



Universidad
Carlos III de Madrid



This is a supplementary material of a postprint version of the following published document:

Perianes-Rodriguez, A.; Ruiz-Castillo, J. (2017). A comparison of the Web of Science and publication-level classification systems of science. [Journal of Infonometrics](#), v. 11, n. 1, pp. 32-45. Available in <https://doi.org/10.1016/j.joi.2016.10.007>.

© Elsevier



This work is licensed under a Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International License.

SUPPLEMENTARY MATERIAL

A. PERCENTAGE DISTRIBUTION OF ARTICLES IN THE WOS SYSTEM BY THE NUMBER OF SUBJECT CATEGORIES TO WHICH THEY ARE ASSIGNED

Number of categories	Distinct articles	%	Articles in the extended count	%
1	1,982,353	54.8	1,982,353	33.3
2	1,084,414	30.0	2,168,828	36.5
3	429,765	11.9	1,289,295	21.7
4	94,135	2.6	376,540	6.3
5	15,151	0.4	75,755	1.3
6	8,627	0.2	51,762	0.9

B. A NUMERICAL EXAMPLE OF THE CITATION DISTRIBUTIONS DISCUSSED IN THE TEXT

1. The G8 classification system

Assume that there are eleven distinct articles classified in four clusters in the G8 system, say $\mathbf{c}_j, j = 1, 2, 3, 4$:

$$\begin{aligned}
 \mathbf{c}_1 &= (0, 3, 12) && \text{with } \mu_1 = 15/3 = 5, \\
 \mathbf{c}_2 &= (1, 2, 6) && \text{with } \mu_2 = 9/3 = 3, \\
 \mathbf{c}_3 &= (4, 8) && \text{with } \mu_3 = 12/2 = 6, \\
 \mathbf{c}_4 &= (5, 7, 9) && \text{with } \mu_4 = 21/3 = 7.
 \end{aligned}$$

The overall citation distribution under the G8 system, $\mathbf{C} = \cup_j \{\mathbf{c}_j\}$, is the following:

$$\mathbf{C} = (0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 12).$$

The normalized cluster distributions are:

$$\begin{aligned}
 \mathbf{c}^*_1 &= (0/\mu_1, 3/\mu_1, 12/\mu_1) = (0, 3/5, 12/5), \\
 \mathbf{c}^*_2 &= (1/\mu_2, 2/\mu_2, 6/\mu_2) = (1/3, 2/3, 2), \\
 \mathbf{c}^*_3 &= (4/\mu_3, 8/\mu_3) = (2/3, 4/3), \\
 \mathbf{c}^*_4 &= (5/\mu_4, 7/\mu_4, 9/\mu_4) = (5/7, 1, 9/7).
 \end{aligned}$$

The overall G8-normalized citation distribution under the G8 system, $\mathbf{C}^* = \cup_j \{\mathbf{c}^*_j\}$, is the following:

$$\mathbf{C}^* = (0, 1/3, 3/5, 2/3, 2/3, 5/7, 1, 4/3, 9/7, 2, 12/5)$$

with $|\mathbf{C}^*| = |\mathbf{C}| = 11$.

2. The G6 classification system

Assume that, in addition to the previous eleven distinct articles, there are two more articles with 0 and 11 citations in system G6. Assume also that the thirteen distinct articles are classified in two clusters, say $\mathbf{d}_g, g = 1, 2$:

$$\begin{aligned}
 \mathbf{d}_1 &= (0, 1, 2, 7, 9, 11) && \text{with } M_1 = 30/6 = 5, \\
 \mathbf{d}_2 &= (0, 3, 4, 5, 6, 8, 12) && \text{with } M_2 = 38/7.
 \end{aligned}$$

The overall citation distribution under the G6 system, $\mathbf{D} = \cup_g \{\mathbf{d}_g\}$, is the following:

$$\mathbf{D} = (0, 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 11, 12).$$

The normalized cluster distributions are:

$$\begin{aligned}\mathbf{d}^*_1 &= (0, 1/5, 2/5, 7/5, 9/5, 11/5), \\ \mathbf{d}^*_2 &= (0, 21/38, 14/19, 35/38, 21/19, 28/19, 42/19).\end{aligned}$$

The overall G6-normalized citation distribution under the G6 system, $\mathbf{D}^* = \cup_g \{\mathbf{d}^*_g\}$, is the following:

$$\mathbf{D}^* = (0, 0, 1/5, 2/5, 21/38, 14/19, 35/38, 21/19, 7/5, 28/19, 9/5, 11/5, 42/19)$$

with $|\mathbf{D}^*| = |\mathbf{D}| = 13$.

3. The construction of the citation distribution \mathbf{C}^{**}

In the construction of the citation distribution \mathbf{C}^{**} to assess the G6-normalization procedure using the G8 classification system for evaluation purposes, the four cluster citation distributions \mathbf{c}^{**}_j are as follows:

$$\begin{aligned}\mathbf{c}^{**}_1 &= (0, 21/38, 42/19), \\ \mathbf{c}^{**}_2 &= (1/5, 2/5, 21/19), \\ \mathbf{c}^{**}_3 &= (14/19, 28/19), \\ \mathbf{c}^{**}_4 &= (35/38, 7/5, 9/5).\end{aligned}$$

The overall G6-normalized citation distribution under the G8 system, $\mathbf{C}^{**} = \cup_j \{\mathbf{c}^{**}_j\}$, is the following:

$$\mathbf{C}^{**} = (0, 1/5, 2/5, 21/38, 14/19, 35/38, 21/19, 7/5, 28/19, 9/5, 42/19).$$

Since there are two distinct articles in citation distribution \mathbf{D}^* that do not belong to \mathbf{C}^{**} , we have $|\mathbf{C}^{**}| = 11 < |\mathbf{D}^*| = 13$.

4. The construction of the citation distribution \mathbf{D}^{**}

To assess the performance of the G8-normalization procedure using the G6 system for evaluation purposes, consider again the intersection $\mathbf{C} \cap \mathbf{D}$. For every article v in this set with citations d_{gv} in cluster g in system G6, there exists a distinct article l in some cluster j with the same number of citations, i.e. $c_{jl} = d_{gv}$. The normalized score of this article according to the G6 system becomes $d^{**}_{gv} = c_{jl}/\mu_j$, where μ_j is the mean citation in cluster j . In this way, we construct an overall G8-normalized citation distribution under the G6 system $\mathbf{D}^{**} = \cup_g \{\mathbf{d}^{**}_g\}$. Since $c_{jl}/\mu_j = c^*_{jl}$, so that $d^{**}_{gv} = c^*_{jl}$, every article in \mathbf{D}^{**} belongs to \mathbf{C}^* . If there are some articles in the G8 system that belong to a small cluster with less than 250 articles in the G6 system, then the citation distribution \mathbf{D}^{**} is strictly contained in \mathbf{C}^* . In turn, since there are some articles in the G6 system that belong to a cluster with less than 250 articles in the G8 system, we have $\mathbf{D}^{**} \subset \mathbf{C}^* \subset \mathbf{D}^*$, so that $|\mathbf{D}^{**}| < |\mathbf{C}^*| = 3.4$ million articles $< |\mathbf{D}^*| = 3.6$ million articles.

In the example, the two cluster citation distributions \mathbf{d}^{**}_g are as follows:

$$\begin{aligned}\mathbf{d}^{**}_1 &= (0, 1/3, 2/3, 1, 9/7), \\ \mathbf{d}^{**}_2 &= (3/5, 2/3, 5/7, 4/3, 2, 12/5).\end{aligned}$$

Since for each article in \mathbf{D} there exists a distinct article in \mathbf{C} , the overall G8-normalized citation distribution under the G6 system, $\mathbf{D}^{**} = \cup_g \{\mathbf{d}^{**}_g\}$, coincides with \mathbf{C}^* , so that

$$\mathbf{D}^{**} = (0, 1/3, 3/5, 2/3, 2/3, 5/7, 1, 4/3, 9/7, 2, 12/5)$$

with $|D^{**}| = |C^*| = 11 < |D^*| = 13$.

5. The extended count in the example

Assume that there are three WoS categories, and that three distinct articles in C and D are assigned to several categories: the two articles with five and seven citations are assigned to the first and the third categories, while the article with three citations is assigned to the three categories. We assume that the seventeen articles in the extended count are distributed as follows:

$$\begin{aligned} e_1 &= (1, 3, 5, 7, 11) & \text{with } \mu_1 &= 27/5, \\ e_2 &= (0, 0, 3, 4, 6, 8) & \text{with } \mu_2 &= 21/6, \\ e_3 &= (2, 3, 5, 7, 9, 12) & \text{with } \mu_3 &= 38/6 = 19/3. \end{aligned}$$

The overall extended citation distribution under the WoS system, $E = \cup_k \{e_k\}$, is the following:

$$E = (0, 0, 1, 2, 3, 3, 3, 4, 5, 5, 6, 7, 7, 8, 9, 11, 12)$$

The normalized extended category citation distributions are:

$$\begin{aligned} e^*_1 &= (5/27, 15/27, 25/27, 35/27, 55/27), \\ e^*_2 &= (0, 0, 18/21, 24/21, 36/21, 48/21), \\ e^*_3 &= (6/19, 9/19, 15/19, 21/19, 27/19, 36/19). \end{aligned}$$

The overall WoS-normalized citation distribution under the WoS system, $E^* = \cup_k \{e^*_k\}$, is the following:

$$E^* = (0, 0, 9/19, 5/27, 15/27, 6/19, 18/21, 25/27, 21/19, 24/21, 35/27, 27/19, 15/9, 36/21, 36/19, 55/27, 48/21)$$

with $|E^*| = |E| = 17$.

6. The construction of the E^{**} citation distribution

In order to assess the performance of the G8-normalization procedure when the WoS system is used for evaluation purposes, consider the intersection between C and the set of distinct articles in E . Consider an article u in this set with citations c_{ju} in cluster j in system G8. Assume, for example, that there are two articles v and w in categories l and h in the WoS system with this same number of citations, i.e. $c_{lv} = c_{hw} = c_{ju}$. In this case, the distinct article with raw citations c_{ju} will give rise to two normalized articles according to the WoS system, say e^{**}_{lv} and e^{**}_{hw} . Since the mean citation in cluster j is μ_j , the two normalized articles become $e^{**}_{lv} = e^{**}_{hw} = c_{ju}/\mu_j$. These articles will be included in the extended normalized citation distributions e^{**}_j and e^{**}_h . In this way, we construct a G8-normalized overall citation distribution under the WoS system $E^{**} = \cup_k \{e^{**}_k\}$. Note, however, that $c_{ju}/\mu_j = c^*_{ju}$, so that $e^{**}_{lv} = e^*_{lv}$, and $e^{**}_{hw} = e^*_{hw}$. Therefore, the citation distribution C^* is strictly contained in E^{**} . On the other hand, since there are some articles in the WoS system that belong to a cluster with less than 250 articles in the G8 system, E^{**} is strictly contained in E^* . Therefore, we have $|C^*| = 3.4$ million articles $< |E^{**}| < |E^*| = 5.9$ million articles.

In the example, the three category citation distributions e^{**}_k are as follows:

$$\begin{aligned} e^{**}_1 &= (1/3, 3/5, 5/7, 1), \\ e^{**}_2 &= (0, 3/5, 2/3, 4/3, 2), \\ e^{**}_3 &= (3/5, 2/3, 5/7, 1, 9/7, 12/5). \end{aligned}$$

The extended overall G8-normalized citation distribution under the WoS system, $E^{**} \cup_k \{e^{**}_k\}$, becomes

$$E^{**} = (0, 1/3, 3/5, 3/5, 3/5, 2/3, 2/3, 5/7, 5/7, 1, 1, 4/3, 9/7, 2, 12/5).$$

Since there are two distinct articles in citation distribution E^* that do not belong to E^{**} , we have $|E^{**}| = 15 < |E^*| = 17$.

7. The construction of the E^{***} citation distribution

In order to assess the performance of the G6-normalization procedure when the WoS system is used for evaluation purposes, we proceed in a similar way. In other words, we construct a G6-normalized overall citation distribution under the WoS system $E^{***} = \cup_k \{e^{***}_k\}$. In this case, we have $|D^*| = 3.6$ million articles $< |E^{***}| < |E^*| = 5.9$ million articles.

In the example, the three category citation distributions e^{***}_k are as follows:

$$\begin{aligned} e^{***}_1 &= (1/5, 21/38, 35/38, 7/5, 11/5), \\ e^{***}_2 &= (0, 0, 21/38, 28/38, 21/19, 28/19), \\ e^{***}_3 &= (2/19, 21/38, 35/38, 7/5, 9/5, 42/19). \end{aligned}$$

The extended overall G6-normalized citation distribution under the WoS system, $E^{***} \cup_k \{e^{***}_k\}$, becomes $E^{***} = (0, 0, 1/5, 2/5, 21/38, 21/38, 21/38, 14/19, 35/38, 35/38, 21/19, 7/5, 7/5, 28/19, 9/5, 11/5, 42/19)$.

Since for any distinct article in citation distribution E^* there exists a distinct article in D^* , we have $|E^{***}| = |D^*| = |E^*| = 17$.

8. The construction of the C^{***} citation distribution

In order to assess the performance of the WoS-normalization procedure when the G8 system is used for evaluation purposes, consider the intersection between C and the set of distinct articles in E . Note that for every distinct article u in this set receiving c_{ju} citations in the cluster citation distribution c_j , there must be one or more articles in the original WoS system with the same number of citations. Assume, for example, that there are two articles v and w in categories k and h , respectively, with $e_{kv} = e_{hw} = c_{ju}$. Since the mean citations in categories k and h are μ_k and μ_h , the two normalized articles according to the G8 system become $c^{***}_{jv} = c_{ju}/\mu_k$ and $c^{***}_{jw} = c_{ju}/\mu_h$. Both of these articles will be included in the extended normalized citation distribution c^{***}_j . In this way, we construct a WoS-normalized overall citation distribution under the G8 system $C^{***} = \cup_j \{c^{***}_j\}$. Note, however, that $c_{ju}/\mu_k = e_{kv}/\mu_k = e^*_{kv}$ and $c_{ju}/\mu_h = e_{hw}/\mu_h = e^*_{hw}$, so that $c^{***}_{jv} = e^*_{kv}$, and $c^{***}_{jw} = e^*_{hw}$. Therefore, every element in C^{***} belongs to the overall extended citation distribution E^* . Since there are some articles in the WoS system whose corresponding distinct articles belong to a small cluster with less than 250 articles, the citation distribution C^{***} is strictly contained in E^* , that is, $C^{***} \subset E^*$. On the other hand, since for any distinct article in C there might be two or more articles in the WoS the citation distribution C^{***} strictly contains C^* , so that $|C| = 3.4$ million articles $< |C^{***}| < |E^*| = 5.9$ million articles.

In the example, the four cluster citation distributions c^{***}_j are as follows:

$$\begin{aligned} c^{***}_1 &= (0, 9/19, 15/27, 18/21, 36/19), \\ c^{***}_2 &= (5/27, 6/19, 36/21), \\ c^{***}_3 &= (24/21, 16/7), \\ c^{***}_4 &= (25/27, 15/19, 21/19, 35/27, 27/19). \end{aligned}$$

The overall WoS-normalized citation distribution under the G8 system, $C^{***} = \cup_j \{c^{***}_j\}$, is the following:

$$C^{***} = (0, 9/19, 5/27, 15/27, 6/19, 15/19, 18/21, 25/27, 21/19, 24/21, 35/27, 27/19, 12/7, 36/19, 16/7).$$

Since there are two distinct articles in citation distribution E^* that do not belong to C^{***} , we have $|C^{***}| = 15 < |E^*| = 17$.

9. The construction of the D^{***} citation distribution

In order to assess the performance of the WoS-normalization procedure when the G6 system is used for evaluation purposes, consider the intersection between D and the set of distinct articles in E . For every distinct article u in this set receiving d_{gu} citations in the cluster citation distribution d_g in the G6 system, there must be one or more articles in the WoS system with the same number of citations. Assume, for example, that there are two articles v and w in categories k and b with $e_{kv} = e_{bw} = d_{gu}$. In this case, the distinct article with raw citations d_{gu} will give rise to two normalized articles according to the WoS system, say d^{***}_{gv} and d^{***}_{gw} . Since the mean citations in categories k and b are μ_k and μ_b , the two normalized articles become $d^{***}_{gv} = e_{kv}/\mu_k = e^*_{kv}$ and $d^{***}_{gw} = e_{bw}/\mu_b = e^*_{bw}$. Both of these articles will be included in the extended normalized citation distribution d^{***}_g . In this way, we construct a WoS-normalized overall citation distribution under the G6 system $D^{***} = \cup_g \{d^{***}_g\}$. The citation distribution D^{***} is strictly contained in E^* , and strictly contains D^* , so that $|D| = 3.6$ million articles $< |D^{***}| < |E^*| = 5.9$ million articles.

In the example, the two cluster citation distributions d^{***}_g are as follows:

$$\begin{aligned} d^{***}_1 &= (0, 5/27, 6/19, 35/27, 21/19, 27/19, 55/27), \\ d^{***}_2 &= (0, 15/27, 18/21, 9/19, 24/21, 25/27, 15/19, 36/21, 48/21, 36/19). \end{aligned}$$

Since for any article in citation distribution E there exists a distinct article in citation distribution D , the overall WoS-normalized citation distribution under the G6 system, $D^{***} = \cup_j \{d^{***}_j\}$, coincides with E^* , so that

$$D^{***} = (0, 0, 9/19, 5/27, 15/27, 6/19, 18/21, 25/27, 21/19, 24/21, 35/27, 27/19, 15/9, 36/21, 36/19, 55/27, 48/21)$$

with $|D^{***}| = |E^*| = 17$.

C. A METHOD TO EVALUATE THE DIFFERENCES BETWEEN A PAIR OF CLASSIFICATION SYSTEMS

Let x_j , x_g , and x_k be the sets of the top $X\%$ most cited articles in cluster citation distributions c_j and c_g , and category citation distribution c_k . Denote the union of these sets by $X^{G8} = \cup_j \{x_j\}$, $X^{G6} = \cup_g \{x_g\}$, and $X^{Wos} = \cup_k \{x_k\}$. There are two ways of comparing the WoS system with systems G8 or G6: in terms of distinct articles, or in terms of extended articles. Let us begin by defining the set $X^{Wos}(\text{Dist})$ of distinct articles in X^{Wos} . We can now form the following two intersections of distinct articles:

$$X^{W8}(\text{Dist}) = X^{Wos}(\text{Dist}) \cap X^{G8}$$

and

$$X^{W6}(\text{Dist}) = X^{Wos}(\text{Dist}) \cap X^{G6}.$$

The difference between the top $X\%$ most cited articles in the G8 and WoS systems can be measured by the percentage that the articles in $X^{G8} - X^{W8}(\text{Dist})$ represent in X^{G8} . Similarly, the difference between the top $X\%$ most cited articles in the G6 and WoS systems can be measured by the percentage that the articles in $X^{G6} - X^{W6}(\text{Dist})$ represent in X^{G6} . The results for these two expressions are presented in Tables 2.a and 2.b under the heading ‘‘Difference in distinct articles’’.

On the other hand, let $X^{G6}(\text{Ext})$ and $X^{G8}(\text{Ext})$ be the sets of extended articles in X^{G6} and X^{G8} . We can now form the following two intersections of extended articles:

$$\mathbf{X}^{W8}(\text{Ext}) = \mathbf{X}^{WoS} \cap \mathbf{X}^{G8}(\text{Ext})$$

and

$$\mathbf{X}^{W6}(\text{Ext}) = \mathbf{X}^{WoS} \cap \mathbf{X}^{G6}(\text{Ext}).$$

The difference between the top $X\%$ most cited articles in the G8 and WoS systems can be measured by the percentage that the articles in $\mathbf{X}^{WoS} - \mathbf{X}^{W8}(\text{Ext})$ represents in \mathbf{X}^{WoS} . Similarly, the difference between the top $X\%$ most cited articles in the G6 and WoS systems can be measured by the percentage that the articles in $\mathbf{X}^{WoS} - \mathbf{X}^{W6}(\text{Ext})$ represents in \mathbf{X}^{WoS} . The results for these two expressions are presented in Tables 2.a and 2.b under the heading “Difference in extended articles”.

For the comparison between the G6 and G8 systems, let \mathbf{X}^{G8} be the set of distinct articles common to both systems, namely, let $\mathbf{X}^{G8} = \mathbf{X}^{G6} \cap \mathbf{X}^{G8}$. There are two ways of measuring the difference between the top $X\%$ most cited articles in both systems: through the percentage that the articles in $\mathbf{X}^{G8} - \mathbf{X}^{G6}$ represent in \mathbf{X}^{G8} , and through the percentage that the articles in $\mathbf{X}^{G6} - \mathbf{X}^{G8}$ represent in \mathbf{X}^{G6} . The results for these two expressions are presented under the headings “Difference in terms of the G8 system” and “Difference in terms of the G6 system” in Table 2.c.