

2. Clustering validity indices

This section explains the internal and external indices that were implemented and used in this work. The indices are presented according to the type of analysis these indices perform. Fig. 1 shows the organization of the indices. Before introducing the each group of indices, preliminary definitions, which some indices have in common, are listed. The indices have in their name the rule to choice a specific partition, i.e.:

- $[\uparrow]$ indicates that partitions with the greatest value of the index are desired.
- $[\downarrow]$ indicates that partitions with the lowest value of the index are desired.
- $[\text{diff } \uparrow]$ indicates that is a difference-like index that grows as long as the number of clusters increase.
- $[\text{diff } \downarrow]$ indicates that is a difference-like index that decreases as long as the number of clusters increase.

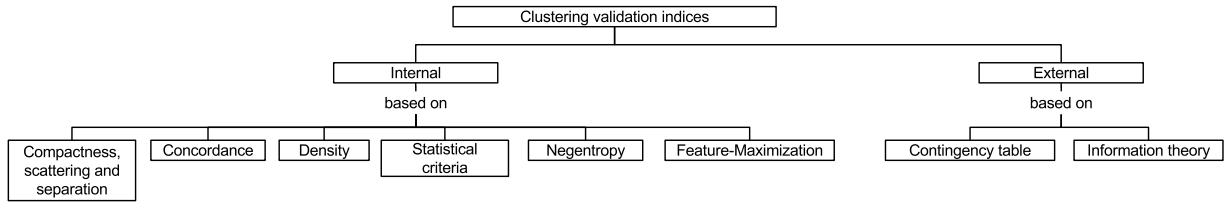


Figure 1: Clasification of internal and external indices according to the analysis these indices perform.

2.1. Internal indices

2.1.1. Indices based on compactness, scattering and separation

• Preliminary definitions:

- **Data matrix $X_{N \times p}$:** Dataset matrix that contains N observations \mathbf{x}_i of p features.
- **Obtained partition P_o :** Clustering of K clusters resulted from the clustering algorithm.
- Each cluster c_k can be represented by a subset $X_{n_k \times p}^{(k)}$, which contains the n_k observations belonging to this cluster.
- **Barycenter μ^k of c_k :**

$$\mu^k = \frac{1}{n_k} \sum_{\mathbf{x}_i \in X^{(k)}} \mathbf{x}_i \quad (1)$$

- **Data barycenter μ :**

$$\mu = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \quad (2)$$

- **Total dispersion:** The scatter matrix T is defined by:

$$T = (X - \mu_n)'(X - \mu_n) \quad (3)$$

Where μ_n a $n \times p$ matrix generated repeating the row vector μ n times along the rows.

The total dispersion or total sum of squares (TSS) is the trace of the matrix T .

- **Within-group dispersion:** The within-group scatter matrix $WG^{(k)}$ for the group k is defined by:

$$WG^{(k)} = (X^{(k)} - \mu_{n_k}^k)'(X^{(k)} - \mu_{n_k}^k) \quad (4)$$

Where $\mu_{n_k}^k$ is a $n_k \times p$ matrix generated repeating the row vector μ^k , n_k times along the rows.

The within-cluster dispersion or within-group sum of squares (WGSS) is the sum of the traces of matrices $WG^{(k)}$.

- **Between-group dispersion:** This metric measures the dispersion of the clusters between each other. The between-group scatter matrix is defined by:

$$BG = B' B \quad (5)$$

Where B is the matrix formed by $(\mu^k - \mu)$ for each cluster in each row, generating a $K \times p$ matrix.

- WGSS, BGSS and TSS can be expressed as a sum of square distances:

$$WGSS = \sum_{k=1}^K WGSS^{(k)} = \sum_{k=1}^K \sum_{x_i \in X^{(k)}} \|x_i - \mu^k\|^2 \quad TSS = \sum_{i=1}^n \|x_i - \mu\|^2 \quad BGSS = \sum_{k=1}^K n_k \|\mu^k - \mu\|^2$$

Where $\|\cdot\|$ represents the Euclidean distance.

- **Columns vectors v_j :** These vector are obtained extracting from $X_{n \times p}$ the p columns vectors.

$$X_{n \times p} = [v_1 \ v_2 \ \dots \ v_p]_{n \times p} \quad (6)$$

- **Columns vectors v_j per cluster k :** These vector are obtained extracting from $X_{n_k \times p}^{(k)}$ the p columns vectors.

$$X_{n_k \times p}^{(k)} = [v_1^{(k)} \ v_2^{(k)} \ \dots \ v_p^{(k)}]_{n_k \times p} \quad (7)$$

- **Number of pairs of points n_T , n_W and n_B :** In total there are n_T distinct pairs of observations in X .

$$n_T = \binom{n}{2} = \frac{n(n-1)}{2} \quad (8)$$

Similarly, if only the observations of the cluster k are considered, there are n_W^k distinct pairs of observations in $X^{\{k\}}$.

$$n_W^{\{k\}} = \binom{n_k}{2} = \frac{n_k(n_k - 1)}{2} \quad (9)$$

n_W is obtained adding all the distinct pairs within the clusters:

$$n_W = \sum_{k=1}^K n_W^{\{k\}} = \sum_{k=1}^K \frac{n_k(n_k - 1)}{2} \quad (10)$$

Finally, n_B , the number of distinct pairs between clusters can be calculated by:

$$n_T = n_W + n_B \quad n_B = n_T - n_W = \sum_{k < k'} n_k n_{k'} \quad (11)$$

– **Sum of distances s_W and s_B :** The sum of the n_W within-cluster distances s_W can be calculated by:

$$s_W = \sum_{k=1}^K \sum_{\substack{\mathbf{x}_i, \mathbf{x}_j \in X^{\{k\}} \\ i < j}} \|\mathbf{x}_i - \mathbf{x}_j\| \quad (12)$$

Similarly, the sum of the n_B between-cluster distances s_B can be calculated by:

$$s_B = \sum_{k < k'}^K \sum_{\substack{\mathbf{x}_i \in X^{\{k\}} \mathbf{x}_j \in X^{\{k'\}} \\ i < j}} \|\mathbf{x}_i - \mathbf{x}_j\| \quad (13)$$

– **Point Symmetry-Distance of a point with respect a cluster k :** Based on the distance between the symmetric point of \mathbf{x}_i with respect to the centroid of a cluster k (i.e., $2\boldsymbol{\mu}^k - \mathbf{x}_i$) and the m -nearest points in the cluster k (Bandyopadhyay and Saha, 2008).

$$d_{ps}(\mathbf{x}_i, c_k) = \frac{1}{m} \sum \min_{\mathbf{x}_j \in X^{\{k\}}} (m) \left\{ \left\| (2\boldsymbol{\mu}^k - \mathbf{x}_i) - \mathbf{x}_j \right\| \right\} \quad (14)$$

Where $\sum \min_{\mathbf{x}_j \in X^{\{k\}}} (m) \{\cdot\}$ computes the sum of the m lowest values of $\{\cdot\}$.

• **Indices based on compactness, scattering and separation**

– **Ball-Hall [diff ↑]:** This index is defined as the average of within-cluster dispersions weighted by the inverse of n_k (Ball and Hall, 1965a):

$$\text{Ball-Hall} = \frac{1}{K} \sum_{k=1}^K \frac{1}{n_k} WGS S^{\{k\}} \quad (15)$$

– **Banfeld-Raftery [↓]:** Banfeld-Raftery is defined as the weighted sum of the logarithms of the within-

cluster dispersion (Banfield and Raftery, 1993):

$$\text{Banfield-Raftery} = \sum_{k=1}^K n_k \log \left(\frac{WGSS^{(k)}}{n_k} \right) \quad (16)$$

- **C index** [\downarrow]: This measure is based on the within-group distances, as well as on their maximum and minimum possible value (Hubert and Schultz, 1976):

$$C = \frac{s_W - s_{\min}}{s_{\max} - s_{\min}} \quad (17)$$

Where s_{\min} and s_{\max} are the sum of the n_W smallest and largest distances between all the pairs of points of $X_{N \times p}$.

- **Calinski-Harabasz** [\uparrow]: The idea behind this index is to evaluate partitions based on the average between- and within-cluster sum of squares (Caliński and Harabasz, 1974):

$$\text{Calinski-Harabasz} = \frac{n - K}{K - 1} \frac{BGSS}{WGSS} \quad (18)$$

- **Davies-Bouldin** [\downarrow]: Two terms are defined in order to estimate this index: δ_k , the mean distance of the cluster k points to their barycenter μ^k and $\Delta_{kk'}$, the distance between the barycenters μ^k and $\mu^{k'}$ (Davies and Bouldin, 1979):

$$\delta_k = \frac{1}{n_k} \sum_{k=1}^K \sum_{x_i \in X^{(k)}} \|x_i - \mu^k\| \quad \Delta_{kk'} = \|\mu^{k'} - \mu^k\| \quad (19)$$

$$\text{Davies-Bouldin} = \frac{1}{K} \sum_{k=1}^K \max_{k' \neq k} \left(\frac{\delta_k + \delta_{k'}}{\Delta_{kk'}} \right) \quad (20)$$

- **Det-Ratio** [$\text{diff } \downarrow$]: This is based the determinant of the total dispersion matrix and the sum of the within-group dispersion matrices (Scott and Symons, 1971).

$$\text{Det-Ratio} = \frac{\det(\mathbf{T})}{\det(\mathbf{WG})} \quad (21)$$

- **Dunn's indices** [\uparrow]: These indices are based on two estimators, one focus on separation and the another in cohesion (compactness).

$$\text{Dunn index} = \min_{1 \leq k \leq c} \left\{ \min_{\substack{1 \leq k' \leq c \\ j \neq i}} \left\{ \frac{\delta_s(X^{(k)}, X^{(k')})}{\max_{1 \leq k \leq c} \Delta_d(X^{(k)})} \right\} \right\} \quad (22)$$

Where δ_s is the separation estimator and Δ_d is the cohesion estimator. In (Bezdek and Pal, 1998) 17

variants of the original Dunn index (Dunn, 1974) were proposed. These variants and the original Dunn index are produced by combining one of the 6 ways to compute δ_s and one of the 3 ways to compute Δ_d :

$$\delta_1(X^{(k)}, X^{(k')}) = \min_{x_i \in X^{(k)}} \left\{ \min_{x_j \in X^{(k')}} \|x_i - x_j\| \right\} \quad (23)$$

$$\delta_2(X^{(k)}, X^{(k')}) = \max_{x_i \in X^{(k)}, x_j \in X^{(k')}} \|x_i - x_j\| \quad (24)$$

$$\delta_3(X^{(k)}, X^{(k')}) = \frac{1}{n_k n_{k'}} \sum_{x_i \in X^{(k)}} \sum_{x_j \in X^{(k')}} \|x_i - x_j\| \quad (25)$$

$$\delta_4(X^{(k)}, X^{(k')}) = \|\mu^k - \mu^{k'}\| \quad (26)$$

$$\delta_5(X^{(k)}, X^{(k')}) = \frac{1}{n_k + n_{k'}} \left(\sum_{x_i \in X^{(k)}} \|x_i - \mu^{k'}\| + \sum_{x_j \in X^{(k')}} \|x_j - \mu^k\| \right) \quad (27)$$

$$\delta_6(X^{(k)}, X^{(k')}) = \max \left\{ \max_{x_i \in X^{(k)}} \left\{ \min_{x_j \in X^{(k')}} \|x_i - x_j\| \right\}, \max_{x_j \in X^{(k')}} \left\{ \min_{x_i \in X^{(k)}} \|x_i - x_j\| \right\} \right\} \quad (28)$$

$$\Delta_1(X^{(k)}) = \max_{x_i \in X^{(k)}} \left\{ \max_{x_j \in X^{(k)}} \|x_i - x_j\| \right\} \quad (29)$$

$$\Delta_2(X^{(k)}) = \frac{1}{n_k(n_k - 1)} \sum_{x_i \neq x_j \in X^{(k)}} \|x_i - x_j\| \quad (30)$$

$$\Delta_3(X^{(k)}) = \frac{2}{n_k} \sum_{x_i \in X^{(k)}} \|x_i - \mu^k\| \quad (31)$$

In this work, we refers a specific combination Δ_i and δ_j by $\text{Dunn}ij$. The original Dunn index is denoted by $\text{Dunn}11$.

– **Ksq DetW [diff ↑]**: Index that is also based in within dispersion matrices (Marriott, 1971).

$$\text{Ksq DetW} = K^2 \det(\mathbf{WG}) \quad (32)$$

– **Log Det Ratio [diff ↓]**: This index is the logarithm of the Det-Ratio index (Scott and Symons, 1971).

$$\text{Log Det Ratio} = n \log (\text{Det-Ratio}) = n \log \left(\frac{\det(\mathbf{T})}{\det(\mathbf{WG})} \right) \quad (33)$$

– **Log SS Ratio [diff ↓]**:

$$\text{Log SS Ratio} = \log \left(\frac{\text{BGSS}}{\text{WGSS}} \right) \quad (34)$$

– **McClain-Rao [↓]**: This measure is defined as the quotient between the mean within-cluster and between-

cluster distances (McClain and Rao, 1975):

$$\text{McClain-Rao} = \frac{\frac{s_W}{n_W}}{\frac{s_B}{n_B}} = \frac{n_B}{n_W} \frac{s_W}{s_B} \quad (35)$$

- **PBM index** [↑]: Also called I index in literature (Pakhira et al., 2004).

$$\text{PBM} = \left(\frac{1}{K} \times \frac{E_T}{E_W} \times D_B \right) \quad (36)$$

Where D_B is maximum distance between two cluster barycenters. E_W is the sum of the distances of the cluster objects to its barycenter. E_T is the sum of the distances of the objects to the entire data barycenter. Their formulas are:

$$E_T = \sum_{i=1}^n \|x_i - \mu\| \quad E_W = \sum_{k=1}^K \sum_{x_i \in X^{(k)}} \|x_i - \mu^k\| \quad D_B = \max_{k < k'} \|\mu^k - \mu^{k'}\|$$

- **Point-Biserial** [↑]: Index based on the statistics point-biserial, where the continuous variable A is a set of n_T possible distinct distances of the dataset. The binary variable B equals 1 if the corresponding distance of A is between two points of the same cluster. B equal 0 otherwise. (Milligan, 1981)

$$\text{Point-Biserial} = \left(\frac{s_W}{n_W} - \frac{s_B}{n_B} \right) \frac{\sqrt{n_W n_B}}{n_T} \quad (37)$$

- **Ratkowsky-Lance** [↑]: This index is based on the between-groups and total dispersion matrix (Ratkowsky and Lance, 1978). A term in denominator is included to avoid the monotone increasing as the number of clusters increases.

$$\text{Ratkowsky-Lance} = \sqrt{\frac{\bar{R}}{K}} = \frac{\bar{c}}{\sqrt{K}} \quad (38)$$

$$\bar{c}^2 = \bar{R} = \frac{1}{p} \sum_{j=1}^p \frac{BGSS_j}{TSS_j} \quad (39)$$

Where $BGSS_j = b_{jj}$ and $TSS_j = t_{jj}$ are the j-element of the diagonal of BGSS and of TSS, respectively.

- **Ray-Turi** [↓]: Ray-Turi is based in the same denominator used the Davies-Bouldin index ($\Delta_{kk'}$, the distance between the barycenters μ^k and $\mu^{k'}$) (Ray and Turi, 1999).

$$\text{Ray-Turi} = \frac{1}{n} \frac{WGSS}{\min_{k \neq k'} \Delta_{kk'}} \quad (40)$$

- **Scott-Symons** [↓]: This index is the weighted sum of the logarithms of the determinants of the variance-

covariance matrix of each cluster (Scott and Symons, 1971):

$$\text{Scott-Symons} = \sum_{k=1}^K n_k \log \det \left(\frac{\mathbf{W}G^{(k)}}{n_k} \right) \quad (41)$$

- **SD index** [\downarrow]: This index is based on the average scattering and the total separation of clusters (Halkidi et al., 2000).

$$\text{SD} = \alpha S + D \quad (42)$$

$$S = \frac{\frac{1}{K} \sum_{k=1}^K \|\mathbf{v}^{(k)}\|}{\|\mathbf{v}\|} \quad D = \frac{\max_{k \neq k'} \Delta_{kk'}}{\min_{k \neq k'} \Delta_{kk'}} \times \sum_{k=1}^K \frac{1}{\sum_{\substack{k'=1 \\ k' \neq k}}^K \Delta_{kk'}} \quad (43)$$

Where α equals to the value of D obtained for the partition with the greatest number of cluster, and \mathbf{v} and $\mathbf{v}^{(k)}$ are the variances vector of the columns vector \mathbf{v}_j and $\mathbf{v}_j^{(k)}$, respectively:

$$\mathbf{v}^{(k)} = \left[\text{Var}(\mathbf{v}_1^{(k)}), \dots, \text{Var}(\mathbf{v}_p^{(k)}) \right] \quad \mathbf{v} = \left[\text{Var}(\mathbf{v}_1), \dots, \text{Var}(\mathbf{v}_p) \right] \quad (44)$$

- **Silhouette index** [\uparrow]: Silhouette is based on two terms computed for each sample \mathbf{x}_i : $a_{\text{SIL}}(i)$ and $b_{\text{SIL}}(i)$. $a_{\text{SIL}}(i)$ represents the mean distance of object i to the other objects in its cluster, while $b_{\text{SIL}}(i)$ is the smallest mean distance of an object i to the other clusters (Ratkowsky and Lance, 1978).

$$\text{Silhouette} = \frac{1}{K} \sum_{k=1}^K \frac{1}{n_k} \sum_{\mathbf{x}_i \in X^{(k)}} \frac{b_{\text{SIL}}(i) - a_{\text{SIL}}(i)}{\max(a_{\text{SIL}}(i), b_{\text{SIL}}(i))} \quad (45)$$

$$a_{\text{SIL}}(i) = \frac{1}{n_k - 1} \sum_{\substack{\mathbf{x}_i, \mathbf{x}_j \in X^{(k)} \\ \mathbf{x}_i \neq \mathbf{x}_j}} \|\mathbf{x}_i - \mathbf{x}_j\| \quad b_{\text{SIL}}(i) = \min_{k' \neq k} \left\{ \frac{1}{n_{k'}} \sum_{\mathbf{x}_j \in X^{(k')}} \|\mathbf{x}_i - \mathbf{x}_j\| \right\} \quad (46)$$

- **Sym** [\uparrow]: Adaptation of the PBM index using the point symmetric distance (Bandyopadhyay and Saha, 2008).

$$\text{Sym} = \frac{\max_{k < k'} \|\boldsymbol{\mu}^k - \boldsymbol{\mu}^{k'}\|}{K \sum_{k=1}^K \sum_{\mathbf{x}_i \in X^{(k)}} d_{ps}(\mathbf{x}_i, c_k)} \quad (47)$$

- **Modified Davies–Bouldin** [\downarrow]: This index is a modification of the Davies–Bouldin index (Saha and Bandyopadhyay, 2009), which uses the point symmetric distance.

$$\delta_k = \frac{1}{n_k} \sum_{\mathbf{x}_i \in X^{(k)}} d_{ps}(\mathbf{x}_i, c_k) \quad (48)$$

- **Modified Dunn11** [\uparrow]: This index is a modification of the Dunn11 (Saha and Bandyopadhyay, 2009),

which uses the point symmetric distance.

$$\Delta(X^{(k)}) = \max_{x_i \in X^{(k)}} d_{ps}(x_i, c_k) \quad (49)$$

- **Modified Dunn33** [↑]: This index is a modification of the Dunn33 (Saha and Bandyopadhyay, 2009), which uses the point symmetric distance.

$$\Delta_3(X^{(k)}) = \frac{2}{n_k} \sum_{x_i \in X^{(k)}} d_{ps}(x_i, c_k) \quad (50)$$

- **Trace W** [diff ↑]: Index previously studied in (Edwards and Cavalli-Sforza, 1965).

$$\text{Trace W} = \text{Tr}(\mathbf{W}\mathbf{G}) = \mathbf{W}\mathbf{G}\mathbf{S}\mathbf{S} \quad (51)$$

- **Trace WiB** [diff ↑]: This index was studied in (Friedman and Rubin, 1967).

$$\text{Trace W} = \text{Tr}(\mathbf{W}\mathbf{G}^{-1}\mathbf{B}\mathbf{G}) \quad (52)$$

- **Wemmert-Gańczarski** [↑]: The Wemmert-Gańczarski index is built using proportions of distances between the points and the clusters barycenters (equation 54).

$$\text{Wemmert-Gańczarski} = \frac{1}{n} \sum_{k=1}^K n_k J_k \quad (53)$$

$$J_k = \max \left\{ 0, 1 - \frac{1}{n_k} \sum_{x_i \in X^{(k)}} R(x_i) \right\} \quad R(x_i) = \frac{\|x_i - \mu^k\|}{\min_{k' \neq k} \|x_i - \mu^{k'}\|} \quad (54)$$

- **Xie-Beni index** [↓]: It is defined as the quotient between the mean quadratic error and the minimum of the minimal squared distances between the points in the clusters. (Xie and Beni, 1991)

$$\text{Xie-Beni} = \frac{1}{n} \frac{\mathbf{W}\mathbf{G}\mathbf{S}\mathbf{S}}{\min_{k < k'} \delta_1(X^{(k)}, X^{(k')})} \quad (55)$$

- **CVNN** [↓]: This is a more recent index, which has a intercluster separation based on L nearest neighbors (in the original paper was used the letter k instead of L to denote the nearest neighbors) and a compactness measure that does not take into account the barycenter in order to deal with arbitrary-shaped data (Liu

et al., 2013).

$$\text{Sep}(K, v) = \max_k \left\{ \frac{1}{n_k} \sum_{j=1}^{n_k} \frac{q_j}{L} \right\} \quad \text{Comp}(K, v) = \Delta_2(X^{(k)}) = \frac{1}{n_k(n_k - 1)} \sum_{x_i \neq x_j \in X^{(k)}} \|x_i - x_j\| \quad (56)$$

Where q_j is the number of nearest neighbors which are not in the same cluster k of object j . CVNN is then defined as the sum normalized compatness and separation:

$$\text{CVNN} = \text{Sep}_{\text{norm}}(K, v) + \text{Comp}_{\text{norm}}(K, v) \quad (57)$$

$$\text{Sep}_{\text{norm}}(K, v) = \frac{\text{Sep}(K, v)}{\max_{K_{\min} < K' < K_{\min}} \text{Sep}(K', v)} \quad \text{Comp}_{\text{norm}}(K, v) = \frac{\text{Comp}(K, v)}{\max_{K_{\min} < K' < K_{\min}} \text{Comp}(K', v)} \quad (58)$$

– **WB** [\downarrow]: Index proposed by (Zhao et al., 2009a).

$$\text{WB} = K \frac{WGS S}{BGS S} \quad (59)$$

– **Xu** [\downarrow]: Index proposed by (Xu, 1997b).

$$\text{Xu} = \log \frac{K \times WGS S^{\frac{p}{2}}}{BGS S^{\frac{p}{2}}} \quad (60)$$

– **CS** [\downarrow]: Index whose numerator estimates the cohesion by the cluster diameters. The separation is estimated in the denominator by nearest neighbor distances. (Chou et al., 2004)

$$\text{CS} = \frac{\sum_{k=1}^K \frac{1}{n_k} \sum_{x_i \in X^{(k)}} \max_{x_j \in X^{(k)}} \|x_i - x_j\|}{\sum_{k=1}^K \min_{\mu^{k'}} \|\mu^k - \mu^{k'}\|} \quad (61)$$

– **Score function** [\uparrow]: Index maily based on a compactness term call wcd and on a separation term called bcd (Saitta et al., 2007). As the term of separation becomes larger than the term of compactness, the upper exponent becomes infinite, like the inferior exponent, resulting in an index equal to 1.

$$\text{SF} = 1 - \frac{1}{\exp^{\exp^{\text{bcd} - \text{wcd}}}} \quad (62)$$

$$\text{bcd} = \frac{\sum_{k=1}^K n_k \|\mu^k - \mu\|}{n \times K} \quad \text{wcd} = \sum_{k=1}^K \frac{1}{n_k} \sum_{x_i \in X^{(k)}} \|x_i - \mu^k\| \quad (63)$$

– **COP** [\downarrow]: Index whose compactness term is based on the distance from the objects to their cluster, and

whose separation term is the furthest distance between points in different clusters (Gurrutxaga et al., 2010).

$$\text{COP} = \frac{1}{n} \sum_{k=1}^K n_k \frac{\text{intra}_{\text{COP}}(k)}{\text{inter}_{\text{COP}}(k)} \quad (64)$$

$$\text{intra}_{\text{COP}}(k) = \frac{1}{n_k} \sum_{x_i \in X^{(k)}} \|x_i - \mu^k\| \quad \text{inter}_{\text{COP}}(k) = \min_{x_i \notin X^{(k)}} \max_{x_j \in X^{(k)}} \|x_i - x_j\| \quad (65)$$

- **SV** [\uparrow]: This index estimates the separation criterion in the numerator by finding the closest distance between clusters centroids. Its compactness term is based on the 10 percent of the lowest distances from the cluster objects to their centroid (Žalik and Žalik, 2011). Operators $\sum \min_{x_j \in X^{(k)}}(n)$ and $\sum \max_{x_j \in X^{(k)}}(n)$ are used to compute the compactness term.

$$\text{SV} = \frac{\sum_{k=1}^K \min_{\mu^{k'}} \{\|\mu^k - \mu^{k'}\|\}}{\sum_{k=1}^K \frac{1}{n_k} \sum \max_{x_j \in X^{(k)}} (0.1n_k) \{\|x_i - \mu^k\|\}} \quad (66)$$

Where $\sum \max_{x_j \in X^{(k)}}(m) \{\cdot\}$ computes the sum of the m greatest values of $\{\cdot\}$.

- **OS** [\uparrow]: This index uses the same compactness term that the SV index, but uses a more complex separation term that involves a function called ov (Žalik and Žalik, 2011). This function requires two terms: $a_{OS}(x_i, \mu^k)$ and $b_{OS}(x_i, \mu^k)$. $a_{OS}(x_i, \mu^k)$ is the sum of distances between pairs in the cluster k and $b_{OS}(x_i, \mu^k)$ is the sum of the n_k minimum distances between a point in a cluster k and points in other clusters.

$$\text{OS} = \frac{\sum_{k=1}^K \sum_{x_i \in X^{(k)}} \text{ov}(x_i, \mu^k)}{\sum_{k=1}^K \frac{1}{n_k} \sum \max_{x_j \in X^{(k)}} (0.1n_k) \{\|x_i - \mu^k\|\}} \quad (67)$$

Where:

$$\text{ov}(x_i, \mu^k) = \begin{cases} \frac{a_{OS}(x_i, \mu^k)}{b_{OS}(x_i, \mu^k)} & \text{if } \frac{b_{OS}(x_i, \mu^k) - a_{OS}(x_i, \mu^k)}{b_{OS}(x_i, \mu^k) + a_{OS}(x_i, \mu^k)} < 0.4 \\ 0 & \text{otherwise} \end{cases} \quad (68)$$

$$b_{OS}(x_i, \mu^k) = \frac{1}{n_k} \sum \min_{x_j \notin X^{(k)}} (n_k) \{\|x_i - x_j\|\} \quad a_{OS}(x_i, \mu^k) = \frac{1}{n_k} \sum_{x_j \in X^{(k)}} \|x_i - x_j\| \quad (69)$$

- **AR_{SD}** [\uparrow]: This index is based on a compactness measure (CM_{SD}), which is obtained using the standard deviation and the maximum diameter of each clusters, and separation measure ($\text{DM}_{\text{clusters}}$) based on a penalty function that is calculated using recursive algorithm (Wani and Riyaz, 2016b).

$$\text{AR}_{\text{SD}} = \text{CM}_{\text{SD}} + \text{DM}_{\text{clusters}} = \text{CM}_{\text{SD}} - Fn \quad (70)$$

Where:

$$CM_{SD} = \frac{1}{n_k} \sum_{k=1}^K \frac{md_k - SD_k}{md_k} \quad (71)$$

$$md_k = \max_{x_i, x_j \in X^{(k)}} \{\|x_i - x_j\|\} \quad (72)$$

Where SD_k is the standard deviation of the cluster k and the steps to compute F_n are:

1. Set the penalty function $F_n = 0$ and a variable $i=K$.
2. Select any cluster from the list of K clusters.
3. Compute the average of minimum distance with the cluster:

$$d_{\min(\text{intra})} = \frac{1}{n_k} \sum_{x_i \in X^{(k)}} \min_{x_j \in X^{(k)}} \{\|x_i - x_j\|\} \quad (73)$$

4. Compute between clusters distances:

$$d_{\min(\text{inter})} = \min_{x_i \in X^{(k)} \cap x_j \notin X^{(k)}} \{\|x_i - x_j\|\} \quad (74)$$

5. Compute the penalty function for the selected cluster:

$$fn = \begin{cases} 0 & \text{if } d_{\min(\text{inter})} > 2d_{\min(\text{intra})} \\ a & \text{if } d_{\min(\text{inter})} \leq 2d_{\min(\text{intra})} \\ a \times \frac{d_{\min(\text{inter})} - d_{\min(\text{intra})}}{d_{\min(\text{intra})}} & \text{otherwise} \end{cases} \quad (75)$$

Where a was set experimentally to 0.1.

6. Accumulate the penalty function in F_n . $F_n = F_n + fn$
7. Remove the selected cluster from the list of clusters and decrement i .
8. Repeat from step 2 until $i < 2$.

- **STR** [\uparrow]: This index is based on two components: cluster compactness $E(K)$ and cluster separation $D(K)$ (Starczewski, 2017b). $E(K)$ is defined as the ratio of the total scatter to the within scatter of the clusters. Its formula is:

$$STR(K) = [E(K + 1) - E(K)] \times [D(K + 2) - D(K + 1)] \quad (76)$$

$$E(K) = \frac{E_0}{E_K} = \frac{\sum_{j=1}^n \|x_j - \mu\|}{\sum_{k=1}^K \sum_{x_i \in X^{(k)}} \|x_i - \mu^k\|} \quad D(K) = \frac{D_{Kmax}}{D_{Kmin}} = \frac{\max_{\mu^k, \mu^{k'} \in P_o} \{\|\mu^k - \mu^{k'}\|\}}{\min_{\mu^k, \mu^{k'} \in P_o} \{\|\mu^k - \mu^{k'}\|\}} \quad (77)$$

- **Bhargavi-Gowda Index** [$\text{diff } \uparrow$]: This index is based on two ratios, the ratio of WGSS and BGSS called "sum-of-squares ratio", and the ratio of intra-cluster distances and inter-cluster distances called "intra-inter

ratio” (Bhargavi and Gowda, 2015).

$$\text{Bhargavi-Gowda} = \left| \frac{WGS S}{BGS S} \times TSS - \frac{C_{\text{Intra}}}{C_{\text{Inter}}} - (n - k) \right| \quad (78)$$

$$C_{\text{Intra}} = \sum_{k=1}^K \sqrt{\sum_{i=1}^{n_k} \sum_{j=1}^p (x_{ij} - \mu_j^k)^2} \quad C_{\text{Inter}} = \frac{1}{K} \sum_{k=1}^K \left[\frac{1}{K} \sum_{k'=1, k' \neq k}^K \sqrt{\sum_{j=1}^p (\mu_j^k - \mu_j^{k'})^2} \right] \quad (79)$$

- **CVM** [\uparrow]: This index favours solutions where the size of the cluster cores is maximized together with their separation (De Morsier et al., 2015b). It is defined as:

$$\text{CVM} = \left(\sum_{k=1}^K \frac{d_{\text{in}}(k)}{R_{\text{in}}(k)} \right) \left(\sum_{k=1}^K d_{\text{inter}}(k) \right) \quad (80)$$

Where $d_{\text{in}}(k)$, $d_{\text{inter}}(k)$ and $R_{\text{in}}(k)$ are, respectively, the maximal distance among the samples in the cluster k , the core separability defined as the minimum distance between a cluster sample to the closest sample from another cluster and the core homogeneity that is the ratio between the maximum and the average smallest inner distances inside a cluster.

$$d_{\text{in}}(k) = \max_{x_i, x_j \in X^{(k)}} \|x_i - x_j\|^2 \quad d_{\text{inter}}(k) = \max_{x_i \in X^{(k)}, x_j \in X^{(k')}} \|x_i - x_j\|^2 \quad (81)$$

$$R_{\text{in}}(k) = \frac{\max_{x_i \in X^{(k)}} \left(\min_{x_j \in X^{(k)}} \|x_i - x_j\|^2 \right)}{\frac{1}{n_k} \sum_{x_i \in X^{(k)}} \left(\min_{x_j \in X^{(k)}, x_i \neq x_j} \|x_i - x_j\|^2 \right)} \quad (82)$$

$R_{\text{in}}(k)$ takes values from 1 to ∞ .

2.1.2. Indices based on concordance:

• Preliminary definitions:

- **Number of concordant pairs of objects** s_+ : The number of times that the distance between a pair of objects from the same cluster is lower than the distance between a pair of objects from different cluster:

$$s_+ = \frac{1}{2} \sum_{k=1}^K \sum_{\substack{x_i, x_j \in X^{(k)} \\ x_i \neq x_j}} \frac{1}{2} \sum_{k'=1}^K \sum_{\substack{x_a \in X^{(k')} \\ x_b \notin X^{(k')}}} \delta(\|x_i - x_j\| < \|x_a - x_b\|) \quad (83)$$

Where $\delta(\cdot)$ equals 1 if the inner inequality is satisfied and equals 0 if not.

- **Number of discordant pairs of objects** s_- : The number of times that the distance between a pair of

objects from the same cluster is lower than the distance between a pair of objects from different cluster:

$$s_- = \frac{1}{2} \sum_{k=1}^K \sum_{\substack{x_i, x_j \in X^{(k)} \\ x_i \neq x_j}} \frac{1}{2} \sum_{k'=1}^K \sum_{\substack{x_a \in X^{(k')} \\ x_b \notin X^{(k')}}} \delta(\|x_i - x_j\| > \|x_a - x_b\|) \quad (84)$$

• **Indices based on concordance**

– **Baker-Hubert Gamma** [\uparrow]: (Rohlf, 1974).

$$\text{Baker-Hubert Gamma} = \frac{s_+ - s_-}{s_+ + s_-} \quad (85)$$

– **G plus** [\downarrow]: This index is defined as the proportion of discordant pairs among all the pairs of distinct points (Rohlf, 1974).

$$G (+) = \frac{2s_-}{n_T(n_T - 1)} \quad (86)$$

– **Tau** [\uparrow]: s_+ and s_- do not count if a between-cluster distance and a within-cluster distance are equal. This index includes in the denominator a term that corrects this issue (Rohlf, 1974).

$$\tau = \frac{s_+ - s_-}{\sqrt{n_B n_W \left(\frac{n_T(n_T - 1)}{2} \right)}} \quad (87)$$

2.1.3. *Density-based indices:*

• **Indices based on density**

– **S_{Dbw} validity index** [\downarrow]:

The density $\gamma_{kk'}$ of a given point, relative to two clusters c_k and $c_{k'}$, is equal to the number of points in these two clusters whose distance to this point is less than σ (Halkidi et al., 2001).

$$\gamma_{kk'}(\mathbf{x}) = \sum_{\mathbf{x}_i \in X^{(k)} \cup \mathbf{x}_i \in X^{(k')}} \delta(\|\mathbf{x}_i - \mathbf{x}\| > \sigma) \quad (88)$$

$$\sigma = \frac{1}{K} \sqrt{\sum_{k=1}^K \|\mathbf{p}^{(k)}\|} \quad (89)$$

Where $\delta(\cdot)$ equals 1 if the inner inequality is satisfied and equals 0 if not. S_{Dbw} is defined by:

$$S_{Dbw} = S + \mathbb{G} \quad (90)$$

$$\mathbb{G} = \frac{2}{K(K-1)} \sum_{k < k'} R_{k,k'} \quad R_{k,k'} = \frac{\gamma_{kk'}(\mathbf{H}_{kk'})}{\max(\gamma_{kk'}(\boldsymbol{\mu}^k), \gamma_{kk'}(\boldsymbol{\mu}^{k'}))} \quad (91)$$

Where $H_{kk'}$ is the midpoint if the barycenters μ^k and $\mu^{k'}$.

2.1.4. Statistical-based indices:

- **Preliminary definitions:**

- **Assumption:** The data is modeled by n gaussians, each with identical variance σ and different means μ^k (This assumption is also called "identical spherical assumption") (Pelleg et al., 2000).
- **Probability of an object i is an element of cluster c_k under the maximum likelihood:**

$$P(\mathbf{x}_i \in c_k) = \frac{n_k}{n} \quad P(\mathbf{x}_i) = \frac{n_k}{n} \frac{1}{(2\pi\sigma^2)^{\frac{p}{2}}} \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{x}_i - \mu^k\|^2\right) \quad (92)$$

μ^i is the mean of the cluster to which \mathbf{x}_i belongs.

- **Log-likelihood:**

$$l(D) = \log \prod_i P(\mathbf{x}_i) \quad (93)$$

- **Log-likelihood detailed formula:**

$$l(D) = \sum_{k=1}^K \left[n_k \left(\log \frac{n_k}{n} - \frac{p}{2} \log(2\pi\sigma^2) \right) - \frac{1}{2\sigma^2} \sum_{\mathbf{x}_i \in X^{(k)}} \|\mathbf{x}_i - \mu^k\|^2 \right] \quad (94)$$

- **Statistical-based Indices**

- **BIC** [\uparrow]: The version from (Pelleg et al., 2000) was implemented.

$$\text{BIC} = l(D) - \frac{(p+1)K}{2} \log(n) \quad (95)$$

- **AIC** [\uparrow]: The version from (Pelleg et al., 2000) was implemented.

$$\text{AIC} = l(D) - (p+1)K \quad (96)$$

2.1.5. Negentropy-based indices:

- **Preliminary definitions:**

- **Assumption:** the concept of negentropy measures the average distance to normality of the clusters in the partition (Lago-Fernández and Corbacho, 2010a).

- **Indices based on negentropy**

- 245 – **Negentropy increment** ($\Delta J_B(\Omega, X)$) [\downarrow]:

$$\text{Negentropy increment} = \frac{1}{2} \sum_{k=1}^K (p_k \log |\Sigma_k|) - \sum_{k=1}^K p_k \log p_k \quad (97)$$

246 Where $p_k = n_k/n$ and Σ_k is the covariance matrix of the cluster k .

- **Negentropy increment corrected** ($\Delta J_U(\Omega, X)$) [\downarrow]: When the previous index is estimated with few samples produces an biased value. With a correction term this issue is supplied:

$$\text{Negentropy increment corrected} = \text{Negentropy increment} + \frac{1}{2} \sum_{k=1}^K p_k C(n_k, p) \quad (98)$$

$$C(n_k, p) = -p \log \frac{2}{n_k - 1} - \sum_{j=1}^p \psi\left(\frac{n_k - j}{2}\right) \quad (99)$$

247 Where $\psi(x) = \frac{d}{dx} (\ln \Gamma(x))$ is the digamma function.

- **Negentropy increment + variance** [\downarrow]: The variance of $\Delta J_U(\Omega, X)$ can be estimated as:

$$\sigma_s^2(\Delta J_U) \approx \frac{1}{4} \sum_{k=1}^K p_k^2 \sum_{j=1}^p \psi'\left(\frac{n_k - j}{2}\right) \quad (100)$$

Where $\psi'(x)$ is the first derivative of the digamma function.

$\Delta J_U(\Omega, X)$ does not take into account the variance and because of it is estimated on a finite sample, the estimation could satisfy $\Delta J(\Omega_1, X) < \Delta J(\Omega_2, X)$ for the true distribution of X , but using the finite sample of X the estimation could satisfy $\Delta J_U(\Omega_1, X) > \Delta J_U(\Omega_2, X)$. To solve this problem we include in the analysis the variance of $\Delta J_U(\Omega_1, X)$. Two partitions are considered equivalent if:

$$\Delta J_U(\Omega_2, X) + \sigma_s^2(\Delta J_U(\Omega_2, X)) > \Delta J_U(\Omega_1, X) - \sigma_s^2(\Delta J_U(\Omega_1, X)) \quad (101)$$

248 In such cases, the simplest (lower number of clusters) partition is selected.

249 2.1.6. Feature-Maximization-based indices:

250 • **Preliminary definitions:**

- 251 – **Assumption:** Feature maximization is an unbiased measure which can be used to estimate the quality of
 252 a classification whether it be supervised or unsupervised. In unsupervised classification (i.e. clustering),
 253 this measure exploits the properties (i.e. the features) of cluster associated data (Lamirel et al., 2016b).
- **Feature F-measure** $FF_k(f)$ **associated to cluster k :** is defined as the harmonic mean of the feature recall $FR_k(f)$ and of the feature predominance $FP_k(f)$ (Where subindex k refers to the feature f evaluated for

the cluster k).

$$FF_k(f) = \frac{2 \times FP_k(f) \times F2_k(f)}{FP_k(f) + FR_k(f)} \quad FP_k(f) = \frac{\sum_{x_i \in X^{(k)}} x_i^f}{\sum_{f' \in F_k} \sum_{x_i \in X^{(k)}} x_i^{f'}} \quad FR_k(f) = \frac{\sum_{x_i \in X^{(k)}} x_i^f}{\sum_{k=1}^K \sum_{x_i \in X^{(k)}} x_i^f}$$

Where x_i^f is the value of the feature f for the object i and F_k contains all the features associated with the cluster k . $FP_k(f)$ measures the ability of f to describe cluster k . $FR_k(f)$ measures the ability of f to discriminate k from other clusters.

– **Set S_k of relevant characteristic features of the cluster k :**

$$S_k = \{f \in F_k \mid FF_k(f) > \overline{FF}(f) \text{ and } FF_k(f) > \overline{FF}_D\} \quad (102)$$

Where $\overline{FF}(f)$ is the arithmetic mean of the F-measure for the feature f and \overline{FF}_D is the arithmetic mean of the $\overline{FF}(f)$ of all features in data.

$$\overline{FF}(f) = \sum_{k=1}^K FF_k(f) \quad \overline{FF}_D = \sum_{f=1}^F \overline{FF}(f) \quad (103)$$

Features with $FF_k(f) < \overline{FF}_D$ are discarded.

– **Set S of relevant characteristic features for the partition:**

$$S = \bigcup_{k=1}^K S_k \quad (104)$$

– **Constrast $G_k(f)$ of the feature f for a given cluster k :**

$$G_k(f) = \frac{FF_k(f)}{\overline{FF}(f)} \quad (105)$$

• **Indices based on feature-maximization**

– **PC [\uparrow]:** This index is based on the maximization of arithmetic mean of the contrast of selected features:

$$PC_K = \frac{1}{K} \sum_{k=1}^K \frac{1}{n_k} \sum_{f \in S_k} G_k(f) \quad (106)$$

– **EC [\uparrow]:** EC index is based on the maximization of the average weighted compromise between the contrast of active features and the inverted contrast of passive features:

$$EC_K = \frac{1}{K} \sum_{k=1}^K \left(\frac{\frac{|s_k|}{n_k} \sum_{f \in S_k} G_k(f) + \frac{|\overline{s}_k|}{n_k} \sum_{f \in \overline{S}_k} \frac{1}{G_k(f)}}{|s_k| + |\overline{s}_k|} \right) \quad (107)$$

Where $|s_k|$ and $|\overline{s_k}|$ are the number of selected and not selected features associated with the cluster k , respectively.

2.2. External indices

2.2.1. External indices based on a contingency table between the obtained partition and the ground partition

Preliminary definitions:

- **Obtained partition P_o :** Clustering resulted from the clustering algorithm.
- **Ground partition P_{GT} :** Predefined classes labels of the dataset.
- **Contingency table:** Matrix \mathbf{M} used to compare the partitions P_o and P_{GT} . P_o and P_{GT} are ways of partitioning S into disjointed subsets $U = \{U_1, U_2, \dots, U_R\}$ and $V = \{V_1, V_2, \dots, V_C\}$, respectively. In this way, \mathbf{M} is a $R \times C$ matrix and each element $n_{i,j}$ represents the number of elements that U_i and V_j have in common (Hubert and Arabie, 1985a).

Table 1: Contingency table

$U \setminus V$	V_1	V_2	...	V_C	Sums
U_1	$n_{1,1}$	$n_{1,2}$...	$n_{1,C}$	a_1
U_2	$n_{2,1}$	$n_{2,2}$...	$n_{2,C}$	a_2
...
U_R	$n_{R,1}$	$n_{R,2}$...	$n_{R,C}$	a_R
Sums	b_1	b_2	...	b_C	N

- **Number of pairs yy:** the number of pair of points that were grouped in the same cluster in P_o and belongs to the same cluster in P_{GT} .

$$yy = \frac{1}{2} \sum_{i=1}^R \sum_{j=1}^C n_{i,j}(n_{i,j} - 1) \quad (108)$$

- **Number of pairs yn:** the number of pair of points that were grouped in the same cluster in P_o and in different clusters in P_{GT} .

$$yn = \frac{1}{2} \left(\sum_{j=1}^R a_j^2 - \sum_{i=1}^R \sum_{j=1}^C n_{i,j}^2 \right) \quad (109)$$

- **Number of pairs ny:** the number of pair of points that were grouped in different clusters in P_o and in the same cluster in P_{GT} .

$$ny = \frac{1}{2} \left(\sum_{j=1}^C b_j^2 - \sum_{i=1}^R \sum_{j=1}^C n_{i,j}^2 \right) \quad (110)$$

- **Number of pairs nn:** the number of pair of points that were grouped in different clusters in P_o and in different

clusters in P_{GT} .

$$nn = \frac{1}{2} \left(n^2 + \sum_{i=1}^R \sum_{j=1}^C n_{i,j}^2 - \left(\sum_{i=1}^R a_i^2 + \sum_{j=1}^C b_j^2 \right) \right) \quad (111)$$

- **Total number of pairs of points n_T :**

$$n_T = \frac{n(n+1)}{2} = yy + yn + ny + nn \quad (112)$$

- **Indicator variable X_a associated to the partition P_a :** Defined as a binary random variable that indicates if a pair of points i and j are in the same cluster in P_a . This variable takes the value of 1 if a pair of points are grouped in the same cluster and the value of 0 otherwise.

- **Mean μ_{X_a} of the indicator variable X_a :**

$$\mu_{X_a} = \frac{1}{n_T} \sum_{i < j} X_a(i, j) \quad (113)$$

- **Standard deviation σ_{X_a} of the indicator variable X_a :**

$$\sigma_{X_a} = \frac{1}{n_T} \sum_{i < j} [X_a(i, j)^2 - \mu_{X_a}] \quad (114)$$

- **Indicator variables for P_o and P_{GT} :**

$$\mu_{X_o} = \frac{1}{n_T} \sum_{i < j} X_o(i, j) = \frac{yy + yn}{n_T} \quad \sigma_{X_o}^2 = \frac{1}{n_T} \sum_{i < j} [X_o(i, j)^2 - \mu_{X_o}^2] = \frac{yy + yn}{n_T} - \left(\frac{yy + yn}{n_T} \right)^2 \quad (115)$$

$$\mu_{X_{GT}} = \frac{1}{n_T} \sum_{i < j} X_{GT}(i, j) = \frac{yy + ny}{n_T} \quad \sigma_{X_{GT}}^2 = \frac{1}{n_T} \sum_{i < j} [X_{GT}(i, j)^2 - \mu_{X_{GT}}^2] = \frac{yy + ny}{n_T} - \left(\frac{yy + ny}{n_T} \right)^2 \quad (116)$$

- **Precision:** is defined as the ratio of yy (the number of pairs of points that were grouped in the same cluster in P_o and belongs to the same cluster in P_{GT}) to $yy + ny$ (the number of pairs of points that in P_{GT} belongs to the same cluster). In other words, precision is the portion of point pairs rightly grouped according to P_{GT} .

$$\text{Precision} = \frac{yy}{yy + ny} \quad (117)$$

- **Recall:** is defined as the ratio of yy to $yy + yn$ (the number of pairs of points that were grouped in the same cluster in P_o). In other words, from all pairs that were grouped in the same cluster in P_o ($yy + yn$) what is the

proportion that belongs to the same cluster in P_{GT} .

$$\text{Recall} = \frac{yy}{yy + yn} \quad (118)$$

Indices based on the contingency table

- **F-measure:** This index is based on the well-known metric of the supervised learning, F-score. The F-measure is the harmonic mean of the precision and recall coefficients (Larsen and Aone, 1999):

$$\text{F-Measure} = \frac{2}{\frac{1}{\text{Precision}} + \frac{1}{\text{Recall}}} = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \quad (119)$$

- **Folkes-Mallows:** This index is defined as (Fowlkes and Mallows, 1983):

$$\text{Folkes-Mallows} = \frac{yy}{\sqrt{(yy + yn) \times (yy + ny)}} = \sqrt{\text{Precision} \times \text{Recall}} \quad (120)$$

- **Hubert Γ :** This index is based on the correlation between the indicator variables of P_o and P_{GT} . This index can be defined in terms of the confusion matrix (Hubert, 1979).

$$\text{Hubert } \Gamma = \frac{n_T \times yy - (yy + yn)(yy + ny)}{\sqrt{(yy + yn)(yy + ny)(nn + yn)(nn + ny)}} \quad (121)$$

- **Jaccard:** This index is defined as (Rousseeuw and Kaufman, 1990):

$$\text{Jaccard} = \frac{yy}{yy + yn + ny} \quad (122)$$

- **Kulczynski:** This index is defined as (Murguía and Villaseñor, 2003):

$$\text{Kulczynski} = \frac{1}{2} \left(\frac{yy}{yy + ny} + \frac{yy}{yy + yn} \right) = \frac{1}{2} (\text{Precision} + \text{Recall}) \quad (123)$$

- **McNemar:** It is an adaptation of the non-parametric test of McNemar for the comparison of frequencies between two paired samples (Eliasziw and Donner, 1991).

$$\text{McNemar} = \frac{nn - ny}{\sqrt{nn + yy}} \quad (124)$$

- **Phi:** Phi index is based on the correlation between two dichotomic variables (Cramér, 2016).

$$\text{Phi} = \frac{yy \times nn - yn \times ny}{\sqrt{(yy + yn)(yy + ny)(nn + yn)(nn + ny)}} \quad (125)$$

306

- **Rand:** (Rand, 1971)

$$\text{Rand} = \frac{yy + nn}{n_T} \quad (126)$$

- **Adjusted Rand index (ARI):** Because Rand have a score different to 0.0 when the label assignments were random (uniform), an adjusted version of the rand index (ARI) was proposed. ARI is the difference of the Rand index and its expected value under the null hypothesis, which assumes a generalized hypergeometric distribution. To see more details refer to (Hubert and Arabie, 1985b).

$$\text{ARI} = \frac{\text{Rand} - E[\text{Rand}]}{\max(\text{Rand}) - E[\text{Rand}]} \quad (127)$$

$$E[\text{Rand}] = \frac{\sum_i \binom{a_i}{2} \sum_j \binom{b_j}{2}}{\binom{N}{2}} \quad \max(\text{Rand}) = \frac{1}{2} \frac{\sum_i \binom{a_i}{2} + \sum_j \binom{b_j}{2}}{\binom{N}{2}} \quad (128)$$

307

- **Rogers-Tanimoto:**

$$\text{Rogers-Tanimoto} = \frac{yy + nn}{yy + nn + 2(yn + ny)} \quad (129)$$

308

- **Russel-Rao:** is defined as the proportion of concordances between the two partitions (Rao, 1948):

$$\text{Russel-Rao} = \frac{yy}{\sqrt{n_T}} \quad (130)$$

309

- **Sokal-Sneath:** Sokal-Sneath is defined as the proportion of concordances between the two partitions:

$$\text{Sokal-Sneath} = \frac{yy}{yy + 2(yn + ny)} \quad (131)$$

310 *2.2.2. External indices based on information theory and intuitive concepts:*

311

Preliminary definitions:

312

- **Entropy $H(A)$ of partition P_a :** P_a is a way of partitioning S into disjointed subsets $A = \{A_1, A_2, \dots, A_M\}$ and the entropy of P_a is the amount of uncertainty.

313

$$H(A) = - \sum_{i=1}^M P(i) \log(P(i)) \quad (132)$$

314

Where $P(i) = |A_i|/N$ is the probability that an element of S falls into cluster A_i .

315

- **Conditional entropy $H(V/U)$:** is the conditional entropy of the C classes given the R cluster assignments.

$$H(V/U) = - \sum_{k=1}^R \frac{n_k}{n} \sum_{c=1}^C \frac{n_{c,k}}{n_k} \log \left(\frac{n_{c,k}}{n_k} \right) \quad (133)$$

Where $n_{c,k}$ is the number of elements of the class c in the cluster k and n_k is the number of elements of the cluster k .

- **Probability** $P(i, j)$: is the probability that an element falls into a cluster A_i of a partition P_a and a cluster B_j of a partition P_b .

$$P(i, j) = \frac{|A_i \cap B_j|}{n} \quad (134)$$

Indices based on information theory

- **Conditional Entropy:** (Shannon, 2001)

$$\text{Conditional Entropy} = H(V/U) \quad (135)$$

- **Homogeneity:** Based on that a cluster must contain only elements of a single class (Rosenberg and Hirschberg, 2007).

$$\text{Homogeneity} = 1 - \frac{H(V/U)}{H(V)} \quad (136)$$

- **Completeness:** This index is symmetrical to homogeneity. This metric assumes that a class must be represented only by a single cluster (Rosenberg and Hirschberg, 2007).

$$\text{Completeness} = 1 - \frac{H(U/V)}{H(U)} \quad (137)$$

- **V-Measure:** defined as the harmonic mean of homogeneity and completeness (Rosenberg and Hirschberg, 2007).

$$\text{V-Measure} = 2 \times \frac{\text{Homogeneity} \times \text{Completeness}}{\text{Homogeneity} + \text{Completeness}} \quad (138)$$

- **Mutual information:** (Strehl and Ghosh, 2002)

$$\text{MI}(U, V) = \sum_{i=1}^R \sum_{j=1}^C P(i, j) \log \left(\frac{P(i, j)}{P(i)P(j)} \right) \quad (139)$$

- **Normalized mutual information:** (Strehl and Ghosh, 2002)

$$\text{NMI}(U, V) = \frac{\text{MI}(U, V)}{\sqrt{H(U)H(V)}} \quad (140)$$

- **Adjusted mutual information:** Proposed by (Vinh et al., 2009)

$$\text{AMI}(U, V) = \frac{\text{MI} - E[\text{MI}]}{\max[H(U), H(V)] - E[\text{MI}]} \quad (141)$$

Where $E[MI]$ is the expected mutual information and is defined by:

$$E[MI] = \sum_{i=1}^R \sum_{j=1}^C \sum_{n_{i,j}=\max(a_i+b_j-n,0)}^{\min(a_i,b_j)} \frac{n_{i,j}}{n} \log \left(\frac{n \times n_{i,j}}{a_i \times b_j} \right) \frac{a_i! b_j! (n - a_i)! (n - b_j)!}{n! n_{i,j}! (a_i - n_{i,j})! (b_j - n_{i,j})! (n - a_i - b_j - n_{i,j})!} \quad (142)$$

3. Experimental setup

This section explains the comparative study performed to analyze the clustering validation indices previously introduced. First, a recompilation of 26 synthetic and 23 real datasets have been chosen as benchmark datasets, aiming to consider the more general problems in cluster analysis with quantitative data. Then, six clustering algorithms have been used to create the candidate partitions: K-means (Jain, 2010), DBSCAN (Ester et al., 1996), Expectation Maximization (Dempster et al., 1977), Birch (Zhang et al., 1996), Sting (Wang et al., 1997), and Lamda (Bedoya et al., 2014). A brief explanation of the algorithms is presented in section 3.2. Each algorithm has at least one input parameter that the user must specify. A combination of parameters was performed per algorithm, generating the multiple partitions. Additional to this, three standardizations of features were used to analyze the effect of the scale in clustering validation. In total, 178737 partitions were generated. Both external and internal clustering validity indices were computed for each candidate partition. However, indices based on concordance and discordance were only estimated for datasets that contains less than 500 samples due to their high computational requirements. Then, a comparative methodology was carried out, analyzing the matching among internal and external indices. Following the ideas of (Gurrutxaga et al., 2011; Arbelaiz et al., 2013), external indices work as reference metrics, and internal indices should select as best partitions those that the external indices have selected.

3.1. Datasets

Datasets from several studies were recollected and those that were used in major number of times were selected (Arbelaiz et al., 2013; Wani and Riyaz, 2016b; Zhao and Fränti, 2014; Liu et al., 2013; Gurrutxaga et al., 2011). We tried to include the more general cases that can be found in clustering literature of quantitative data. The characteristics of these synthetic and real datasets are shown in table 2. Most of the synthetic datasets are two-dimensional datasets, therefore is possible to observe graphically the results. Real datasets were mainly chosen from the work of (Arbelaiz et al., 2013) due to one of the main goal of the present work is to extend that analysis.

3.2. Clustering algorithms

The first five algorithms were selected based on a clustering algorithm categorization done by (Fahad et al., 2014), who analyzed the ability of clustering algorithms to handle big data and divided the algorithms into five categories (partitional, hierarchical, density-based, grid-based, and model-based). A representative algorithm was chosen per each category. In this work, the same algorithms of the partitional, hierarchical and model-based categories were kept, but the density-based and grid-based algorithms were replaced by DBSCAN and STING. This decision was taken due