

Hashtag Discernability – Competitiveness Study of Graph Spectral and Other Clustering Methods

Bartłomiej Starosta, Mieczysław A. Kłopotek, Sławomir T. Wierzchoń, Dariusz Czernski
 Institute of Computer Science, Polish Academy of Sciences
 ul. Jana Kazimierza 5, 01-248 Warsaw, Poland
 Email: b.starosta,kłopotek,stw,d.czerski@ipipan.waw.pl

Abstract—Spectral clustering methods are claimed to possess ability to represent clusters of diverse shapes, densities etc. They constitute an approximation to graph cuts of various types (plain cuts, normalized cuts, ratio cuts). They are applicable to unweighted and weighted similarity graphs. We perform an evaluation of these capabilities for clustering tasks of increasing complexity.

I. INTRODUCTION

DOCUMENT clustering (or text clustering) has a multitude of applications, including topic extraction, fast information retrieval, filtering, authorship discovery, topic drift detection in news streams and social media, automatic document organization etc. ([1], [2], [3], [4])

Two clustering methods are of particular interest in this area, the Graph Spectral Clustering (GSC) and spherical k -means.

Graph Spectral Clustering methods [1] are generally praised for possessing ability to represent clusters of diverse shapes, densities etc. They constitute an approximation to graph cuts of various types (plain cuts, normalized cuts, ratio cuts). They are applicable to unweighted and weighted similarity graphs.

Spherical k -means algorithm [5] is a variant of k -means algorithm that measures similarity of documents based on their cosine similarity, that is quite popular in the domain of text analysis (e.g. for search engines).

In this paper we pose the question: If the grouping method correctly groups certain datasets, can we expect that a combination of these datasets will also be correctly clustered? We will examine the following problem in more detail. Assume that a clustering method can cluster correctly documents from categories $[A, B]$, $[B, C]$, and $[C, A]$. Can we expect the algorithm to cluster correctly data from the mixed set $[A, B, C]$? Let us illustrate this with three datasets, tweets, marked with (single) tags 'lolinginlove', 'tejran', 'anjisalvacion'.

We used standard Python implementation of spectral clustering from scikit-learn library.¹ The affinity matrix was constructed from a k -nearest neighbors connectivity matrix, with the default value of $k = 10$.

In one of the experiments the clustering illustrated in Fig. 1 was obtained for the hashtags 'lolinginlove', 'tejran'. For the hashtags 'tejran', 'anjisalvacion' the nearest neighbor spectral clustering achieves the best clustering agreement visible in Fig. 2. For the hashtags 'lolinginlove', 'anjisalvacion', the

¹Consult <https://scikit-learn.org/stable/modules/generated/sklearn.cluster.SpectralClustering.html> for details.

	T\P	0	1	
lolinginlove:	0	1258	0	1258
tejran:	1	8	337	345
		1266	337	1603

F-score: 0.990046

Fig. 1. Spectral clustering with affinity "nearest neighbors" example 1; row labels - "true" clusters, column labels - clustering result

	T\P	0	1	
tejran:	0	324	21	345
anjisalvacion:	1	5	727	732
		329	748	1077

F-score: 0.968385

Fig. 2. Spectral clustering with affinity "nearest neighbors" example 2; row labels - "true" clusters, column labels - clustering result

nearest neighbor spectral clustering achieves the clustering agreement visible in Fig. 3.

So, for each pair of the three hashtags we see a very good agreement of clusterings with the target (hashtags). If we look at the hashtags ['lolinginlove', 'tejran', 'anjisalvacion'], we get clustering agreement visible in Fig. 4. We see that more errors are committed here than for each pair of hashtags presented in Figs. 1, 2 and 3, though the increase does not seem to be large in absolute numbers. We will return to this issue in the next section.

Here and in further sections, F-score is computed as follows. We assume that the clustering is to predict the hashtag. The "true" hashtag is identified as the majority hashtag in the cluster. For a given hashtag H we proceed as follows. True positives (TP) are those cases when cluster membership agrees

	T\P	0	1	
lolinginlove:	0	1258	0	1258
anjisalvacion:	1	0	732	732
		1258	732	1990

F-score: 1.000000

Fig. 3. Spectral clustering with affinity "nearest neighbors" example 3; row labels - "true" clusters, column labels - clustering result

T\P	0	1	2		
lolinginlove:	0	1258	0	0	1258
tejran:	1	7	314	24	345
anjisalvacion:	2	0	5	727	732
		1265	319	751	2335

F-score: 0.970334

Fig. 4. Spectral clustering with affinity "nearest neighbors" example 4; row labels - "true" clusters, column labels - clustering result

with this hashtag. False positives (FP) are the cases which belong to the cluster for which the hashtag H is the true hashtag, but the hashtag for the given document is different from H. True negatives (TN) are the cases which belong to the cluster for which the hashtag H is not the true hashtag, and the hashtag for the given document is different from H. False negatives (FN) are the cases which belong to the cluster for which the hashtag H is not the true hashtag, but the hashtag for the given document is the hashtag H. Computation of precision and recall follows the standard pattern and the F-score is computed for each hashtag separately, and then the average is taken as the F-score for the clustering.

In this paper we study the extent to which this behaviour extends to larger number of clusters. This study is a starting point for a future revision of the studied clustering algorithms.

II. CONCEPTUAL CONSIDERATIONS

Despite the example shown above, it is not entirely obvious that given a grouping method that allows to correctly group documents from the categories $[A, B]$, $[B, C]$, $[C, A]$, we can expect that the algorithm will correctly group data from the mixed set $[A, B, C]$.

If the sets $A \cup B$, $B \cup C$ and $C \cup A$ have block diagonal document similarity matrices (after proper reordering the documents), and the blocks are actually within A, B, C then in fact the $[A, B, C]$ similarity matrix will be block diagonal too so that GSC algorithm will cluster A, B, C correctly. This can be seen immediately by inspection of block matrix structure, i.e.

$$S_{A,B} = \begin{bmatrix} S_{A,A} & 0 \\ 0 & S_{B,B} \end{bmatrix} \quad S_{B,C} = \begin{bmatrix} S_{B,B} & 0 \\ 0 & S_{C,C} \end{bmatrix}$$

$$S_{A,C} = \begin{bmatrix} S_{A,A} & 0 \\ 0 & S_{C,C} \end{bmatrix}$$

implies

$$S_{A,B,C} = \begin{bmatrix} S_{A,A} & 0 & 0 \\ 0 & S_{B,B} & 0 \\ 0 & 0 & S_{C,C} \end{bmatrix}$$

Recall that combinatorial Laplacian is computed as $L = D - S$, where S is the similarity matrix and D is the diagonal matrix with elements being sums of corresponding rows of S . Hence

$$L_{A,B} = \begin{bmatrix} L_{A,A} & 0 \\ 0 & L_{B,B} \end{bmatrix}, \text{ etc.}$$

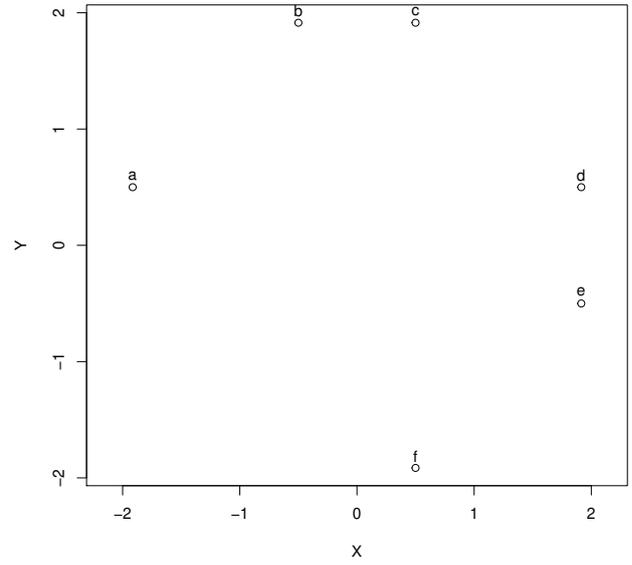


Fig. 5. Visualization of datapoints used to illustrate the increasing clustering problem for k -means

and

$$L_{A,B,C} = \begin{bmatrix} L_{A,A} & 0 & 0 \\ 0 & L_{B,B} & 0 \\ 0 & 0 & L_{C,C} \end{bmatrix}$$

Eigenvalues of $L_{A,B}$, $L_{B,C}$, $L_{A,C}$ will become eigenvalues of $L_{A,B,C}$ with corresponding eigenvectors being only extended with zeros appropriately. So theoretically it should be easy to separate the sets A, B, C based on eigenvectors of $L_{A,B,C}$. However, this enthusiasm needs to be mitigated because such a pure block structure rarely occurs, see our example Fig. 1, Fig. 2, Fig. 3, so the "noise" is inherited in sets with more hashtags as visible in Fig. 4. But there are also further concerns. Spectral clustering is based on lowest eigenvalue eigenvectors of respective Laplacians. But as shown in [6], the two lowest eigenvectors of $L_{A,B}$, $L_{B,C}$, $L_{A,C}$ do not need to be lowest three eigenvectors of $L_{A,B,C}$. For higher number of clusters, the situation may be more complex.

If the dataset $A \cup B \cup C$ is well separated in the sense of k -means algorithm, so that a clustering with k -means will yield A, B, C as clusters, then its application to $A \cup B$, $B \cup C$ or $C \cup A$ will also return correct pairs of clusters. But this is not necessarily true for k -means in the reverse direction. Well-separatedness of $A \cup B$, $B \cup C$ and $C \cup A$ does not imply well-separatedness of $A \cup B \cup C$. Let us illustrate this point with a bit artificial example. Consider the datapoints $\mathbf{a} = (-(0.5 + \sqrt{2}), 0.5)$, $\mathbf{b} = (-0.5, 0.5 + \sqrt{2})$, $\mathbf{c} = (0.5, 0.5 + \sqrt{2})$, $\mathbf{d} = (0.5 + \sqrt{2}, 0.5)$, $\mathbf{e} = (0.5 + \sqrt{2}, -0.5)$, $\mathbf{f} = (0.5, -(0.5 + \sqrt{2}))$, see Fig. 5 for visualization. Consider "hashtags" with their "documents": $A = \{\mathbf{a}, \mathbf{b}\}$, $B = \{\mathbf{c}, \mathbf{d}\}$, $C = \{\mathbf{e}, \mathbf{f}\}$. Clustering with k -means of $A \cup C$ into two clusters

TABLE I
TWT.10 DATA SET - HASHTAGS AND CARDINALITIES OF THE SET OF RELATED TWEETS USED IN THE EXPERIMENTS

No.	hashtag	count
0	90dayfiance	316
1	tejrán	345
2	ukraine	352
3	tejasswiprakash	372
4	nowplaying	439
5	anjisalvacion	732
6	puredoctrinesofchrist	831
7	1	1105
8	lolinginlove	1258
9	bbnaira	1405

will yield A, C , similarly any two hashtag combinations. But clustering with k -means of $A \cup B \cup C$ will yield three clusters $\{\mathbf{a}\}, \{\mathbf{b}, \mathbf{c}\}, \{\mathbf{d}, \mathbf{e}, \mathbf{f}\}$, not A, C, E .

In all these cases, if some noise is added to fuzzify the well-separatedness, the noise can be more destructive for the set A, B, C than for any of the three mentioned subsets – this affects GSC as well as k -means clustering. This is easily imagined by considering k -means algorithm. The cluster center of A when clustering fuzzified A and B may lie in a different position than when clustering fuzzified A and C .

This behavior will be subsequently illustrated by a series of experiments.

III. DATA

We used tweets retrieved from the stream endpoint of Twitter API (a random sample of about 1% of English tweets), collected by one of the Authors for the time period from mid September 2019 till end of November 2022. From this set we extracted the subset TWT.10 used in experiments. It is a collection of top thread tweets related to hashtags listed in Table I. While selecting the data, we imposed the restriction that the tweets had to have one single hashtag (which we treated as an indication of being devoted to a single theme).

IV. METHODS

We study two standard versions of Graph Spectral Clustering, available from scikit-learn, and the 6 versions of spherical k -means and 6 versions of our proprietary so-called K -embedding based clustering algorithm.

More precisely the clustering experiments were performed with popular Python libraries: numpy [7], scipy [8], scikit-learn [9] and soyclustering [10] which is an implementation of spherical k -means [11]. In particular, we used

- 1) `SpectralClustering` class from scikit-learn with two distinct settings of the `affinity` parameter: `precomputed` (affinity from similarity matrix) and `nearest_neighbors` (affinity from graph of nearest neighbors) - as a representative of the spectral clustering, and
- 2) `SphericalKMeans` class from soyclustering with the following combinations of (`init`, `sparsity`) parameter pairs (the mentioned 6 versions, short names given for reference): `"sc.n"`: ('similar_cut', None), `"sc.sc"`:

(`'similar_cut'`, `'sculley'`), `"sc.md"`: ('similar_cut', 'minimum_df'), `"k++.n"`: ('k-means++', None), `"k++.sc"`: ('k-means++', 'sculley'), `"k++.md"`: ('k-means++', 'minimum_df'), and

- 3) K -embedding clustering (our implementation, exploiting spherical k -means – see subsection IV-C). Same combinations of parameter pairs (versions) were used as for `SphericalKMeans` above. The following numbers of eigenvectors were tried: $r = 12$ and higher.

The advantages and disadvantages of these methods are briefly discussed below.

A. Spectral analysis

In fact spectral clustering algorithms constitute a large family, see e.g. [12], [13], [14], which have numerous desirable properties (like detection of clusters with various shapes, applicability to high dimensional datasets, capability to handle categorical variables), yet they suffer from various shortcomings, common to other sets of algorithms, including multiple possibilities of representation of the same dataset, producing results in a space different from the space of original problem, curse of dimensionality, etc. These shortcomings are particularly grievous under large and sparse data set scenario, like in Twitter data.

Let us briefly recall the typical spectral clustering algorithm in order to make it understandable, how distant the clustering may be from the applicator's comprehension [12]. The first step consists in creating a similarity matrix of objects (in case of documents based on tf, tfidf, in unigram or n-gram versions, or some transformer based embeddings are the options – consult e.g. [15] for details), then mixing them in case of multiple views available. The second step is to calculate a Laplacian matrix. There are at least three variants to use: combinatorial, normalized, and random-walk Laplacian, [12]. But other options are also possible, like: some kernel-based versions, non-backtracking matrix [16], degree-corrected versions of the modularity matrix [17] or the Bethe-Hessian matrix [18]. Then computing eigenvectors and eigenvalues, eigenvector smoothing (to remove noise and/or achieve robustness against outliers) choice of eigenvectors, and finally clustering in the space of selected eigenvectors (via e.g. k -means). The procedure may be more complex, e.g. one may add loops back to preceding steps based on feedback from quality analysis, like degree of deviation from block-structure of the Laplacian.

From this diversified set we chose the two mentioned implementations available from scikit-learn.

B. Spherical k -means

Spherical k -means was developed in [5] by observing that the squared Euclidean distance between two vectors, $\|\mathbf{x}_i - \mathbf{x}_j\|^2 = \|\mathbf{x}_i\|^2 - 2\mathbf{x}_i^T \mathbf{x}_j + \|\mathbf{x}_j\|^2$, in case of normalized vectors reduces to

$$\|\mathbf{x}_i - \mathbf{x}_j\|^2 = 2(1 - \mathbf{x}_i^T \mathbf{x}_j), \quad (1)$$

and $\mathbf{x}_i^T \mathbf{x}_j = \cos \angle(\mathbf{x}_i, \mathbf{x}_j)$. This makes it very efficient in case of sparse vectors, a typical representation of text documents. Such a variant of k -means suffers dependence on initialization, thus further improvements are proposed, e.g. [19], [20], [21] and [22].

C. K -embedding

K -embedding has the following underlying idea. Let us think for a moment about a particular embedding of the nodes of the graph, based on [23]. Let A be a matrix of the form:

$$A = \mathbf{1}\mathbf{1}^T - I - S, \quad (2)$$

where S stands for an affinity matrix, I is the identity matrix, and $\mathbf{1}$ is the (column) vector consisting of ones, both of appropriate dimensions. (Note that here we have to assume that the diagonal of S consists of zeros). Let K be the matrix of the (double centered) form [24]:

$$K = -\frac{1}{2}\left(I - \frac{1}{n}\mathbf{1}\mathbf{1}^T\right)A\left(I - \frac{1}{n}\mathbf{1}\mathbf{1}^T\right), \quad (3)$$

with $n \times n$ being the dimension of S . $\mathbf{1}$ is an eigenvector of K , with the corresponding eigenvalue equal to 0. All the other eigenvectors must be orthogonal to it as K is real and symmetric, so for any other eigenvector \mathbf{v} of K we have: $\mathbf{1}^T \mathbf{v} = 0$.

Let Λ be the diagonal matrix of eigenvalues of K , and V the matrix where columns are corresponding (unit length) eigenvectors of K . Then $K = V\Lambda V^T$. Let $\mathbf{z}_i = \Lambda^{1/2}V_i^T$, where V_i stands for i -th row of V . Let $\mathbf{z}_i, \mathbf{z}_\ell$ be the embeddings of the nodes i, ℓ , resp. This embedding shall be called K -embedding. Then

$$\|\mathbf{z}_i - \mathbf{z}_\ell\|^2 = 1 - S_{i\ell} \quad (4)$$

for $i \neq \ell$. Hence upon performing k -means clustering in this space we *de facto* try to maximize the sum of similarities within a cluster. Note that $K = V\Lambda V^T$ may be quite well approximated if we drop from Λ low eigenvalues and from V their corresponding eigenvectors (which we do in our experiments).

V. EVALUATION

For each of the algorithms we perform the following tests. For each pair of datasets associated with two hashtags from Table I (45 pairs in all) the clustering will be performed by each of the mentioned algorithms 10 times (due to stochastic nature of these algorithms) and the average F-score will be computed. Ten pairs with the highest average F-scores will be taken for the next phase. Now datasets associated with 3 hashtags will be created out of these selected pairs plus each of the hashtags not present in the selected pairs. This process is continued till all 10 hashtags are exhausted. In figures, the average value of F over all computations with the given hashtag cardinality is presented plus the average of the top 10 groups of hashtags. The results are summarized in Figs. 6–19.

The next experiment was to compare the F-score obtained by a given set of hashtags, considered in the preceding experiment, and its subsets obtained by removing one of

Average F-score – blue – over all hashtag sets, green – 10 top values

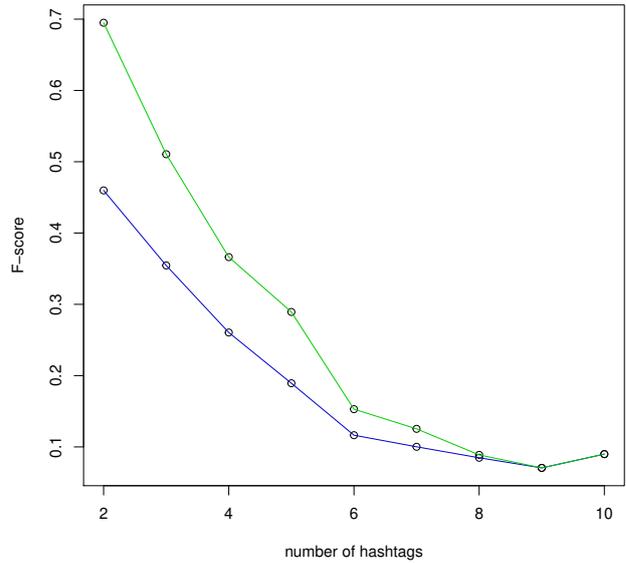


Fig. 6. F-scores for various numbers of hashtags; spectral clustering with affinity nearest_neighbors

TABLE II
CORRELATION BETWEEN THE F-SCORE OF A GIVEN GROUP OF HASHTAGS AND THEIR SUBGROUPS OF CARDINALITY LOWER BY ONE.

algorithm	pearson	p.val	spearman	p.val
spectral nearest_neighbors	0.7745	0	0.8358	0
spectral precomputed	0.7374	0	0.7437	0
spherical sc.md	0.7036	0	0.7711	0
spherical sc.sc	0.8306	0	0.8538	0
spherical k++.n	0.7647	0	0.8189	0
spherical sc.n	0.7778	0	0.8167	0
spherical k++.md	0.7796	0	0.8129	0
spherical k++.sc	0.8099	0	0.8502	0
K-embedding.12plus sc.md	0.6057	0	0.6041	0
K-embedding.12plus sc.sc	0.6948	0	0.6678	0
K-embedding.12plus k++.n	0.7975	0	0.8294	0
K-embedding.12plus sc.n	0.6901	0	0.7113	0
K-embedding.12plus k++.md	0.7976	0	0.8483	0
K-embedding.12plus k++.sc	0.7924	0	0.8460	0

the hashtags. For example, we considered the F-score of ['lolinginlove', 'tejran', 'anjisalvacion'] and the average of F-scores for the subsets ['lolinginlove', 'tejran'], ['lolinginlove', 'anjisalvacion'] and ['tejran', 'anjisalvacion']. We computed the Spearman and Pearson correlations for such pairs (F-score of a set of hashtags and average of its subsets) and presented the results in Table II for each of the analysed clustering algorithms. We have created also a more detailed view for one of the algorithms: spectral clustering with affinity nearest neighbors. Fig. 20 presents the histogram of differences between the average F-score of subgroups and the F-score of the group. Fig. 21 presents the relation between the average F-score of subgroups and the F-score of the group as a scatterplot.

Average F-score – blue – over all hashtag sets, green – 10 top values

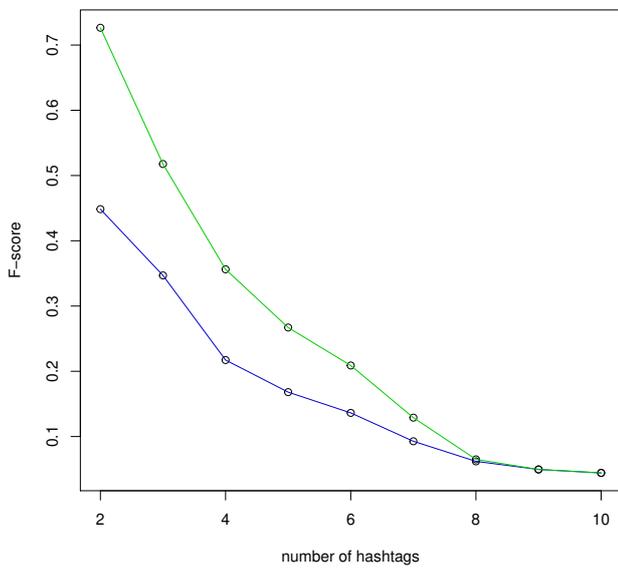


Fig. 7. F-scores for various numbers of hashtags; spectral clustering with affinity precomputed

Average F-score – blue – over all hashtag sets, green – 10 top values

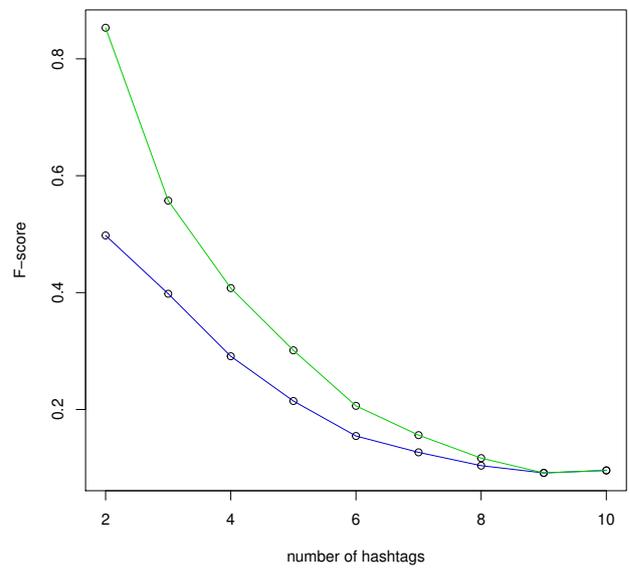


Fig. 9. F-scores for various numbers of hashtags; spherical k -means clustering with sc.sc configuration

Average F-score – blue – over all hashtag sets, green – 10 top values

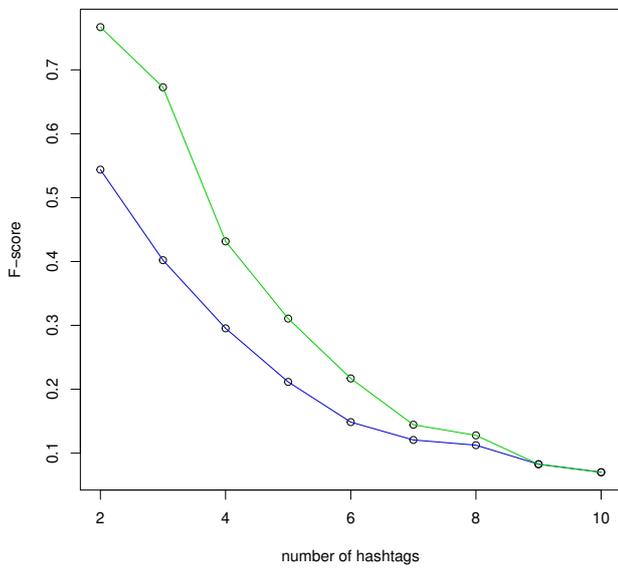


Fig. 8. F-scores for various numbers of hashtags; spherical k -means clustering with sc.md configuration

Average F-score – blue – over all hashtag sets, green – 10 top values

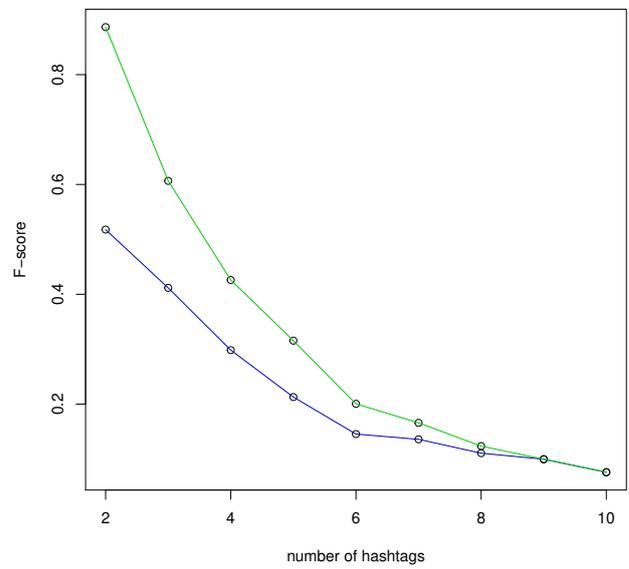


Fig. 10. F-scores for various numbers of hashtags; spherical k -means clustering with sc.n configuration

Average F-score – blue – over all hashtag sets, green – 10 top values

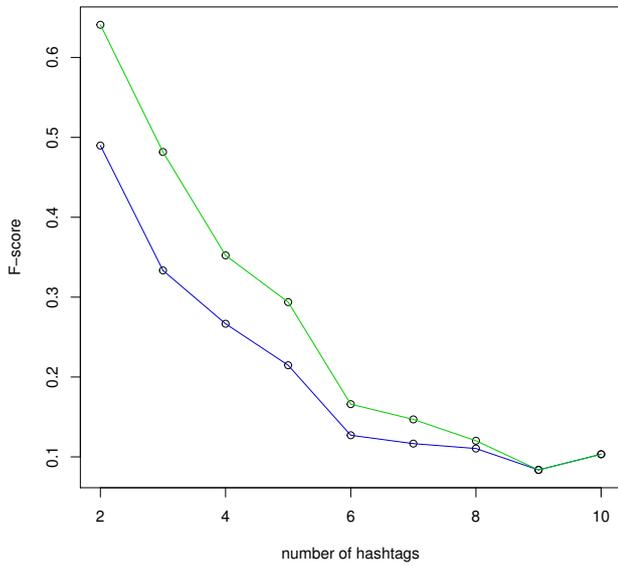


Fig. 11. F-scores for various numbers of hashtags; spherical k -means clustering with k++.md configuration

Average F-score – blue – over all hashtag sets, green – 10 top values

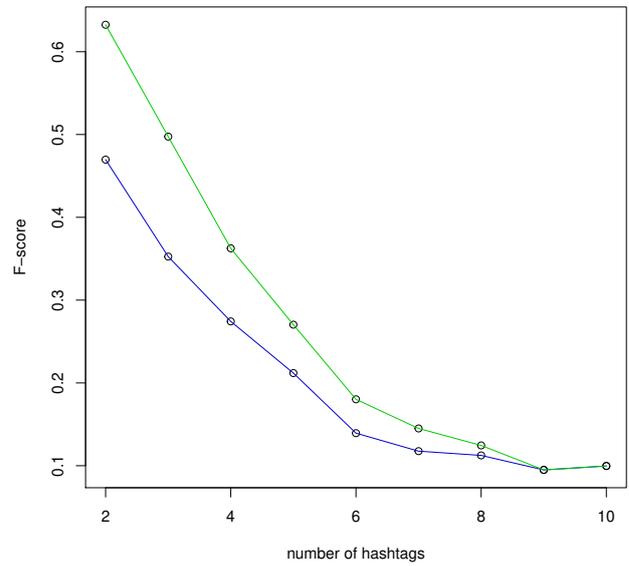


Fig. 13. F-scores for various numbers of hashtags; spherical k -means clustering with k++.sc configuration

Average F-score – blue – over all hashtag sets, green – 10 top values

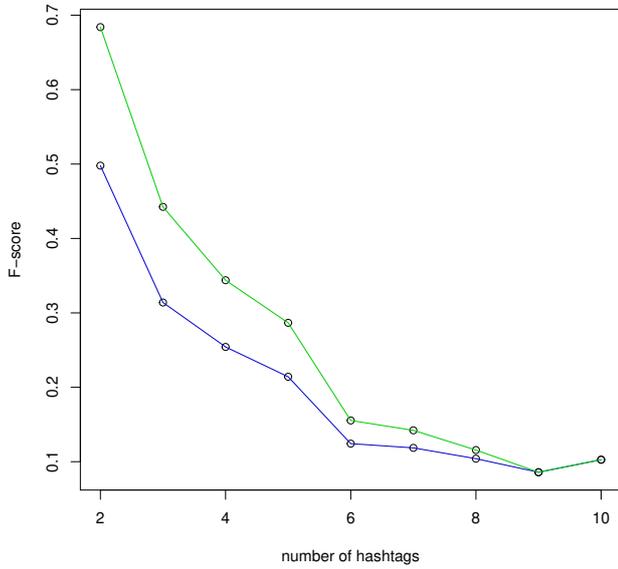


Fig. 12. F-scores for various numbers of hashtags; spherical k -means clustering with k++.n configuration

Average F-score – blue – over all hashtag sets, green – 10 top values

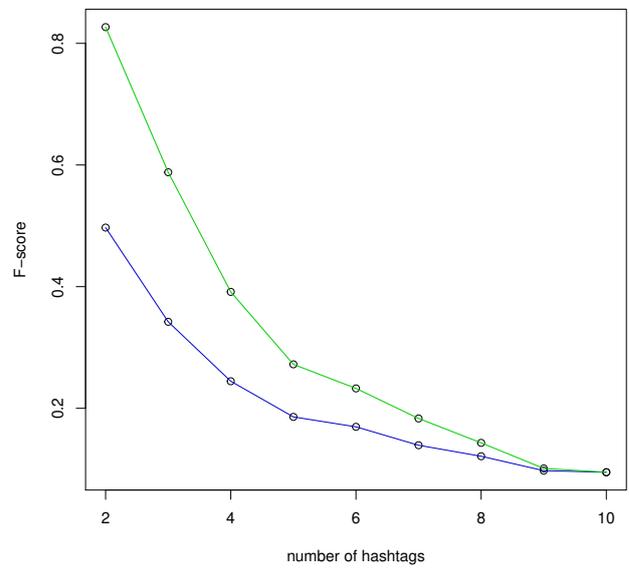


Fig. 14. F-scores for various numbers of hashtags; K-embedding based clustering with sc.md configuration

Average F-score – blue – over all hashtag sets, green – 10 top values

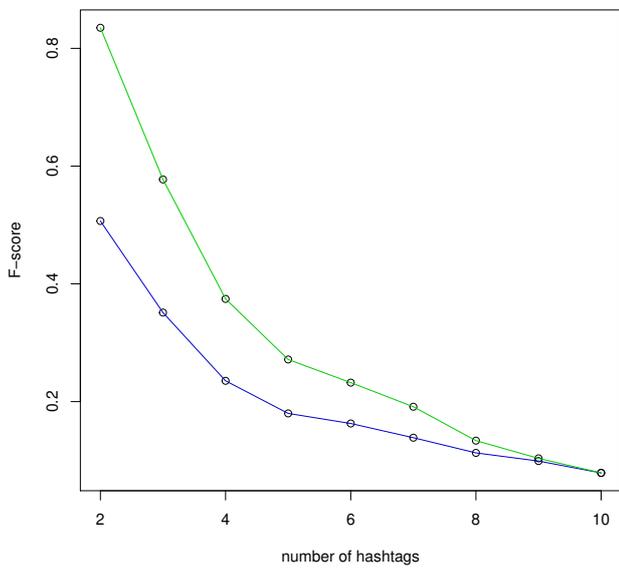


Fig. 15. F-scores for various numbers of hashtags; K-embedding based clustering with sc.sc configuration

Average F-score – blue – over all hashtag sets, green – 10 top values

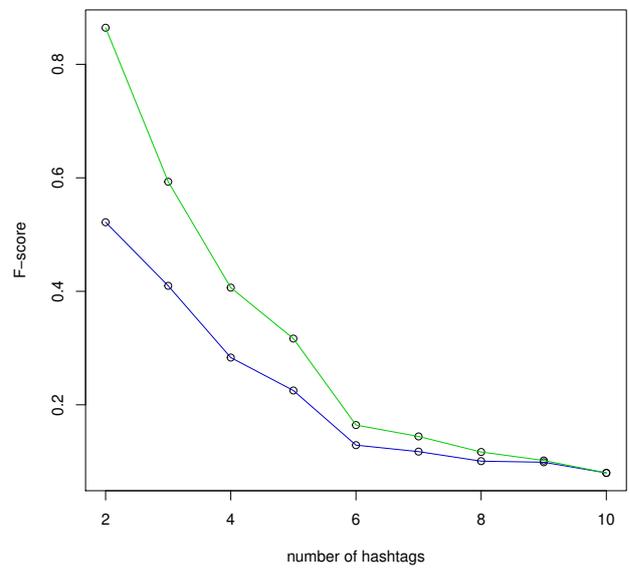


Fig. 17. F-scores for various numbers of hashtags; K-embedding based clustering with k++.md configuration

Average F-score – blue – over all hashtag sets, green – 10 top values

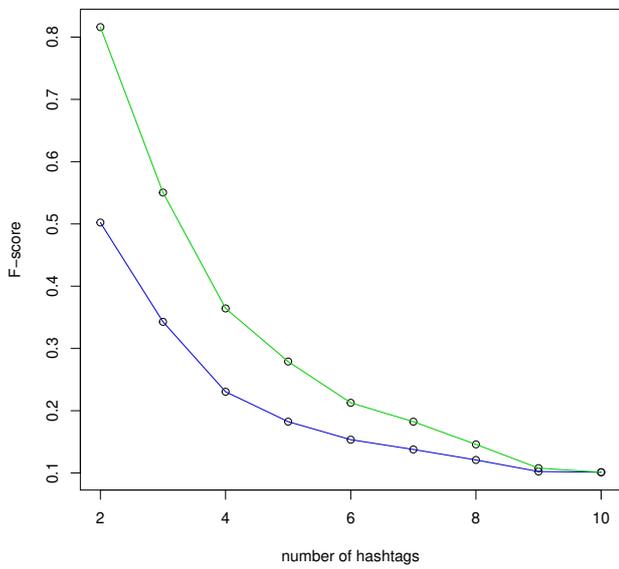


Fig. 16. F-scores for various numbers of hashtags; K-embedding based clustering with sc.n configuration

Average F-score – blue – over all hashtag sets, green – 10 top values

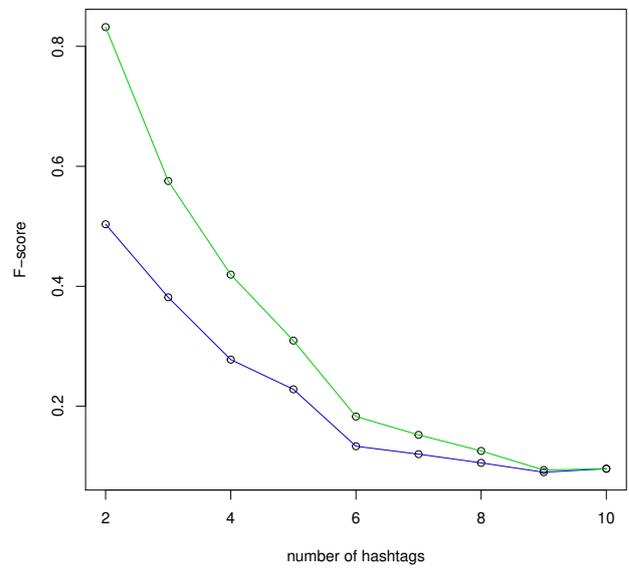


Fig. 18. F-scores for various numbers of hashtags; K-embedding based clustering with k++.n configuration

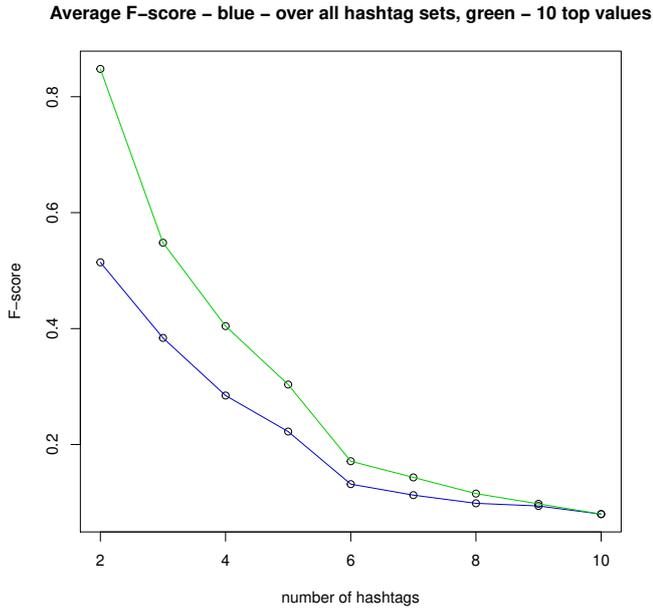


Fig. 19. F-scores for various numbers of hashtags; K-embedding based clustering with k++.sc configuration

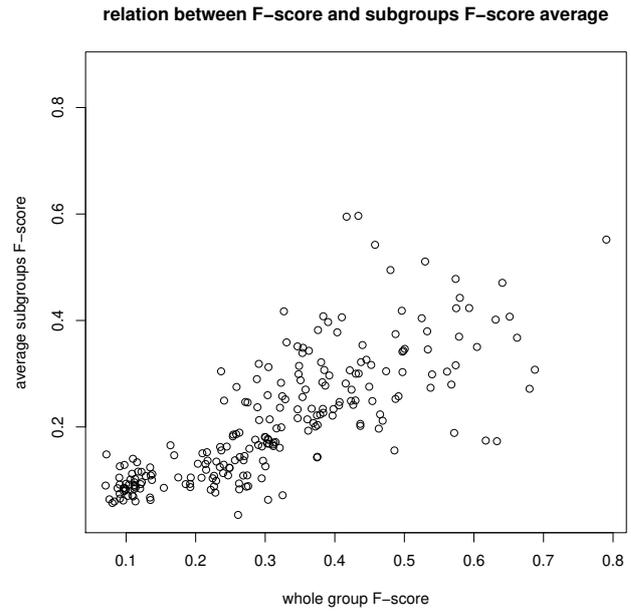


Fig. 21. Relationship between F-score of the given group that was clusters and the average F-score of its subgroups (with one less hashtag); spectral clustering with affinity nearest_neighbors

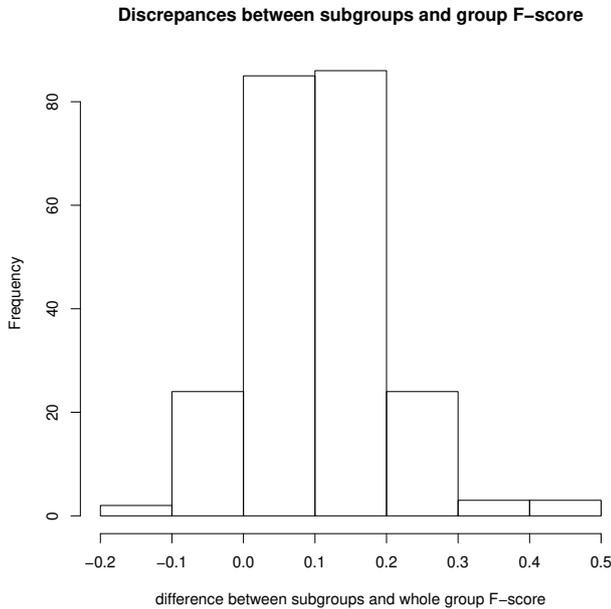


Fig. 20. Difference (negated) between F-score of the given group that was clustered and the average F-score of its subgroups (with one less hashtag); spectral clustering with affinity nearest_neighbors

VI. RESULTS

As visible from Figs. 6–19, the increase of the number of intended clusters to be discovered constitutes a problem for the clustering algorithms, with even 9-fold decrease of F-score when going from 2 to 10 clusters. This behaviour is consistent throughout all the investigated methods though minor variations of the shape of the curves may be observed.

Spherical *k*-means clustering with sc.n configuration appears to perform best for the 10 top pairs of hashtags (Fig. 10) and with sc.sc configuration (Fig. 9), followed by K-embedding based clustering with most configurations (Figs. 14–19, except 16).

In most cases the top average of the F-score for next higher number of cluster is usually higher than the average score for the entire previous number of clusters, which indicates that better separation of subgroups gives some advantage for the capability to separate the entire group.

Table II shows Spearman and Pearson correlations between the F-score achieved by grouping a dataset related to a given set of hashtags and by grouping datasets obtained by removing data of one of the hashtags, split by the clustering algorithm. The correlations are generally high and are statistically very significant. This means that clustering capability of subsets of hashtags can be a good indicator of clustering capability for the set of hashtags. The algorithm spherical sc.sc seems to perform best for such a criterion, followed by spherical k++.sc and in the column on Spearman correlation – K-embedding.12plus k++.md.

A more detailed insight into this relationship for one of the algorithms is presented in Figs 20 and 21. Fig. 20 convinces us, however, that generally this clustering capability decreases (the F-score of a group is usually lower than that of the average of the subgroups). Fig. 21 shows additionally, that the high correlations between group and subgroups of hashtags are to be expected rather for low values of F-score. Higher F-score values are responsible for higher variation in supergroup F-score.

VII. CONCLUSIONS

The performed experiments demonstrate that, in spite of the generally praised properties, graph spectral clustering methods have still a large space for improvements with respect to increasing number of clusters to be detected. Even if all the subsets of intended clusters may be well separated by the algorithms, their mixture does not so. Same observation can be made about the spherical k -means algorithm.

REFERENCES

- [1] S. T. Wierchoń and M. A. Kłopotek, *Modern Clustering Algorithms*, ser. Studies in Big Data. Springer Verlag, 2018, vol. 34. ISBN 978-3-319-69307-1. doi: <https://doi.org/10.1007/978-3-319-69308-8>
- [2] P. Łoziński, D. Czerski, and M. A. Kłopotek, "Grammatical case based IS-A relation extraction with boosting for polish," in *Proceedings of the 2016 Federated Conference on Computer Science and Information Systems, FedCSIS 2016, Gdańsk, Poland, September 11-14, 2016*, ser. Annals of Computer Science and Information Systems, M. Ganzha, L. A. Maciaszek, and M. Paprzycki, Eds., vol. 8. IEEE, 2016. doi: [10.15439/2016F391](https://doi.org/10.15439/2016F391) pp. 533–540. [Online]. Available: <https://doi.org/10.15439/2016F391>
- [3] J. Dörpinghaus, S. Schaaf, J. Fluck, and M. Jacobs, "Document clustering using a graph covering with pseudostable sets," in *Proceedings of the 2017 Federated Conference on Computer Science and Information Systems, FedCSIS 2017, Prague, Czech Republic, September 3-6, 2017*, ser. Annals of Computer Science and Information Systems, M. Ganzha, L. A. Maciaszek, and M. Paprzycki, Eds., vol. 11, 2017. doi: [10.15439/2017F84](https://doi.org/10.15439/2017F84) pp. 329–338. [Online]. Available: <https://doi.org/10.15439/2017F84>
- [4] P. Borkowski, M. A. Kłopotek, B. Starosta, S. T. Wierchoń, and M. Sydow, "Eigenvalue based spectral classification," *PLoS ONE*, vol. 18, no. 4, p. e0283413, 2023. doi: <https://doi.org/10.1371/journal.pone.0283413>
- [5] I. S. Dhillon and D. S. Modha, "Concept decompositions for large sparse text data using clustering," *Machine Learning*, vol. 42, no. 1, pp. 143–175, Jan 2001. doi: <https://doi.org/10.1023/A:1007612920971>
- [6] S. T. Wierchoń and M. A. Kłopotek, "Spectral cluster maps versus spectral clustering," in *Computer Information Systems and Industrial Management*, ser. LNCS, vol. 12133. Springer, 2020. doi: [10.1007/978-3-030-47679-3_40](https://doi.org/10.1007/978-3-030-47679-3_40) pp. 472–484.
- [7] C. R. Harris, K. J. Millman, S. J. van der Walt, R. Gommers, P. Virtanen, D. Cournapeau, E. Wieser, J. Taylor, S. Berg, N. J. Smith, R. Kern, M. Picus, S. Hoyer, M. H. van Kerkwijk, M. Brett, A. Haldane, J. Fernández del Río, M. Wiebe, P. Peterson, P. Gérard-Marchant, K. Sheppard, T. Reddy, W. Weckesser, H. Abbasi, C. Gohlke, and T. E. Oliphant, "Array programming with NumPy," *Nature*, vol. 585, p. 357–362, 2020. doi: [10.1038/s41586-020-2649-2](https://doi.org/10.1038/s41586-020-2649-2) <https://numpy.org>
- [8] P. Virtanen, R. Gommers, T. E. Oliphant, M. Haberland, T. Reddy, D. Cournapeau, E. Burovski, P. Peterson, W. Weckesser, J. Bright, S. J. van der Walt, M. Brett, J. Wilson, K. J. Millman, N. Mayorov, A. R. J. Nelson, E. Jones, R. Kern, E. Larson, C. J. Carey, Í. Polat, Y. Feng, E. W. Moore, J. VanderPlas, D. Laxalde, J. Perktold, R. Cimrman, I. Henriksen, E. A. Quintero, C. R. Harris, A. M. Archibald, A. H. Ribeiro, F. Pedregosa, P. van Mulbregt, and SciPy 1.0 Contributors, "SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python," *Nature Methods*, vol. 17, pp. 261–272, 2020. doi: [10.1038/s41592-019-0686-2](https://doi.org/10.1038/s41592-019-0686-2) <https://scipy.org>
- [9] L. Buitinck, G. Louppe, M. Blondel, F. Pedregosa, A. Mueller, O. Grisel, V. Niculae, P. Prettenhofer, A. Gramfort, J. Grobler, R. Layton, J. VanderPlas, A. Joly, B. Holt, and G. Varoquaux, "API design for machine learning software: experiences from the scikit-learn project," in *ECML PKDD Workshop: Languages for Data Mining and Machine Learning*, 2013. doi: [10.48550/arXiv.1309.023](https://doi.org/10.48550/arXiv.1309.023) pp. 108–122, <https://scikit-learn.org>
- [10] H. Kim and H. K. Kim, "clustering4docs github repository," 2020, <https://pypi.org/project/soyclustering/>. [Online]. Available: <https://github.com/lovit/clustering4docs>
- [11] H. Kim, H. K. Kim, and S. Cho, "Improving spherical k-means for document clustering: Fast initialization, sparse centroid projection, and efficient cluster labeling," *Expert Systems with Applications*, vol. 150, p. 113288, 2020. doi: <https://doi.org/10.1016/j.eswa.2020.113288>. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0957417420301135>
- [12] U. von Luxburg, "A tutorial on spectral clustering," *Statistics and Computing*, vol. 17, no. 4, pp. 395–416, 2007. doi: <https://doi.org/10.48550/arXiv.0711.0189>
- [13] P. Macgregor and H. Sun, "A tighter analysis of spectral clustering, and beyond," in *Proceedings of the 39th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, K. Chaudhuri, S. Jegelka, L. Song, C. Szepesvari, G. Niu, and S. Sabato, Eds., vol. 162. PMLR, 17–23 Jul 2022. doi: <https://doi.org/10.48550/arXiv.2208.01724> pp. 14 717–14 742. [Online]. Available: <https://proceedings.mlr.press/v162/macgregor22a.html>
- [14] Y. Xu, A. Srinivasan, and L. Xue, *A Selective Overview of Recent Advances in Spectral Clustering and Their Applications*. Cham: Springer International Publishing, 2021, pp. 247–277. ISBN 978-3-030-72437-5. doi: [10.1007/978-3-030-72437-5_12](https://doi.org/10.1007/978-3-030-72437-5_12)
- [15] C. Manning, P. Raghavan, and H. Schütze, *Introduction to Information Retrieval*. New York, NY, USA: Cambridge University Press, 2008. doi: <https://doi.org/10.1017/CBO9780511809071>
- [16] F. Krzakala, C. Moore, E. Mossel, J. Neeman, A. Sly, L. Zdeborov, and P. Zhang, "Spectral redemption in clustering sparse networks," in *Proc. the National Academy of Sciences*, vol. 110[50], 2013. doi: [10.48550/arXiv.1306.5550](https://doi.org/10.48550/arXiv.1306.5550) pp. 20 935–20 940.
- [17] H. T. Ali and R. Couillet, "Improved spectral community detection in large heterogeneous networks," *Journal of Machine Learning Research*, vol. 18, no. 225, pp. 1–49, 2018. [Online]. Available: <http://jmlr.org/papers/v18/ali17-247.html>
- [18] A. Saade, F. Krzakala, and L. Zdeborová, "Spectral clustering of graphs with the bethe hessian," 2014. [Online]. Available: <https://arxiv.org/abs/1406.1880>. doi: [10.48550/ARXIV.1406.1880](https://doi.org/10.48550/ARXIV.1406.1880)
- [19] Y. Endo and S. Miyamoto, "Spherical k-means++ clustering," in *Modeling Decisions for Artificial Intelligence*, V. Torra and T. Narukawa, Eds. Cham: Springer International Publishing, 2015. doi: https://doi.org/10.1007/978-3-319-23240-9_9. ISBN 978-3-319-23240-9 pp. 103–114.
- [20] S. Ji, D. Xu, L. Guo, M. Li, and D. Zhang, "The seeding algorithm for spherical k-means clustering with penalties," *J. Comb. Optim.*, vol. 44, no. 3, p. 1977–1994, oct 2022. doi: [10.1007/s10878-020-00569-1](https://doi.org/10.1007/s10878-020-00569-1). [Online]. Available: <https://doi.org/10.1007/s10878-020-00569-1>
- [21] J. Knittel, S. Koch, and T. Ertl, "Efficient sparse spherical k-means for document clustering," in *Proceedings of the 21st ACM Symposium on Document Engineering, DocEng '21*. ACM, New York, NY, United States, 2021. doi: <https://doi.org/10.1145/3469096.3474937> pp. 1–4.
- [22] R. Pratap, A. Deshmukh, P. Nair, and T. Dutt, "A faster sampling algorithm for spherical k-means," in *Proceedings of The 10th Asian Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, J. Zhu and I. Takeuchi, Eds., vol. 95. PMLR, 14–16 Nov 2018, pp. 343–358. [Online]. Available: <https://proceedings.mlr.press/v95/pratap18a.html>
- [23] R. A. Kłopotek, M. A. Kłopotek, and S. T. Wierchoń, "A feasible k-means kernel trick under non-euclidean feature space," *International Journal of Applied Mathematics and Computer Science*, vol. 30, no. 4, pp. 703–715, 2020. doi: <https://doi.org/10.34768/amcs-2020-0052> Online publication date: 1-Dec-2020.
- [24] J. Gower, "Some distance properties of latent root and vector methods used in multivariate analysis," *Biometrika*, vol. 53(3-4), pp. 325–338, 1966. doi: <https://doi.org/10.1093/biomet/53.3-4.325>