. Moreover, by comparing Fig. 6c and Fig. 6d we may observe that the influence of varying n is significant.

Additionally, in Fig. 7 (cf. Fig. 4) we observe that the method for generating n is of substantial influence.

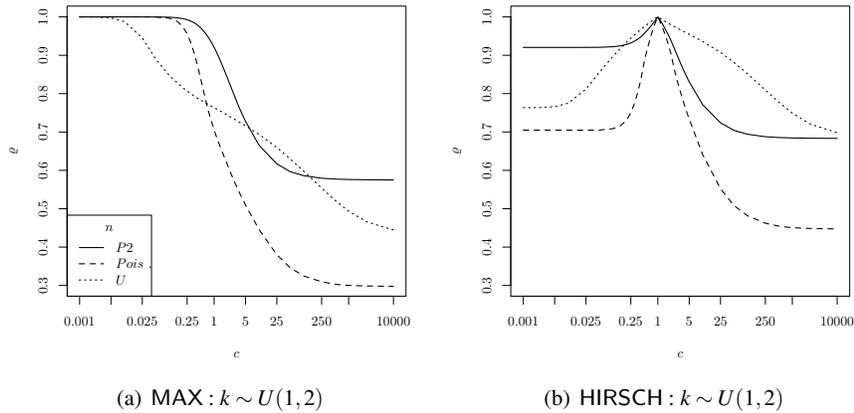


Fig. 7 Spearman's ρ between OM3-operators and, respectively, MAX- and HIRSCH-based rankings for different n generation methods.

4 Conclusions

The main aim of our simulation study was to assess the behavior of some OM3 operators under various transformations of the input data. We focused our investigation on input vectors calibration, since we have shown that, as far as the ranking problem is concerned, the form of the coefficients' triangle may be fixed.

To evaluate the impact of data scaling we examined the correlations between OM3-based rankings for different scaling parameters. Moreover, we paid special attention to some popular operators such as the generalized h -index, sample maximum and sample median. In our study we considered input vectors of fixed and random lengths. Moreover, we examined the correlation for vectors of equal lengths and their expanded versions.

First of all, we noted that a choice of the scaling parameter has a significant impact on OM3 operators. Hence, in practical applications we should be very careful while selecting an appropriate aggregation method. The issue of automated generation of w definitely should be investigated much more deeply. Thus, we leave this for our future research.

What is more, we observed high sensitivity of the operators to the formulation of the model describing real-world phenomena being considered.

Acknowledgments. The contribution of Marek Gagolewski was partially supported by FNP START Scholarship from the Foundation for Polish Science.

Please cite this paper as:

Cena A., Gagolewski M., *OM3: ordered maxitive, minitive, and modular aggregation operators — Part II: A simulation study*, In: Bustince H. et al, *Aggregation Functions in Theory and in Practise* (AISC 228), 2013, pp. 105–115, doi:10.1007/978-3-642-39165-1_14.

References

- [1] Alonso S, Cabrerizo FJ, Herrera-Viedma E, Herrera F (2009) *h*-index: A review focused on its variants, computation and standardization for different scientific fields. *Journal of Informetrics* 3:273–289
- [2] Barcza K, Telcs A (2009) Paretian publication patterns imply Paretian Hirsch index. *Scientometrics* 81(2):513–519
- [3] Beliakov G, Pradera A, Calvo T (2007) *Aggregation Functions: A Guide for Practitioners*. Springer-Verlag
- [4] Cena A, Gagolewski M (2013) OM3: ordered maxitive, minitive, and modular aggregation operators – Part I: Axiomatic analysis under arity-dependence. In: Bustince H et al (eds), *Aggregation Functions in Theory and in Practise (AISC 228)*, Springer-Verlag, pp 93–103.
- [5] Dubois D, Prade H, Testemale C (1988) Weighted fuzzy pattern matching. *Fuzzy Sets and Systems* 28:313–331
- [6] Franceschini F, Maisano DA (2009) The Hirsch index in manufacturing and quality engineering. *Quality and Reliability Engineering International* 25:987–995
- [7] Gagolewski M (2013) On the relationship between symmetric maxitive, minitive, and modular aggregation operators. *Information Sciences* 221:170–180
- [8] Gagolewski M, Cena A (2013) *agop*: Aggregation Operators Package for R. <http://www.ibspan.waw.pl/~gagolews/agop/>
- [9] Gagolewski M, Grzegorzewski P (2010) Arity-monotonic extended aggregation operators. In: Hüllermeier E, Kruse R, Hoffmann F (eds) *Information Processing and Management of Uncertainty in Knowledge-Based Systems*, vol 80, Springer-Verlag, pp 693–702
- [10] Gagolewski M, Grzegorzewski P (2011) Possibilistic analysis of arity-monotonic aggregation operators and its relation to bibliometric impact assessment of individuals. *International Journal of Approximate Reasoning* 52(9):1312–1324
- [11] Gagolewski M, Mesiar R (2012) Aggregating different paper quality measures with a generalized *h*-index. *Journal of Informetrics* 6(4):566–579
- [12] Glänzel W (2008) *H*-index concatenation. *Scientometrics* 77(2):369–372
- [13] Glänzel W (2008) On some new bibliometric applications of statistics related to the *h*-index. *Scientometrics* 77(1):187–196
- [14] Grabisch M, Marichal JL, Mesiar R, Pap E (2009) *Aggregation functions*. Cambridge
- [15] Hirsch JE (2005) An index to quantify individual’s scientific research output. *Proceedings of the National Academy of Sciences* 102(46):16,569–16,572