

Analysis of Clustering Methods with CBIR- A Survey

Romsha Dhamija¹, Prof.Jaimala Jha²

P.G. Student, Department of Computer Science & Engineering, Madhav Institute of Technology & Science, Gwalior,
India¹

Asst Prof ,Department of Computer Science & Engineering, Madhav Institute of Technology & Science, Gwalior,
India²

ABSTRACT: Content based image retrieval is a broadly used method for retrieving images from the enormous image collection. With the progression of new techniques in recent years, users are not satisfied with traditional technologies. In a charge to find the relationship between data points and dataset of pixels of an image, clustering supplies testing and suitable confirmation. An unsupervised method for withdrawing unseen patterns and developing the adaptation of CBIR is known as Clustering. In this paper, the various clustering techniques and methods are demonstrated and inspected.

KEYWORDS: CBIR, Clustering, Image Retrieval, Feature Extraction, K-means clustering.

I. INTRODUCTION

Content Based Image Retrieval (CBIR) [1] is a well-adopted efficient and quick retrieval technique. The term CBIR is an application to Computer Vision [2] used for retrieving visual digital images from a vast image database according to user concern [3]. The image features such as color, shape, texture are extracted [4] and compared with the image dataset. The query image is given as an input and, the most similar images are extracted as output by the CBIR engine. We can extract color features of an image like RGB, HSV, HSL, etc by the methods such as Intersection, Image Bitmap, Color Correlogram [5] Color Histograms, Color Descriptors, Color Indexing, etc for color-based retrieval. The texture of an image defines the contents or a region of an image added with color features. The features such as coarseness, regularity, periodicity, randomness, and smoothness [6], etc are extracted by methods like structural, statistical, model-based, transform information, [7] gray level co-occurrence matrix (GLCM), Gabor filter and global neighborhood structure (GNS), etc for texture based retrieval. The shape features of an image like the physical structure of the objects, geometric shapes, visual features, etc are extracted by methods such as Fourier descriptor, template matching, quantized descriptor and elementary descriptor [6], etc for shape-based retrieval.

The CBIR is a serving area in various applications such as fingerprint identification, biodiversity information systems, digital libraries, crime prevention, medical and historical research [8], etc. Some challenging problems with CBIR are partial and semantic similarity, domain variance, multi-object images, fast research, local or global features, properties of the contents, etc.

A. CLUSTERING

Clustering is a popular method to be work in unsupervised learning. The idea behind clustering is to split data into clusters of like data items belong to the same cluster and unlike data items to different clusters. It encloses several different algorithms and methods that are all used for grouping objects of similar types into respective categories [9]. Distance and similarity are the roots for building clustering algorithms. Distance is chosen to recognize the relationship among quantitative data features while similarity is chosen for qualitative data features [10]. To test the validity of an algorithm evaluation indicator is used. Clustering achieves an extremely significant role in CBIR to advance the accurateness in an image retrieval process. The major issue about clustering is to find out some clusters and a good inter-cluster distance. Clustering is useful in many applications as it finds out the hidden features and patterns so broadly used in pattern recognition, spatial data analysis, image processing, economic science, social networks, data mining, biology, market research [11], etc.

B. BENEFITS OF CLUSTERING

- It increases the availability of resources that if an individual smart system in a cluster crashed, the other system server in the cluster can take up the amount of workload. This avoids loss of information and time if the server fails.
- It increases the performance of multiple systems and offers the extent of power processing and has greater scalability that as the user increases the resources can grow.
- It has a deliberate resource usage which distributes projects in any arrangement. All machines did not need to be run parallel which reduces overhead and increases resource flexibility.

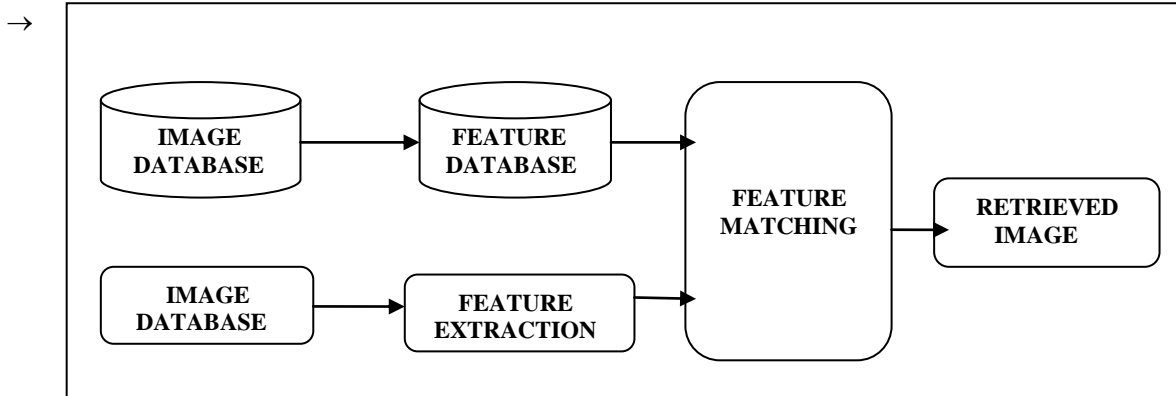


Fig. 1 Functional Diagram of CBIR.

C. BASIC STEPS TO BE PERFORMED WITH CLUSTERING

- The first step checks the representation of type, dimensions and features of available data by feature selection and extraction.
- The second step measures the similarities between data points. Various methods are used such as Euclidean Distance, Mean Square etc.
- The third step collects and groups the data points into clusters based on similarity measures.
- The fourth step performs data abstraction and calculates the centroid of the cluster and uses it for final processing.
- The final step evaluates the clustering result and judges the validity of the output result[12].

II. RELATED WORK

There are numeral techniques in CBIR performed on various clustering methods with color, shape, and texture features in precedent works. The salient works in a CBIR system with clustering is introduced in this section. In [3], presented an approach for content based image retrieval for searching for texture and non-texture images through color and texture features. The system employs the region based segmentation and then performs the k-means clustering of images in the INRIA holiday dataset and get faster retrieval with an accuracy of 97.56%. In [8] presents a method HDK(Hierarchical and Divide and Conquer K means clustering technology) to defeat difficulties of optimizing the limitations of high dimensionality no. of clusters and enhance proficiency and accuracy. Testing is performed on the Corel image database and achieves more than (90%) accuracy for large datasets. In [1], proposed a Bounded Laplace Mixture Model (BLMM) for wavelet coefficients and apply it to image clustering and CBIR. In this experiment, BLMM is evaluated using UIUC, KTH-TIPS, and DTD texture databases. The BLMM(75.48%) has exceeded LMM (72.34%) in data modelling. In [13] presented an efficient retrieval system which includes feature extraction using CCA, clustering, machine learning techniques along with image classifier for grouping the images and to effectively retrieve the medical images. They validate the result on 5 different datasets with SVM, NN, a deep convolutional neural network which shows that retrieval performance is fast and robust by precision (95.9%) and recall (94.6%) with the accuracy of 95.798%. In [14] proposed a combination of several feature types on IMPLIcity database. After this k-means clustering algorithm to give the centroid of image features which gives fruitful result as 100% accuracy of retrieval. In[15] they presented a pre processing step of wavelet transform on their own image dataset. They perform k-means clustering and achieves accuracy of 77%. In [16] implemented narrative techniques for image retrieval using

clustering features. The proposed techniques Row Mean Clustering (RMC), Row Mean Wavelet Clustering (RMWC) are compared with the hierarchical clustering technique. The approach on the Air Freight dataset shows 80% average precision and recall of RMC and RMWC over hierarchical (56.41 %) clustering technique. In [17], followed a new CBIR that uses t-SNE data reduction method with density-based clustering. The system combines CBIR with a clustering technique and dimensionality reduction that clusters reduced images with introducing data parameters which is performed on wang and zubud database results fast and robust system with 81% precision and requires less storage space. In [18], adopted a CBIR system for computing similarity among numerous images. They perform feature extraction with the color histogram method on the Wang database and then apply clustering on the extracted features. K-means clustering is to be used for image indexing. They used Euclidean distance to measure the similarity among databases and finds results with precision (0.68) and recall. Table I and II shows the analysis of work done on the basis of their algorithm, dataset and results in accuracy and precision recall. Table III shows work done on their theoretical results.

Table I- Comparison of various algorithms based on accuracy.

S.NO.	AUTHORS & YEAR	ALGORITHM USED	DATASET	ACCURACY
1.	Lakshmi R.Nair, et.al, [2019]	CCA with Fuzzy c-means	Medical images	95.798%
2.	Monika Jain, Dr. S.K.Singh [2018]	Hierarchical and Divide and Conquer K-means clustering technology (HDK)	Corel image database	90%
3.	Muhammad Azam, Nizar Bouguila, [2018]	Bounded Laplace Mixture Model (BLMM)	UIUC, KTH-TIPS and DTD texture databases	75.48%
4.	Vishal Lonarkar, Ashwath Rao,[2017]	K-means clustering with region segmentation.	INRIA holiday dataset	97.65%
5.	S. Vadivel Murugan, R. Chihra [2017]	K-means Clustering	Image dataset	77%
6.	Mostafa G.Saeed,et.al [2017]	K-means Clustering	IMPLIcity database	100%
7.	M.Karthikeyan, P.Aruna [2012]	dbscan and K-means Clustering	Image dataset	94.4%

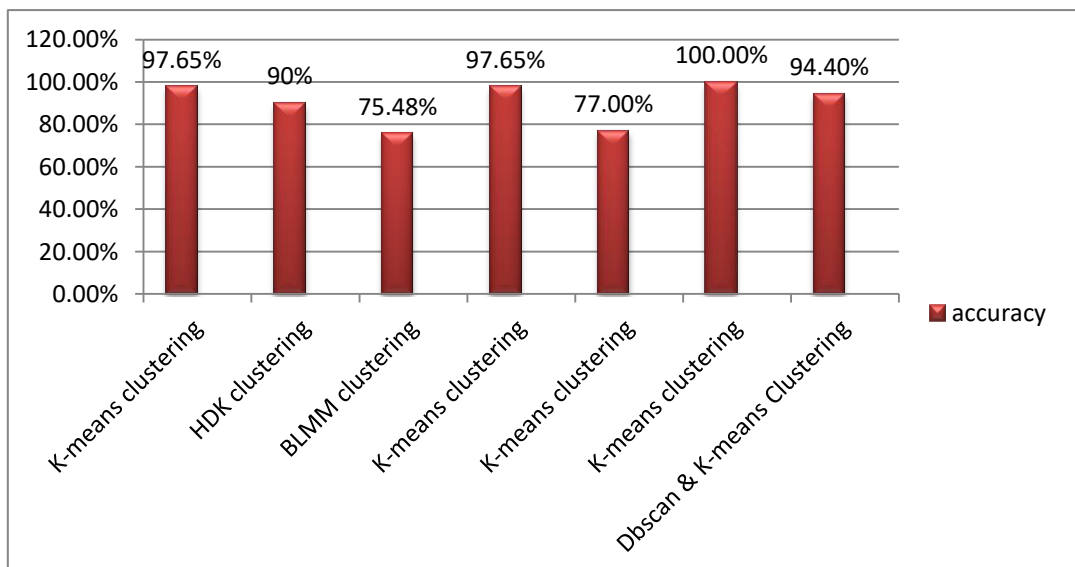


Fig. 2 Comparison of various clustering methods based on previous work done.



Table- II Comparison of various algorithms based on Precision/ Recall

S.NO.	AUTHOR & YEAR	ALGORITHM USED	DATASET	PRECISION/ RECALL
1.	Manoj Kumar [2016]	Euclidean distance	Medical images	Precision=63.18% Recall= 57.84%
2.	Hadjer LACHEHEB [2016]	t-SNE with Density Based Clustering	Zubud Database	Precision = 0.81 Recall = 0.67
3.	Snehal Mahajan, et.al, [2014]	CLBP contribution based clustering	771 images dataset	Precision= 0.80 Recall= 0.49
4.	Swapna Borde, Udhav Bhosle,[2012]	Row means clustering(RMC)	Air Freight dataset	Precision= 80% Recall= 80%
5.	CAI-YUNZHAO, et.al, [2010]	ART-2 with Hierarchical k-means clustering	5000 images with texture and color features	Precision= 98.8% Recall= 78.3%

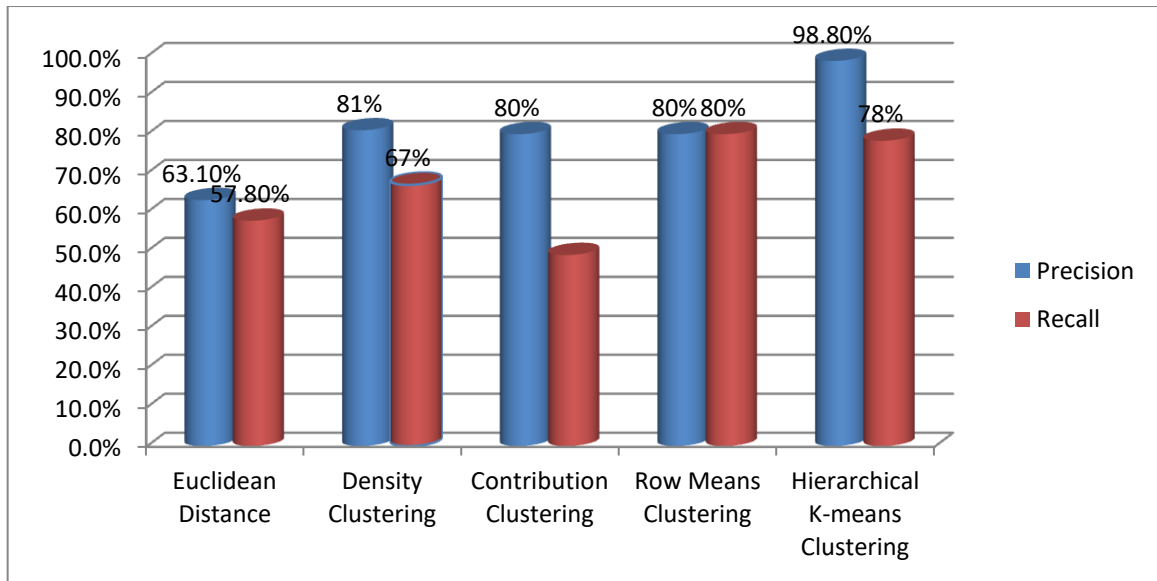


Fig (c) Comparison of various algorithms based on precision and recall.



Table III- Comparison of various algorithms based on theoretical aspects.

S.No.	AUTHOR & YEAR	ALGORITHM USED	DATASET	THEORETICAL RESULTS
1.	N.S. MURTHY, et.al [2010]	Hierarchical and K-means clustering	Large image database	Both techniques increase accuracy and decrease the user's effort.
2.	Vanita.G. Tonge [2011]	CBIR with K-means clustering	Digitized images	K-means gives better and accurate results.
3.	Deepika Nagthane [2013]	K-means with LNM, NDC & GDC method	COREL dataset	Elapsed time reduced effectively.
4.	Harikrishnan Narsimhan, et.al, [2010]	Partitional clustering	Benchmark image dataset	Shows better recall compared to k-means clustering.
5.	V.Ramya [2018]	K-means clustering	CBIR research project image dataset	Improves retrieval quality and reduced elapsed time.

III. CLUSTERING METHODS

Clustering methods classify similar pixels into clusters or groups. These methods also represent pixels cluster and image patches as feature vectors. These methods use a distance function to estimate the distance between clusters and pixels. A distance measure function is different for different applications. The aim of clustering is based on interclass and intraclass similarity maximization. The different methods for clustering were classified as-

A.PARTITION BASED

Partitioning methods separate the data points into a set of disjoint clusters[12]. It constructs k partitions of the data in which each object must belong to exactly one class and each class must contain at least one object. Each partition contains one cluster[19].The partitional clustering arranges data into a single partition instead of a nested structure.[12] The drawback of partitional algorithms is misleading due to the overlapping of data points. The partitional clustering is classified as square error, k-means, mean shift clustering-

Square Error

The mostly used clustering is based on square root error criteria. This method partitions the data into a fixed number of clusters and randomly initialize data points to predefined clusters based on the similarity between the data point and cluster center until convergence criteria[12].Then minimizes the square error where the square error is the sum of Euclidean distances between each pattern and its centroid of clusters. The square error algorithm works well with isolated and compact clusters [12].

K-means Clustering

K-means clustering is an unsupervised type of clustering algorithm which is used with unlabeled data. In this algorithm, the numbers of clusters are initially defined. This randomly initializes each data point to a single cluster[12]. The algorithm works iteratively to assign data points to one of the clusters based on feature similarity which is represented by k groups. The resulting groups examine the centroid features weight and interpret the kind of group each cluster represents. The output of the k-means algorithm depends on the selection of properly chosen clusters either it may result in a false cluster location.

Mean Shift Clustering

The mean shift clustering is a non-parametric feature space analysis algorithm that clusters the given dataset by locating the maxima of a density function. For each point, mean shift iteratively finds the high-density regions and shifting our window in that direction, until we reach convergence, i.e. until the shift is less than the threshold. The Mean shift is quite better then k means because we do not need to specify the number of clusters initially.

B.HIERARCHICAL BASED

Hierarchical clustering is a nested chain of partitions[20] that represents data into a structured format like a tree[12]. In this whole data is represented as a root node, the intermediate node in the tree structure represents the similarity between data points and the leaf node represents the individual data points. This method works on both top-down and

bottom-up approaches[19]. At different levels in structure different classes of clustering will form as per user concern. The distance function and linkage criteria are important factors of hierarchical clustering[20]. Based on the above terms hierarchical clustering is classified as agglomerative and divisive technique-

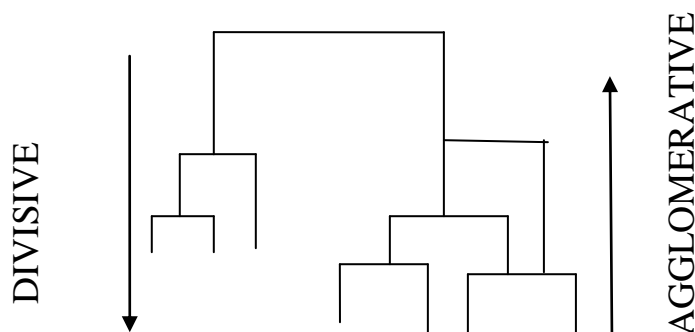


Fig (d) Hierarchical Clustering

Agglomerative

In agglomerative hierarchical clustering initially, each single data point is considered as an individual cluster. At each repetition similar cluster merges with other clusters until k clusters are formed. It starts with a single element with variant patterns and combines it until stopping criteria reached. The output is a tree based representation of objects, named dendrogram. This approach is also known as the bottom-up approach.

Divisive

A Divisive clustering algorithm is called a top-down approach. This clustering works on global information of data and splits the data continuously that have higher dissimilarity among them until individual clusters will form as per user wish. This does not require some clusters to be defined prior. Divisive clustering is more complex, efficient and accurate as compared to agglomerative clustering.

C.GRID BASED

Grid based methods use a multi-resolution grid data structure. These methods allocate the dataset and objects to a set of rectangular grid cells to a specific location. It sets a predefined threshold value t and computes the density of each cell and eliminates those cells whose density is below the threshold value[21]. It then forms clusters from adjacent class cells. Grid based methods are mainly classified as STING and CLIQUE. In the STING algorithm, there are different levels of rectangular cells forms a hierarchical structure corresponds to different resolutions cell. CLIQUE constructs static grids and is a bottom-up approach that operates on multidimensional data.

Graph based

This kind of clustering is observed on the graph where nodes of the weighted graph coincide to the data point and edges reflect the relationship between data points. A graph based clustering is well needed for generating minimum spanning tree cluster analysis.

N-Cut

N-Cut is a hierarchical divisive clustering which accomplished by constructing graph utilizing similarities among neighbouring objects. This arranges the nodes into clusters or groups so that the similarity amongst nodes in one group is high and similarity among the group is low. In each phase the sub graph with maximum number of nodes is additionally split into clusters until stopping criteria reached[12]. The output results in more than two clusters. So, useful for applications where multiple clusters needed.

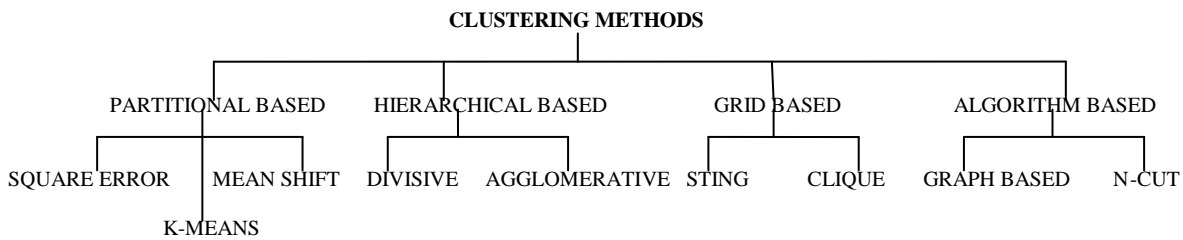


Fig.5 Chart showing different clustering methods

Table IV- Comparison of various clustering methods

CLUSTERING METHODS	ADVANTAGES	DISADVANTAGES
Partition Based	When the value of k is small, it is faster than hierarchical clustering	It is crucial to estimate the value of k.
Hierarchical Based	It is not required to know earlier the no. of clusters	Compassion to noise and outliers
Grid Based	It has quick dispensation	To define the no.of grid is difficult
Density Based	Discovers groups of arrogant shapes and is well-organized for large databases	Algorithms are based on Gaussian, Bernoulli and Poisson distribution.
Model Based	It displays a distinctive explanation for each group, where each group constitutes a class	It is not appropriate for large database

IV. CONCLUSION

The main motive of this survey is to discuss the various techniques and parameters for clustering methods. This paper introduces the basic functionality of CBIR systems and the core idea of each commonly used clustering algorithms. Clustering techniques are useful in finding similarities between images in a collection of unlabeled data. Partition based clustering constructs k-partitions of data and each data belongs to exactly one cluster and evaluates them by some criteria. Hierarchical clustering creates a set of nested clusters arranged in a nested tree. Grids based methods form finite cells in the grid structure and have fast processing time.

REFERENCES

- [1] M. Azam and N. Bouguila, “Bounded Laplace Mixture Model with Applications to Image Clustering and Content Based Image Retrieval,” Proc. - 17th IEEE Int. Conf. Mach. Learn. Appl. ICMLA 2018, pp. 558–563, 2019, doi: 10.1109/ICMLA.2018.00090.
- [2] V. Tyagi and V. Tyagi, “Content-Based Image Retrieval Techniques: A Review,” Content-Based Image Retr., pp. 29–48, 2017, doi: 10.1007/978-981-10-6759-4_2.
- [3] V. Lonarkar and B. A. Rao, “Content-based image retrieval by segmentation and clustering,” Proc. Int. Conf. Inven. Comput. Informatics, ICICI 2017, no. Icici, pp. 771–776, 2018, doi: 10.1109/ICICI.2017.8365241.
- [4] A. J. Afifi and W. M. Ashour, “Image Retrieval Based on Content Using Color Feature,” ISRN Comput. Graph., vol. 2012, pp. 1–11, 2012, doi: 10.5402/2012/248285.
- [5] D. Srivastava, R. Wadhvani, and M. Gyanchandani, “A Review: Color Feature Extraction Methods for Content Based Image Retrieval,” Int. J. Comput. Eng. Manag., vol. 18, no. 3, pp. 9–13, 2015.
- [6] Monika Jain and S.K.Singh, “A Survey On: Content Based Image Retrieval Systems Using Clustering Techniques For Large Data sets,” Int. J. Manag. Inf. Technol., vol. 3, no. 4, pp. 23–39, 2011, doi: 10.5121/ijmit.2011.3403.
- [7] J. O. Olaleke, A. O. Adetunmbi, B. A. Ojokoh, and I. Olaronke, “An Appraisal of Content-Based Image Retrieval (CBIR) Methods,” Asian J. Res. Comput. Sci., no. April, pp. 1–15, 2019, doi: 10.9734/ajrcos/2019/v3i230089.

- [8] M. Jain and S. K. Singh, "An efficient content based image retrieval algorithm using clustering techniques for large dataset," 2018 4th Int. Conf. Comput. Commun. Autom. ICCCA 2018, pp. 1–5, 2018, doi: 10.1109/CCAA.2018.8777591.
- [9] H. Lacheheb and S. Aouat, "SIMIR: New mean SIFT color multi-clustering image retrieval," *Multimed. Tools Appl.*, vol. 76, no. 5, pp. 6333–6354, 2017, doi: 10.1007/s11042-015-3167-3.
- [10] D. Xu and Y. Tian, "A Comprehensive Survey of Clustering Algorithms," *Ann. Data Sci.*, vol. 2, no. 2, pp. 165–193, 2015, doi: 10.1007/s40745-015-0040-1.
- [11] A. Cord, C. Ambroise, and J. P. Cocquerez, "Feature selection in robust clustering based on Laplace mixture," *Pattern Recognit. Lett.*, vol. 27, no. 6, pp. 627–635, 2006, doi: 10.1016/j.patrec.2005.09.028.
- [12] S. Ghosh, S. K. S. Dubey, A. Gulhane, P. L. Paikrao, and D. S. Chaudhari, "A Review of Image Data Clustering Techniques," *Ijacs*, vol. 2, no. 1, pp. 35–38, 2012, doi: 10.14569/IJACSA.2013.040406.
- [13] L. R. Nair, K. Subramaniam, and G. K. D. P. Venkatesan, "An effective image retrieval system using machine learning and fuzzy c- means clustering approach," *Multimed. Tools Appl.*, 2019, doi: 10.1007/s11042-019-08090-2.
- [14] M. G. Saeed, F. L. Malallah, and Z. A. Aljawaryy, "Content-Based Image Retrieval by Multi-Featrus Extraction and K-Means Clustering," *Int. J. Electr. Electron. Comput.*, vol. 2, no. 3, pp. 1–11, 2017, doi: 10.24001/eec.2.3.1.
- [15] M. Taileb, "Content based image retrieval system using NOHIS-tree," *ACM Int. Conf. Proceeding Ser.*, vol. 11, no. December, pp. 103–108, 2012, doi: 10.1145/2428955.2428980.
- [16] B. Mounika, Y. Sowmya, S. Pasala, and A. Sravani, "Content based image retrieval using color," *Int. J. Appl. Eng. Res.*, vol. 11, no. 6, pp. 4331–4334, 2016, doi: 10.5120/ijca2016909633.
- [17] H. Lacheheb and S. Aouat, "A density clustering approach for CBIR system," *Proc. IEEE/ACS Int. Conf. Comput. Syst. Appl. AICCSA*, vol. 0, 2016, doi: 10.1109/AICCSA.2016.7945742.
- [18] J. Rejito, A. S. Abdullahi, Akmal, D. Setiana, and B. N. Ruchjana, "Image indexing using color histogram and k-means clustering for optimization CBIR in image database," *J. Phys. Conf. Ser.*, vol. 893, no. 1, 2017, doi: 10.1088/1742-6596/893/1/012055.
- [19] S. J. Swarndepp and S. Pandya, "An Overview of Partitioning Algorithms in Clustering Techniques," *Int. J. Adv. Res. Comput. Eng. Technol.*, vol. 5, no. 6, pp. 1943–1946, 2016, doi: 10.1017/CBO9781107415324.004.
- [20] N. Rajalingam, "Hierarchical Clustering Algorithm - A Comparative Study," vol. 19, no. 3, pp. 42–46, 2011.
- [21] S. Saini and P. Rani, "A Survey on STING and CLIQUE Grid Based Clustering Methods," *Int. J. Adv. Res. Comput. Sci.*, vol. 8, no. 5, pp. 2015–2017, 2017.
- [22] H. Narasimhan and P. Ramraj, "Contribution-based clustering algorithm for content-based image retrieval," 2010 5th Int. Conf. Ind. Inf. Syst. ICIIS 2010, pp. 442–447, 2010, doi: 10.1109/ICIINFS.2010.5578664.
- [23] V. Ramya, "Content Based Image Retrieval System using Clustering with Combined Patterns," *Int. J. Sci. Res. Comput. Sci. Eng. Inf. Technol.*, vol. 3, no. 1, pp. 1060–1063, 2018.
- [24] U. K K, "Web Image Retrieval Using Clustering Approaches," *Comput. Sci. Conf. Proc.*, vol. 3, pp. 27–36, 2011, doi: 10.5121/csit.2011.1304.
- [25] N. Dalvi, "Content Based Image retrieval system using machine learning," *Int. J. Res. Appl. Sci. Eng. Technol.*, vol. 7, no. 6, pp. 538–541, 2019, doi: 10.22214/ijraset.2019.6095.
- [26] V. G. Tonge, "Content Based Image Retrieval by K-Means Clustering Algorithm," pp. 46–49, 2011.
- [27] C. Y. Zhao, B. X. Shi, M. X. Zhang, and Z. W. Shang, "Image retrieval based on improved hierarchical clustering algorithm," 2010 Int. Conf. Wavelet Anal. Pattern Recognition, ICWAPR 2010, vol. 1, no. July, pp. 154–157, 2010, doi: 10.1109/ICWAPR.2010.5576314.
- [28] H. Zhou, A. H. Sadka, M. R. Swash, J. Azizi, and A. S. Umar, "Content based image retrieval and clustering: A brief survey," *Recent Patents Electr. Eng.*, vol. 2, no. 3, pp. 187–199, 2009, doi: 10.2174/1874476110902030187.
- [29] N. V. R. S. K. Javvadi, "Content Based Image Retrieval using Hierarchical and K-Means Clustering Techniques," *Int. J. Eng. Sci. Technol.*, vol. 2, no. 3, pp. 209–212, 2010.
- [30] N. Silpa and K. Santhi, "Featuresand Clustering Techniques," pp. 33–36, 2014.
- [31] M. A. Bouker and E. Hervet, "Image retrieval based on mean-shift clustering using color descriptor," 2012 11th Int. Conf. Inf. Sci. Signal Process. their Appl. ISSPA 2012, pp. 1422–1423, 2012, doi: 10.1109/ISSPA.2012.6310521.
- [32] C. D. Arthur and S. Vassilvitskii, "What Is Clustering?," pp. 1–37, 2012, doi: 10.1201/b13101-2.
- [33] S. Patel and A. Patel, "Performance Analysis and Evaluation of Clustering Algorithms," *Int. J. Innov. Technol. Explor. Eng.*, no. April, pp. 179–183, 2019.



- [34] S. Mahajan and D. Patil, "Image retrieval using contribution-based clustering algorithm with different feature extraction techniques," Proc. 2014 Conf. IT Business, Ind. Gov. An Int. Conf. by CSI Big Data, CSIBIG 2014, 2014, doi: 10.1109/CSIBIG.2014.7057001.
- [35] Z. Wang, G. S. Xia, C. Xiong, and L. Zhang, "Spectral active clustering of remote sensing images," Int. Geosci. Remote Sens. Symp., no. January 2014, pp. 1737–1740, 2014, doi: 10.1109/IGARSS.2014.6946787.