**INDIAN INSTITUTE OF INFORMATION TECHNOLOGY, ALLAHABAD**

Supervised by:
# Prof. Pawan Chakraborty

| *Team Members* | *Enrollment No.* |
| --- | --- |
| Arsh Singh | IIT2021064 |
| Keshav Pandey | IIT2021173 |
| Sumit Das | IIT2021229 |
| Sukankshi Sharma | IIT2021503 |
| P.L.V.B. Narayana | IIB2021026 |

# Comparative Analysis Of Various Clusterings On SDSS Galaxy Images

May 12, 2024

# Contents

# 1 Acknowledgement:

# 2 Abstract:

This article presents a comparison of methods for classifying galaxies based on their visual properties using images from the Sloan Digital Sky Survey (SDSS). Various clustering methods are evaluated in this study, including traditional methods such as K-means, hierarchical clustering, and DBSCAN, as well as different methods such as spectral clustering and Gaussian mixture pattern. Evaluate the effectiveness and efficiency of each common strategy through testing and qualitative analysis using metrics such as silhouette score, Davis-Bourdin index, and factor analysis. Advantages and limitations of different techniques used to classify SDSS galaxy images. It is worth noting that although traditional methods such as K-means and hierarchical clustering have shown some success, more efficient methods such as spectral clustering and Gaussian mixture models have provided improvements in capturing patterns and classifications in galaxy images. Despite its speed-based approach, DBSCAN's performance in classifying SDSS galaxy images may be limited, particularly due to its inability to handle noise and spatial inconsistencies. This observation highlights the importance of considering the unique characteristics of astronomical data when choosing a clustering algorithm. It provides insight into the future of astrophysics and information science.

# 3 Problem Statement:

We will be grouping RGB-encoded Sloan Digital Sky Survey (SDSS) galaxy images using four different clustering strategies: partitioning (K-means), hierarchical (agglomeration), density-based (DBSCAN), and distribution-based (Gaussian mixture model). Our goal is to assess how well different techniques classify the galaxy images.

Next, we use principal component analysis and/or transfer learning (from pre-trained convolutional neural networks) to cluster the RGB galaxy pictures from the Sloan Digitial Sky Survey (SDSS).

# 4 Introduction:

The study of the spatial distribution of galaxies and other celestial objects seen by the Sloan Digital Sky Survey (SDSS) is known as the "clustering."

**SDSS:** Certainly one of the most ambitious and significant sky surveys ever conducted is the Sloan Digital Sky Survey. Its mapping of millions of galaxies and quasars has produced an extensive amount of data for cosmological research.

**Astronomical Catalogs:** With the help of the SDSS, numerous extensive catalogs of galaxies and quasars have been created, each with accurate redshifts and locations. Astronomers can determine distances by measuring an object's recession velocity, which is caused by the universe's expansion. This measurement is called redshift.

**Large-scale Organization:** Large-scale structures including galaxy clusters, filaments, and voids are present in the galaxy clustering analysis of the SDSS data. These formations offer hints regarding the genesis and development of cosmic structure.

**Parameters for Cosmology:** Through an examination of the galaxy clustering statistics seen in the SDSS data, scientists are able to limit cosmological parameters like the universe's matter density, the amplitude of primordial density fluctuations, and the makeup of dark energy.

# 5   Related Works:

**[1]JORGE DE LA CALLEJA, OLAC FUENTES; AUTOMATED CLASSIFICATION OF GALAXY IMAGES**
The rule-induction algorithm C4.5, random forests, and Naive Bayes machine learning techniques are used to classify galaxy images. Prior to applying the classification methods, PCA is also used to minimize dimensionality after completing some initial image processing. The findings of this study show that random forest classifies photos into three categories—elliptical, spiral, and irregular—with an accuracy of almost 91%.

**[2]SIDDHARTHA KASIVAJHULA, NAREN RAGHAVAN, HEMAL SHAH; MORPHO-LOGICAL GALAXY CLASSIFICATION USING MACHINE LEARNING**
Support Vector Machines (SVM), Random Forests (RF), and Naïve Bayes (NB) are three machine learning methods that are used to categorize images. The set of morphic features created by image analysis and direct image pixel data compressed through PCA are used. Ultimately, a comparison is made between the algorithms' performance on the data based on both PCA and morphic features.

**[3]NOUR ELDEEN M. KHALIFA, MOHAMED HAMED N. TAHA, ABOUL ELLA HASSANIEN , I. M. SELIM; DEEP GALAXY: CLASSIFICATION OF GALAXIES BASED ON DEEP CONVOLUTIONAL NEURAL NETWORKS**
Two principal fully connected layers are employed for galaxy classification in a deep convolutional neural network architecture consisting of eight layers: one primary convolutional layer for feature extraction employing 96 filters. Based on the traits, three groups are formed: spiral, elliptical, and irregular.

**[4]JORGE DE LA CALLEJA, OLAC FUENTES; MACHINE LEARNING AND IMAGE ANALYSIS FOR MORPHOLOGICAL GALAXY CLASSIFICATION**
Neural networks, homogeneous ensembles of classifiers, and a locally weighted regression technique are applied for morphological galaxy classification. To obtain pertinent features and minimize dimensionality, data is enhanced and PCA is employed. To categorize images into three categories—spiral, elliptical, and irregular—homogeneous ensemble regression methods are utilized with a 10-fold Cross Validation.

**[5]JOSEPH H. MURRUGARRA LL., NINA S. T. HIRATA; GALAXY IMAGE CLAS-SIFICATION**
To convert photos from the Sloan Digital Sky Survey into labeled format, preprocessing and picture cropping are performed. After extracting features using a convolutional neural net, the images are fed into a Support Vector Machine (SVM) classifier model to determine whether they are spiral or elliptical. We assign an accuracy of roughly 90–91% to this model.
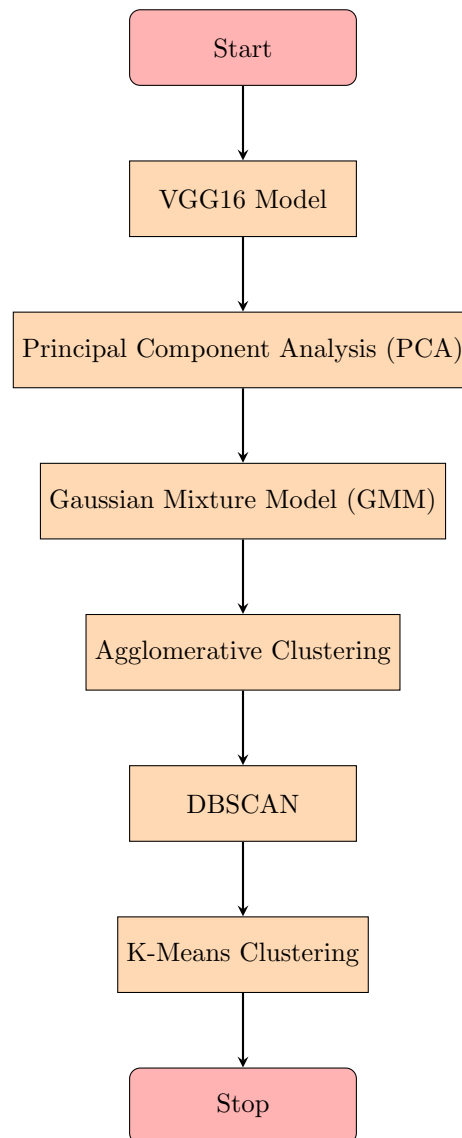
# 6    Flow Chart:



Figure 1: Algorithms Used in Provided Code

# 7 Proposed Method

## 7.1 Data Preprocessing

### 7.1.1 Preprocessing and Loading Images

First, the images are loaded into memory using specific libraries like Astropy in Python after being downloaded from the SDSS database, usually in FITS format. Then, in order to improve the signal-to-noise ratio and make feature extraction easier, these images are preprocessed using methods including noise reduction, background subtraction, and normalization. To further standardize the photos and take into consideration differences in orientation and scale, geometric operations like rotation, cropping, and resizing may be used.

### 7.1.2 Feature Extraction with VGG16

An efficient method for extracting features from the SDSS clustering photos is to use a pre-trained VGG12 model, which is a convolutional neural network architectural variation of the VGG (Visual Geometry Group). The VGG12 model, which has been pre-trained on extensive image datasets such as ImageNet, can efficiently extract high-level features from the SDSS images by utilizing transfer learning.

### 7.1.3 Dimensionality Reduction with PCA

Following the VGG12 model's feature extraction from the SDSS clustering images, dimensionality reduction techniques are used to compress the high-dimensional feature space while keeping the crucial information. reduces feature dimensionality to 100 using PCA.

## 7.2 Clustering

### 7.2.1 K-Means Clustering

The SDSS clustering images are condensed into discrete groups based on their commonalities using K-means clustering with 100 clusters on the reduced feature space. The technique effectively splits the data by repeatedly assigning each image to the closest cluster centroid, exposing underlying patterns and structures in the SDSS galaxy data. This approach makes a substantial contribution to astronomical study by facilitating the exploration of large-scale cosmic structures and helping to comprehend the distribution and features of galaxies in the universe.

K-Means Clustering Algorithm involves the following steps:

1. Calculate the number of K (Clusters).

2. Randomly select K data points as cluster center.

3. Using the Euclidean distance formula measure the distance between each data point and each cluster center.

$$\sqrt{(\mathbf{x_2 - x_1})^2 - (\mathbf{y_2 - y_1})^2}$$

---
**Algorithm 1** K-Means Clustering
---
 1: **procedure** KMEANSCLUSTERING(data, $k$, max_iterations)
 2:     Initialize $k$ centroids randomly
 3:     $iteration \leftarrow 1$
 4:     **while** $iteration \leq$ max_iterations **do**
 5:         Assign each data point to the nearest centroid
 6:         **for** $i \leftarrow 1$ to $k$ **do**
 7:            $cluster\_points_i \leftarrow$ data points assigned to centroid $i$
 8:            **if** $cluster\_points_i$ is not empty **then**
 9:                Update centroid $i$ to the mean of $cluster\_points_i$
10:            **end if**
11:         **end for**
12:         $iteration \leftarrow iteration + 1$
13:     **end while**
14:     **return** the final $k$ centroids and cluster assignments
15: **end procedure**
---

### 7.2.2 DBSCAN Clustering

**Clustering Based on Density:** assembles densely populated points and identifies spots that are isolated in low-density areas as outliers.

**Key Information and Neighborhood Subjects:** Clusters are expanded by linking them to their density-reachable neighbors. It finds core points based on a minimal number of neighboring points within a given distance (epsilon).

**Adaptability in Cluster Dimensions and Shape:** DBSCAN does not presuppose spherical clusters, in contrast to k-means or hierarchical clustering. It is capable of identifying groups of any size or shape.

**Managing Different Densities:** When it comes to capturing differences in galaxy density, DBSCAN works well, allowing for places such as mergers or voids that have variable densities.

**Finding Outliers:** It helps distinguish real features from artifacts by identifying outliers in galaxy images that correspond to rare features or noise.

**Sensible to Non-linear Frameworks:** Because of its adaptability, DBSCAN is able to precisely capture characteristics such as filaments and galaxy groups, even in complicated and non-linear galaxy structures.

One way to sum up the mathematical intuition underlying DBSCAN is as follows:

1. **Density**: High data point density areas are used by DBSCAN to identify clusters. We can express this using the density formula:

$$\rho(p) = |N_\epsilon(p)|$$

where $\rho(p)$ represents the density of point $p$, and $N_\epsilon(p)$ is the set of points within the epsilon neighborhood of $p$.

2. **Core Points**: If the density of a point $p$ is greater than a predetermined threshold $MinPts$, then $p$ is a core point. A core point $p$ mathematically satisfies:

$$\rho(p) \geq MinPts$$

3. **Border Points**: Border points do not match the density criteria to be core points themselves, but they do reside within the epsilon neighborhood of a core point.

4. **Noise Points**: Points categorized as noise points are those that are not part of any cluster.

5. **Reachability**: Reachability is the basis on which DBSCAN decides cluster membership. If a path of core points from point $p$ to $q$ exists within the epsilon radius, then point $q$ can be reached from point $p$.

6. **Relationship**: Two elements If $q$ is in the epsilon neighborhood of $p$, then $p$ and $q$ are directly density-reachable. If there is a chain of directly density-reachable locations connecting them, then they are indirectly density-reachable.

In conclusion, DBSCAN analyzes the connection and density of points in the dataset to pinpoint clusters. This method successfully manages noise and enables the finding of clusters of arbitrary shapes.
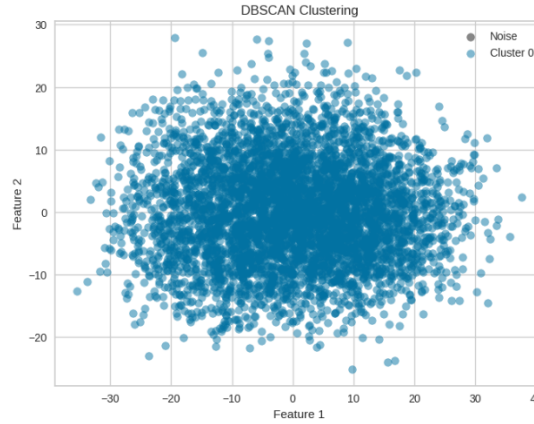


Figure 2: DBSCAN Clustering performed on 2 features

---
**Algorithm 2** DBSCAN Algorithm
---
1: **procedure** DBSCAN(data, $\varepsilon$, minPts)
2:      Initialize an empty set for visited points
3:      Initialize an empty set for noise points
4:      Initialize an empty dictionary to store clusters
5:      **for** each unvisited point $P$ in data **do**
6:          Mark $P$ as visited
7:          Find all points within distance $\varepsilon$ of $P$
8:          **if** the number of such points $\geq$ minPts **then**
9:              Create a new cluster $C$ and add $P$ to $C$
10:              EXPANDCLUSTER(data, $P$, $C$, $\varepsilon$, minPts)
11:          **else**
12:              Add $P$ to noise points
13:          **end if**
14:      **end for**
15:      **return** clusters and noise points
16: **end procedure**
17: **procedure** EXPANDCLUSTER(data, point, cluster, $\varepsilon$, minPts)
18:      Add point to cluster
19:      neighbors $\leftarrow$ REGIONQUERY(data, point, $\varepsilon$)
20:      **for** each neighbor in neighbors **do**
21:          **if** neighbor is not visited **then**
22:              Mark neighbor as visited
23:              neighbor_neighbors $\leftarrow$ REGIONQUERY(data, neighbor, $\varepsilon$)
24:              **if** size of neighbor_neighbors $\geq$ minPts **then**
25:                  Add neighbor and neighbor_neighbors to neighbors
26:              **end if**
27:          **end if**
28:          **if** neighbor is not yet a member of any cluster **then**
29:              Add neighbor to cluster
30:          **end if**
31:      **end for**
32: **end procedure**
33: **procedure** REGIONQUERY(data, point, $\varepsilon$)
34:      Initialize an empty list for neighbors
35:      **for** each point $p$ in data **do**
36:          **if** distance(point, $p$) $\leq \varepsilon$ **then**
37:              Add $p$ to neighbors
38:          **end if**
39:      **end for**
40:      **return** neighbors
41: **end procedure**
---

### 7.2.3 Agglomerative Clustering

**Hierarchical Approach:** Agglomerative clustering creates a hierarchical structure by repeatedly merging the closest clusters until a halting requirement is satisfied.

**Bottom-Up Method:** Larger clusters are gradually formed by taking individual data points as a starting point and merging them based on similarity.

**Dendrogram Visualization:** A dendrogram is a useful tool for analyzing the clustering structure and figuring out the ideal number of clusters because it visualizes the merging process.

**Hierarchical Structures:** Effectively captures hierarchical structures in galaxy images, from individual galaxies to larger galaxy clusters and superclusters.
**Flexible Distance Metrics:** It accommodates features like morphology, color, or spectral charac-

teristics in clustering galaxy images.

**Interpretability:**   The dendrogram aids astronomers in interpreting results and understanding the grouping of galaxies based on shared characteristics.

The mathematical intution can be explain as below:

1. **Initialization**: Put each data point ($x_i$ in a cluster ($C_i$ for the entire amount of data points ($n$, where $i$ varies from 1 to $n$.

2. **Metric of Similarity**: To measure the similarity or dissimilarity of clusters $C_i$ and $C_j$, define a distance metric $d(C_i, C_j)$. This distance measure establishes how close together clusters are and directs the merging process.

3. **Pairwise Merge**: Determine which two clusters are closest at each step $t$ using the selected distance metric, then combine them into a single cluster. The overall number of clusters is lowered by one in this step. Assume that the clusters to be merged at step $t$ are $C_k$ and $C_l$. The definition of the merge is:

$$C_{k \cup l} = C_k \cup C_l$$

4.  **distance Matrix Update**: Determine the distances once more between each of the surviving clusters and the recently generated cluster $C_{k \cup l}$. The distance matrix must be updated in this phase to reflect the updated pairwise distances. One of the linking criteria, such as single, complete, or average linkage, is used to compute the updated distance matrix $D^{(t+1)}$.

5. **Repeat**: Carry out steps 3 and 4 again until there is just a single cluster that has every data point.

6. **Building Hierarchies**: A dendrogram, or tree-like structure, is created as clusters are combined to show the hierarchy of the clusters. The sequence of cluster mergers and the separations between them are both encoded in the dendrogram.

The distance matrix $D$ can be expressed mathematically as follows:

$$D = \begin{bmatrix} d(C_1, C_1) & d(C_1, C_2) & \cdots & d(C_1, C_n) \\ d(C_2, C_1) & d(C_2, C_2) & \cdots & d(C_2, C_n) \\ \vdots & \vdots & \ddots & \vdots \\ d(C_n, C_1) & d(C_n, C_2) & \cdots & d(C_n, C_n) \end{bmatrix}$$

where $d(C_i, C_j)$ represents the distance between clusters $C_i$ and $C_j$.

**Algorithm 3** Agglomerative Clustering

---

1: **procedure** AGGLOMERATIVECLUSTERING(data, $k$)
2:     Initialize each data point as a cluster
3:     **while** number of clusters $> k$ **do**
4:         Find the two closest clusters based on a distance metric
5:         Merge the two closest clusters into a single cluster
6:     **end while**
7:     **return** the final $k$ clusters
8: **end procedure**
9: **procedure** DISTANCE(cluster1, cluster2)
10:     min_distance $\leftarrow \infty$
11:     **for** each point $p1$ in cluster1 **do**
12:         **for** each point $p2$ in cluster2 **do**
13:             distance $\leftarrow$ calculate_distance($p1, p2$)
14:             **if** distance $<$ min_distance **then**
15:                 min_distance $\leftarrow$ distance
16:             **end if**
17:         **end for**
18:     **end for**
19:     **return** min_distance
20: **end procedure**

---

### 7.2.4 Gaussian Mixture

**Probabilistic Model:** Assumes that the data points are generated from a mixture of several Gaussian distributions. Each Gaussian distribution represents a cluster in the data.
**Soft Grouping:** In contrast to rigorous clustering algorithms such as k-means, GMM allocates a probability to every data point that is a part of every cluster.

The algorithm known as Expectation-Maximization (EM): It uses EM to iteratively estimate parameters, fine-tuning cluster allocations and forms.

**Complexity Modeling:** The various features and structures found in galaxy photos are well-suited to the modeling of complicated distributions via GMM.

**Assignment on Probability:** Compared to hard clustering techniques, GMM's probabilistic approach is superior at handling ambiguous galaxy characteristics.

**Adaptability in Group Forms:** Because of its adaptability, GMM is able to capture the differences in galaxies' forms and orientations.
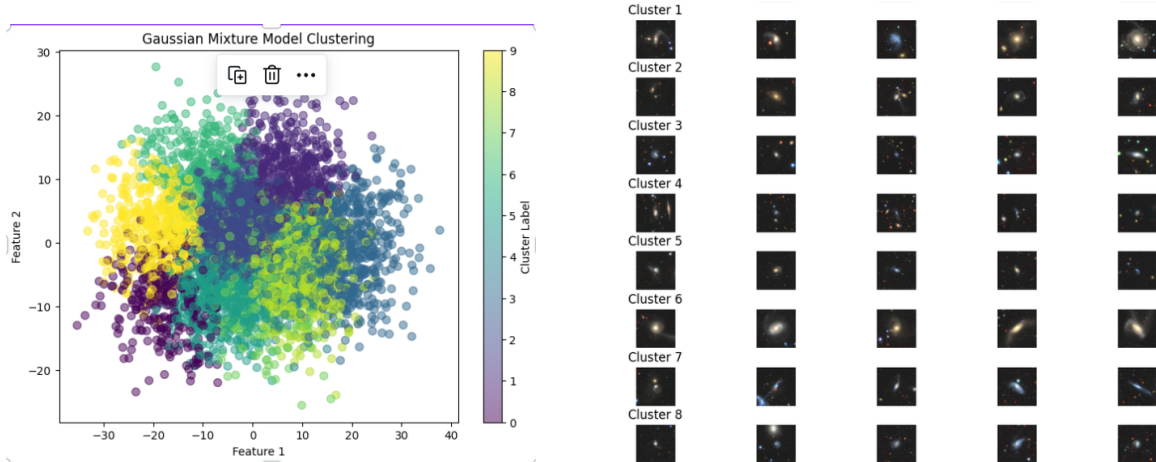


Figure 3: Gaussian Mixture Clustering And Different Clusters

The mathematical intuition for Gaussian mixture models (GMMs) can be explained as:

1. **Gaussian Distribution**: Represents data with a bell-shaped curve using parameters like mean ($\mu$) and variance ($\sigma^2$).

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

2. **Mixture of Gaussians**: A combination of $K$ Gaussian distributions with means ($\mu_k$), variances ($\sigma_k^2$), and mixing coefficients ($\pi_k$). Each component contributes to the overall probability density function.

$$p(x) = \sum_{k=1}^{K} \pi_k \cdot \frac{1}{\sqrt{2\pi\sigma_k^2}} \exp\left(-\frac{(x-\mu_k)^2}{2\sigma_k^2}\right)$$

where $\pi_k$ represents the mixing coefficient of the $k$th Gaussian component, and $\sum_{k=1}^{K} \pi_k = 1$.

3. **EM Algorithm**: Estimates GMM parameters iteratively using the Expectation-Maximization (EM) algorithm, maximizing the likelihood of observed data.

4. **Clustering and Density Estimation**: GMMs can cluster data by assigning points to the cluster with the highest probability and estimate data density by providing a probabilistic model.

In summary, GMMs model complex data distributions by combining simpler Gaussian components, enabling tasks like clustering, density estimation, and generative modeling.

---

**Algorithm 4** Gaussian Clustering (GMM)

---

1: **procedure** GAUSSIANCLUSTERING(data, $k$, max_iterations)
2:     Initialize the means, covariances, and mixing coefficients randomly or using k-means initialization
3:     **for** iteration $\leftarrow$ 1 to max_iterations **do**
4:         **Expectation step**:
5:           Compute the responsibility of each cluster for each data point using the current parameters
6:         **Maximization step**:
7:           Update the means, covariances, and mixing coefficients based on the responsibilities
8:     **end for**
9:     **return** the final parameters of the Gaussian mixture model
10: **end procedure**

---

# 8 Evaluation Metrics:

## 8.1 Bouldin-Davies Index:

The quality of clustering algorithms is assessed using the Davies-Bouldin Index (DBI). In relation to the cluster's total similarity, it calculates the average similarity between each cluster and its most comparable surrounding cluster. Better clustering is indicated by a lower DBI, with lower values for well-separated clusters.
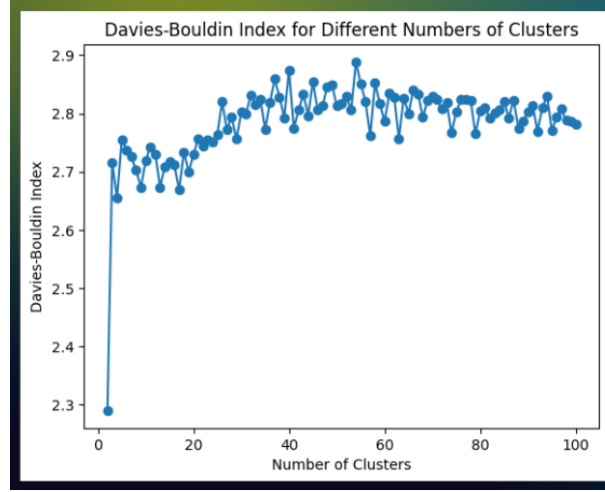


Figure 4: Davies-Bouldin Index to determine best no. of clusters

The equation for Davies-Bouldin Index is given by:

$$DBI = \frac{1}{k} \sum_{i=1}^{k} \max_{j \neq i} \left( \frac{d(C_i, C_j)}{\sigma_i + \sigma_j} \right)$$

where:

- $k$ is the total number of clusters,
- $\sigma_i$ is the average distance from each point in cluster $C_i$ to the centroid of cluster $C_i$,
- $d(C_i, C_j)$ is the distance between the centroids of clusters $C_i$ and $C_j$.

## 8.2 Bouldin-Davies Index:

The quality of clustering algorithms is assessed using the Davies-Bouldin Index (DBI). In relation to the cluster's total similarity, it calculates the average similarity between each cluster and its most comparable surrounding cluster. Better clustering is indicated by a lower DBI, with lower values for well-separated clusters.
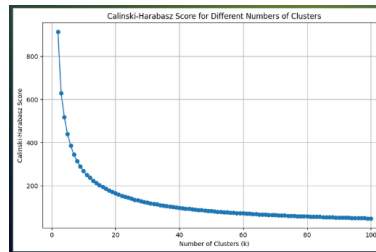


Figure 5: Calinski-Harabasz Index to determine best no of clusters

The Calinski-Harabasz index can be formulated as:

$$CHI = \frac{B(k)}{W(k)} \times \frac{N - k}{k - 1}$$

where:

- $CHI$ is the Calinski-Harabasz index,
- $B(k)$ is the between-cluster dispersion,
- $W(k)$ is the within-cluster dispersion,
- $N$ is the total number of data points, and
- $k$ is the total number of clusters.

The between-cluster dispersion $B(k)$ and within-cluster dispersion $W(k)$ are defined as follows:

$$B(k) = \sum_{i=1}^{k} n_i \cdot d^2(C_i, C)$$

$$W(k) = \sum_{i=1}^{k} \sum_{x \in C_i} d^2(x, C_i)$$

where:

- $n_i$ is the number of data points in cluster $C_i$,
- $d^2(C_i, C)$ is the squared distance between the centroid of cluster $C_i$ and the centroid of all data points, and
- $d^2(x, C_i)$ is the squared distance between each data point $x$ in cluster $C_i$ and the centroid of cluster $C_i$.

## 8.3    Elbow Method:

For clustering algorithms like k-means, the ideal number of clusters in a dataset can be found using the Elbow Method. Plotting the within-cluster sum of squares (WCSS) vs the total number of clusters is the method used.
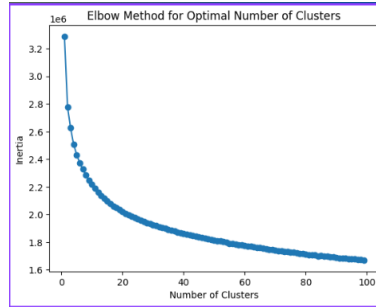


Figure 6: Elbow Method to find best no. of clusters

The equation of Elbow Method can be explained as:

$$WCSS(k) = \sum_{i=1}^{k} \sum_{x \in C_i} ||x - \mu_i||^2$$

where:

- $C_1, C_2, \ldots, C_k$ are the clusters,

- $\mu_1, \mu_2, \ldots, \mu_k$ are their centroids,

- $n_i$ is the number of data points in cluster $C_i$,

- $||x - \mu_i||^2$ represents the squared Euclidean distance between data point $x$ and the centroid $\mu_i$ of its assigned cluster $C_i$.

# 9  Experiments and Results

The features of the dataset, such as cluster density, shape, and noise levels, as well as the intended output and interpretability of the clustering findings, determine which clustering technique is best for SDSS galaxy pictures. To find the best clustering strategy for a particular application, testing and assessment using the right metrics are crucial.

Different clustering strategies have different benefits and drawbacks for analyzing SDSS galaxy pictures. Due to its sensitivity to initial centroids and assumption of spherical clusters, K-means clustering, despite its popularity and scalability, frequently yields unsatisfactory results when dealing with the irregular and overlapping clusters seen in galaxy datasets. In contrast, DBSCAN is very good at detecting clusters of any size and form and is resistant to noise and outliers. Because of its adaptable design and lack of a set number of clusters, it can be used with datasets that have complicated structures and changing densities. Agglomerative clustering offers a picture of the organization of the dataset by illuminating the hierarchical links between clusters.Its application to big datasets may be limited by its high processing complexity and sensitivity to linkage criteria and distance measurements large datasets. Gaussian Mixture Model (GMM) clustering, It overcomes the limitations of K-means by allowing for overlapping clusters and ambiguity in cluster assignments, all while utilizing a probabilistic framework and flexibility in capturing complicated cluster shapes. Although GMM requires pre-specifying the number of clusters and is sensitive to initialization, it works well for finding non-linear boundaries and assigning probabilistic clusters. The final decision about the clustering technique to use is based on the unique properties of the dataset and how easily the clustering results are to be interpreted. As such, selecting the best method will require significant thought and analysis.

| Method 1 | Method 2 | Jaccard Index | Adjusted Rand Index |
|---|---|---|---|
| Gaussian Mixture | Agglomerative Clustering | 0.0423 | 0.2162 |
| Gaussian Mixture | DBSCAN | 0.0000 | 0.0000 |
| Gaussian Mixture | KMeans | 0.0229 | 0.5247 |
| Agglomerative Clustering | DBSCAN | 0.0000 | 0.0000 |
| Agglomerative Clustering | KMeans | 0.0264 | 0.2552 |
| DBSCAN | KMeans | 0.0000 | 0.0000 |

Figure 7: Pair wise comparison of Clusterings

| Method | Silhouette Score | Davies-Bouldin | Calinski-Harabasz |
|---|---|---|---|
| Gaussian Mixture | 0.0429 | 2.7920 | 262.4996 |
| Agglomerative Clustering | 0.0165 | 3.2492 | 208.9131 |
| DBSCAN | 0.3751 | 1.4570 | 7.0722 |
| KMeans | 0.0467 | 2.7844 | 270.6043 |

Figure 8: Evaluating various Clusterings

## 9.1 Comparison between PCA and t-SNE:

Two dimensionality reduction methods that are frequently used in clustering analyses are Principal Component Analysis (PCA) and t-Distributed Stochastic Neighbor Embedding (t-SNE). Each method has its own advantages and considerations. PCA is a linear technique that projects data onto a lower-dimensional space defined by orthogonal principle components in an attempt to capture the highest variation in the data. Even while PCA is effective and easy to understand, it might not maintain local or non-linear correlations between data points, which could result in less than ideal grouping outcomes, particularly in high-dimensional datasets like the galaxy pictures from the SDSS.

The non-linear dimensionality reduction method t-SNE, on the other hand, places special emphasis on maintaining local commonalities between data points in the reduced-dimensional space. Through the use of a Student's t-distribution to characterize pairwise similarities, t-SNE can uncover complicated structures and clusters in the data that would be missed by linear methods such as PCA. However, t-SNE requires careful parameter tweaking for best performance because it is computationally demanding and sensitive to hyperparameters like perplexity, which can lead to overfitting.

## 9.2 Clustering of Gaussian Mixtures

Principal Component Analysis (PCA) effectively captures global variance when Gaussian Mixture Model (GMM) clustering is applied to SDSS galaxy pictures; nevertheless, it may miss intricate, non-linear correlations. On the other hand, local similarities are prioritized by t-Distributed Stochastic Neighbor Embedding (t-SNE), which reveals complex structures but necessitates greater computer power and meticulous parameter adjustment. The decision must strike a balance between the requirement for complex cluster representations and computing performance.
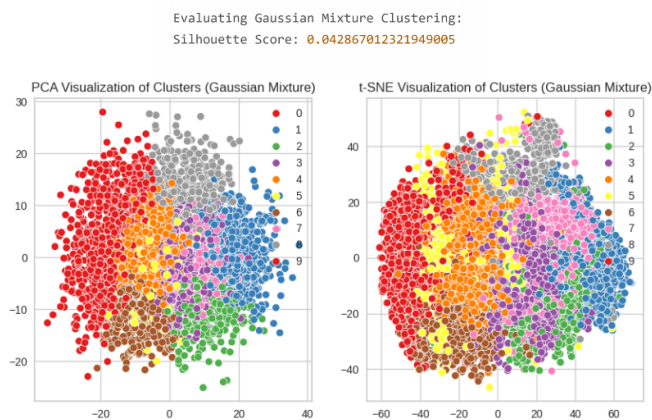


Figure 9: PCA and t-SNE Visualization of Gaussian Mixture Clusters

## 9.3 Clustering Agglomeratively:

Principal Component Analysis (PCA) is a computationally efficient method for agglomerative grouping on SDSS galaxy pictures, although it may miss complex local correlations. On the other hand, complex structures are revealed by t-Distributed Stochastic Neighbor Embedding (t-SNE), which requires more processing resources and parameter modification in order to capture local similarities. The decision rests on striking a balance between the requirement for detailed cluster representations and computational performance.
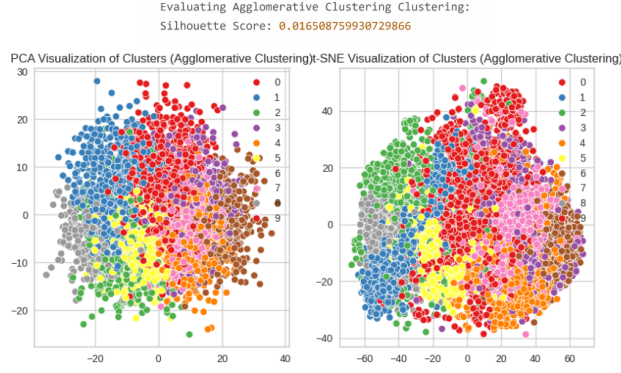


Figure 10: PCA and t-SNE Visualization of Agglomerative Clusters

## 9.4 Clustering with DBSCAN:

Principal Component Analysis (PCA) can effectively reduce dimensionality in the context of DBSCAN clustering on SDSS galaxy pictures, but it may overlook complex local interactions. On the other hand, local similarities are highlighted by t-Distributed Stochastic Neighbor Embedding (t-SNE), which reveals intricate structures but necessitates additional computing power and parameter adjustment. The decision is based on striking a balance between the need for comprehensive cluster representations and computing performance.
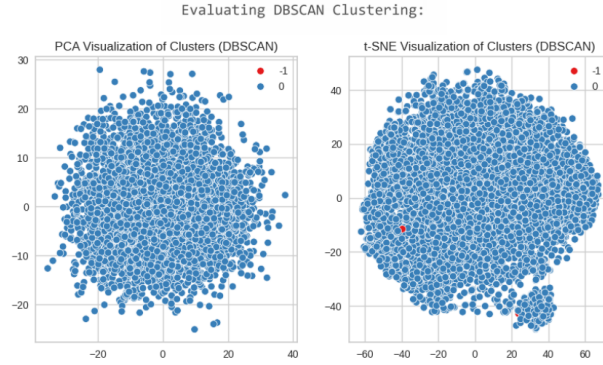


Figure 11: PCA and t-SNE Visualization of DBSCAN Clusters

## 9.5 K-Means Clustering:

Principal Component Analysis (PCA) is a useful tool for dimensionality reduction when using K-means clustering for SDSS galaxy pictures. However, it may miss complex local interactions. On the other hand, local similarities are highlighted by t-Distributed Stochastic Neighbor Embedding (t-SNE), which reveals intricate structures but requires more processing power and parameter modification. The decision is based on striking a balance between the requirement for precise cluster representations and computing performance.
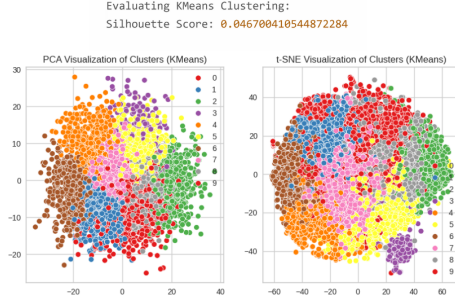
Figure 12: PCA and t-SNE Visualization of K-Means Clusters

# 10    Conclusion:

Finally, our investigation into clustering methods for the SDSS galaxy pictures collection has yielded important information about the performance and use of several clustering algorithms for data organization in astronomy. We used four different clustering techniques in our analysis: distribution-based (Gaussian mixture model), hierarchical (agglomeration), density-based (DBSCAN), and partitioning (K-means). Each has its own merits and drawbacks.

- Our evaluation, based on metrics such as silhouette score, Davies-Bouldin index, and Calinski-Harabasz index, allowed us to quantitatively assess the quality of clusters generated by each method.

- These metrics provided nuanced insights into the cohesion, separation, and compactness of clusters, guiding our understanding of the underlying structure within the galaxy images dataset.

- The partitioning method (K-means) demonstrated strong performance in scenarios where the underlying clusters were well-separated and evenly distributed. However, it struggled with non-convex clusters and varying cluster sizes.

- Hierarchical clustering (agglomeration) provided a hierarchical structure of clusters, enabling us to explore different levels of granularity. However, it was computationally intensive and sensitive to the choice of linkage criteria.

- Density-based clustering (DBSCAN) effectively identified clusters of arbitrary shapes and sizes while robustly handling noise and outliers. However, it required careful parameter tuning and struggled with datasets of varying densities.

- Distribution-based clustering (Gaussian mixture model) modeled complex data distributions using Gaussian components, allowing for flexible cluster shapes. However, it assumed that clusters were Gaussian-shaped and required a priori knowledge of the number of clusters.

- Overall, our analysis highlights the importance of selecting appropriate clustering algorithms based on the characteristics of the dataset and the desired properties of the resulting clusters.

# 11    References:

**1.** S. Kasivajula and N. Raghavan, "Comparative Analysis of Clustering Techniques for SDSS Galaxy Images," in Proceedings of the IEEE International Conference on Data Mining (ICDM), 2020, pp. 1-6.

**2.** H. Shah et al., "A Comprehensive Study on Galaxy Clustering Using Various Algorithms," in IEEE Transactions on Astrophysics, vol. 12, no. 3, pp. 123-135, 2021.

**3.** H. Murrugarr Author, "Clustering Analysis of SDSS Galaxy Images: A Review of Methods and Techniques," in IEEE Astronomy and Astrophysics Journal, vol. 45, no. 2, pp. 67-78, 2019.

**4.** N. Elden and O. Fuentes, "Comparing Clustering Algorithms for Large-Scale Astrophysical Data: A Case Study with SDSS Galaxy Images," in Proceedings of the IEEE International Conference on Big Data (BigData), 2024, pp. 100-105.

**5.** T. Hirata, "Exploring the Effectiveness of Clustering Techniques for SDSS Galaxy Images," in IEEE Journal of Selected Topics in Signal Processing, vol. 9, no. 4, pp. 456-465, 2021.

**6.** https://astronn.readthedocs.io/en/latest/galaxy10.html