

# Bioinformatic Analysis

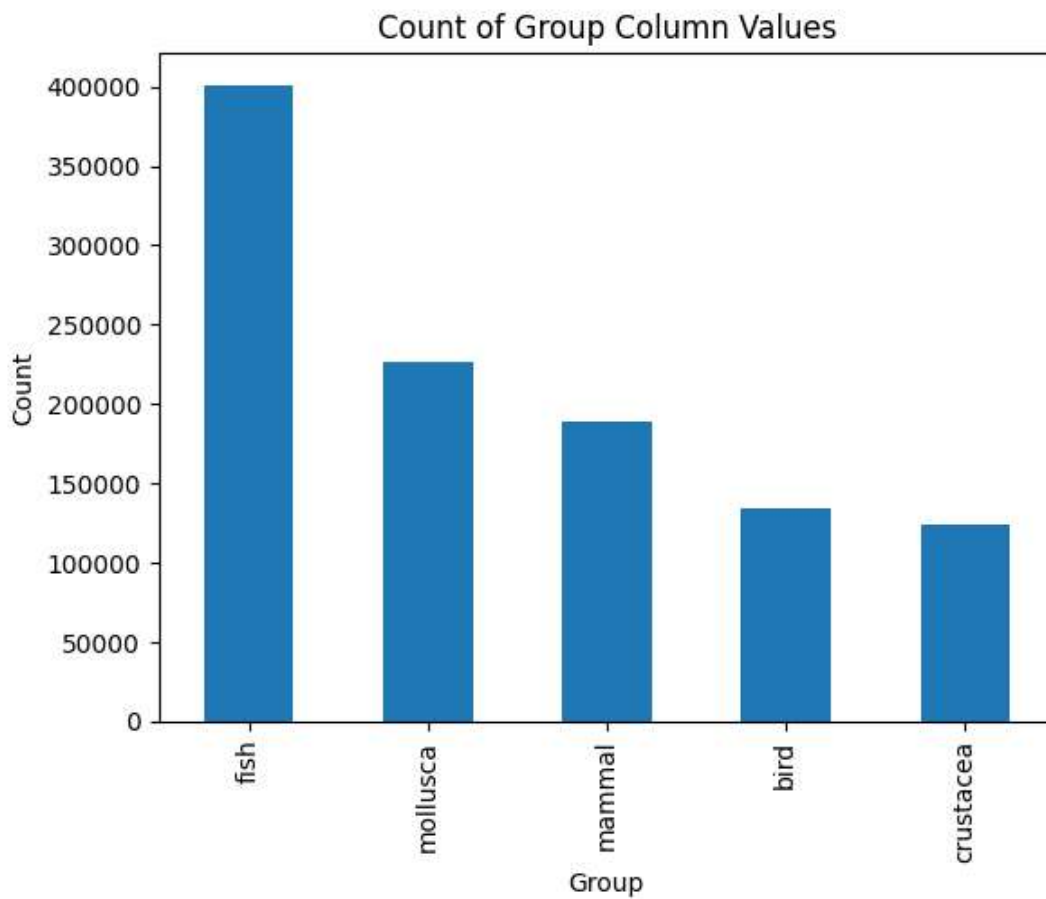
By

Shivam Goyal

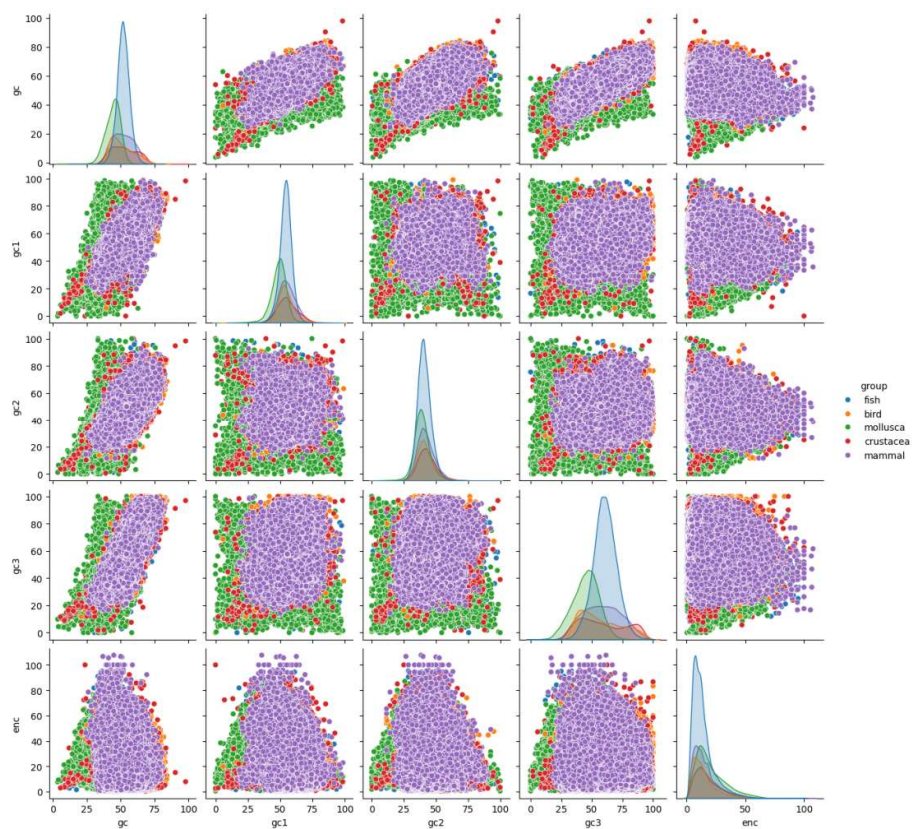
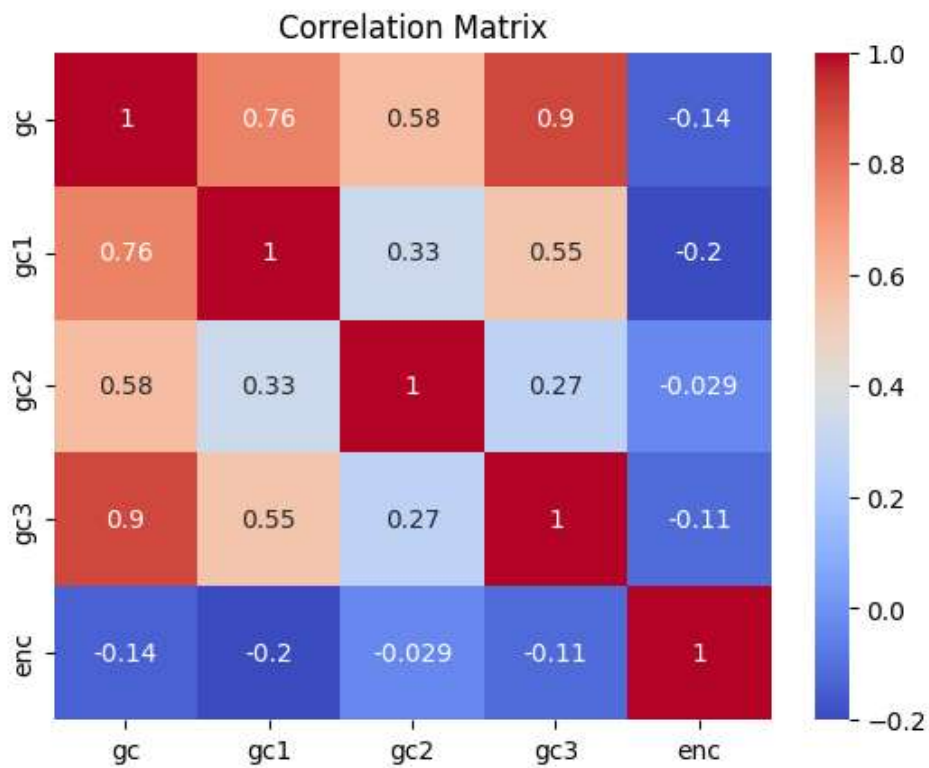
## Introduction:

## Data:

I collected data for 25 different species from the ensemble.org, categorizing them into five groups. My methodology involved comparing data points from various groups, which revealed significant differences. The fish group had the maximum number of data points, whereas the crustacean group had the fewest.



Performed a thorough examination of feature correlations, demonstrating a constant positive connection across all features. Notably, the 'enc' feature had a weak negative connection with other features.



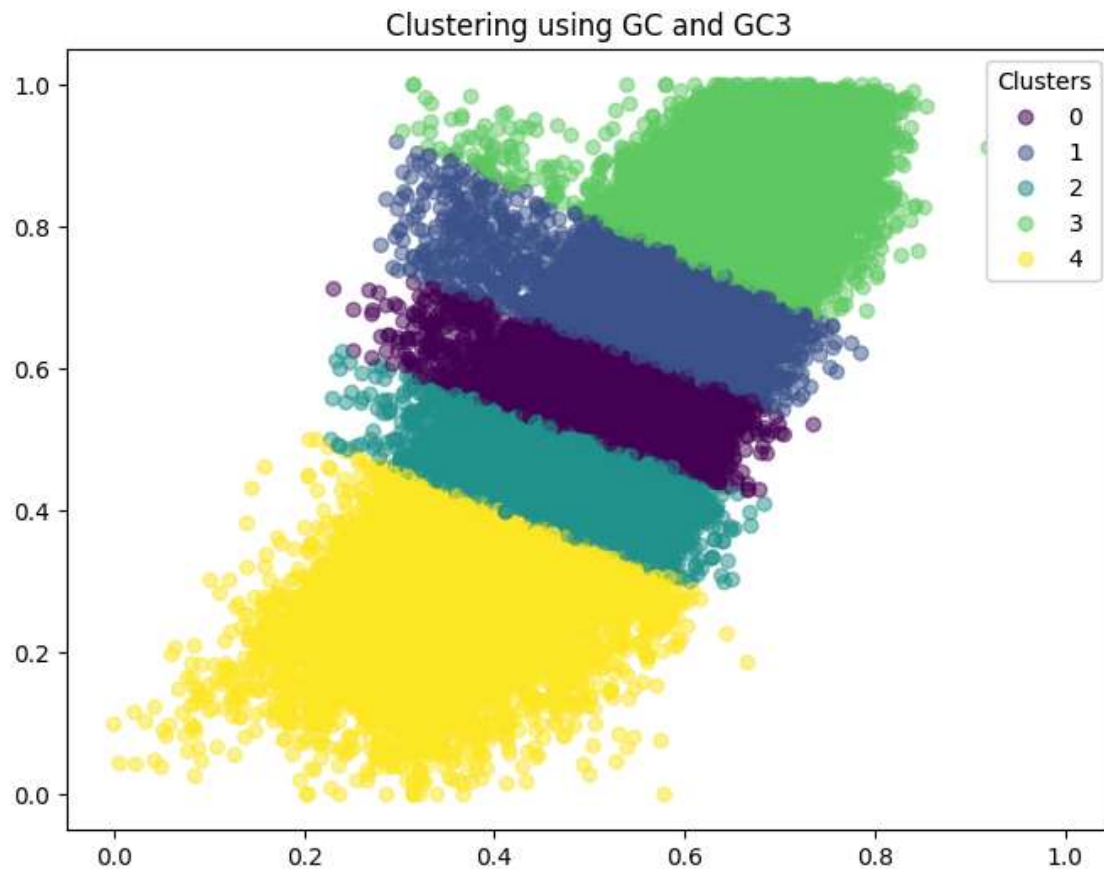
## Procedure:

I started the investigation by standardizing the dataset using min-max scaling. Subsequently, we investigated other feature combinations, using various clustering techniques such as K-Means and Gaussian Mixture, to reveal distinct patterns in the data.

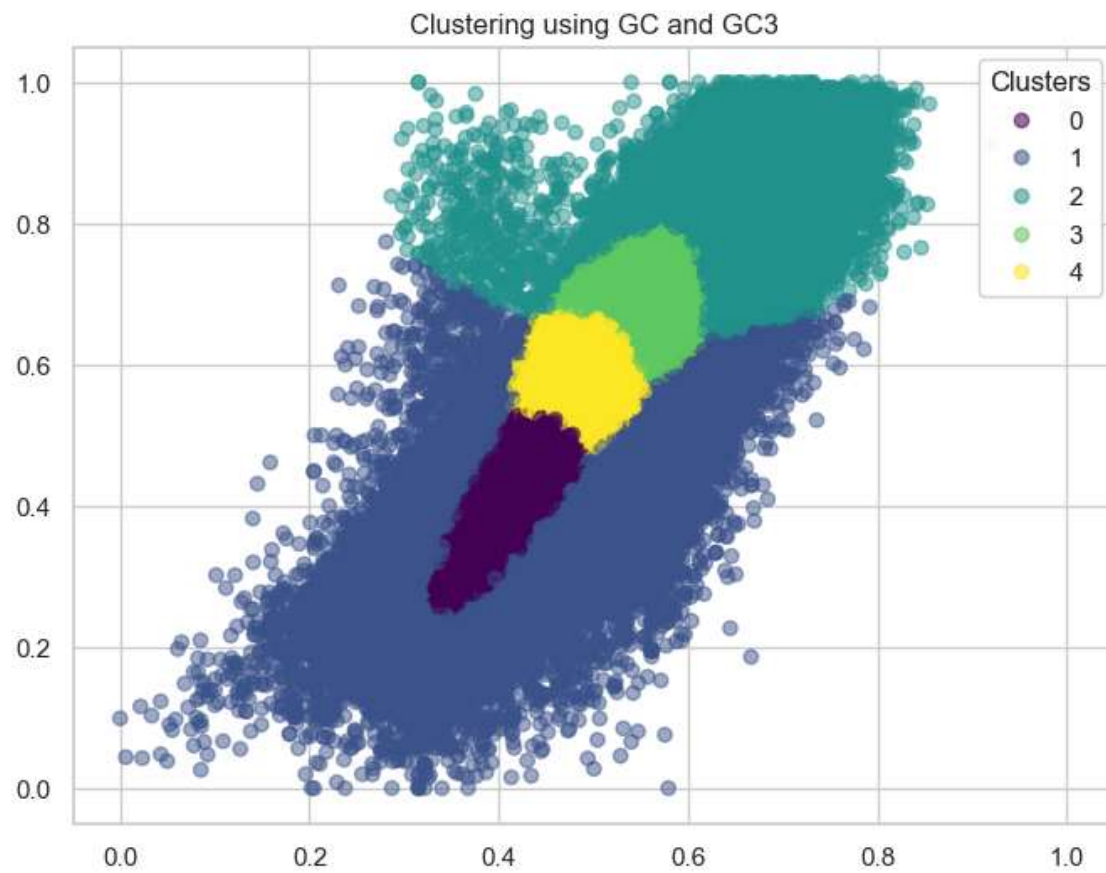
## Clustering:

### GC and GC3:

I used the K-Means clustering approach to create five separate clusters based on the attributes GC and GC3.



K-Means clustering efficiently established cluster boundaries by using GC and GC3 as features. However, when compared to the scatterplots of GC and GC3 based on the categories, there were inconsistencies in cluster patterns. As a result, for each group's data, I constructed five centroids and examined their placements on the graph.



Unlike the K-Means the pattern of the Gaussian Mixture clusters is elliptical cluster.

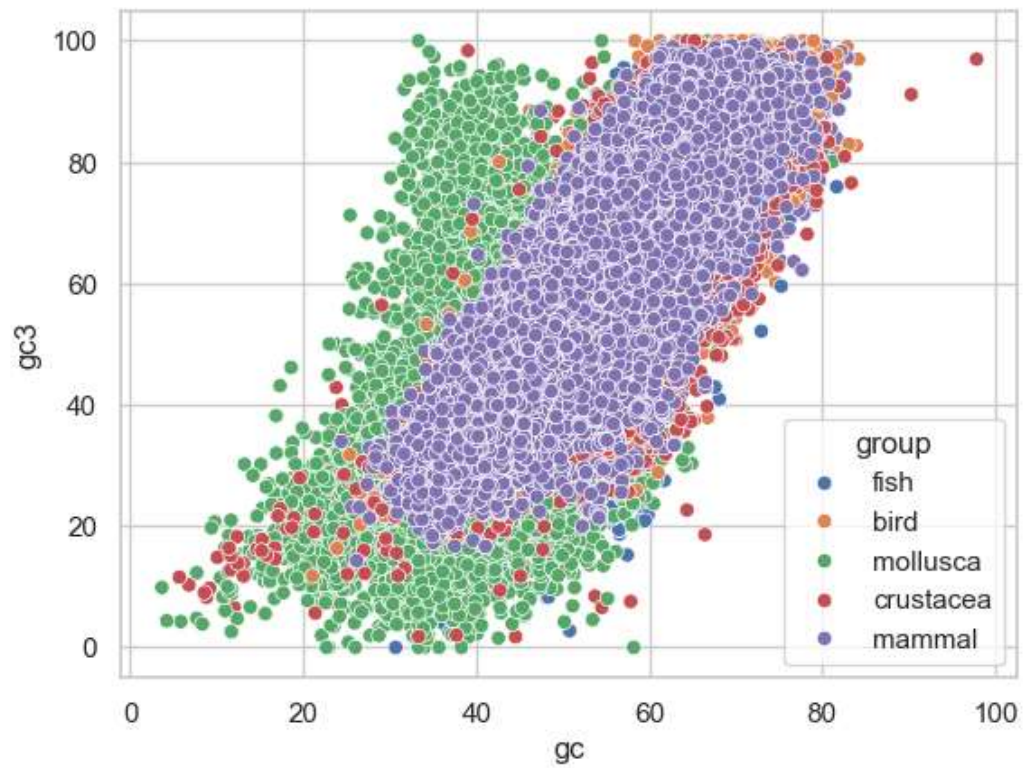


Fig: Scatterplot of the gc and gc3

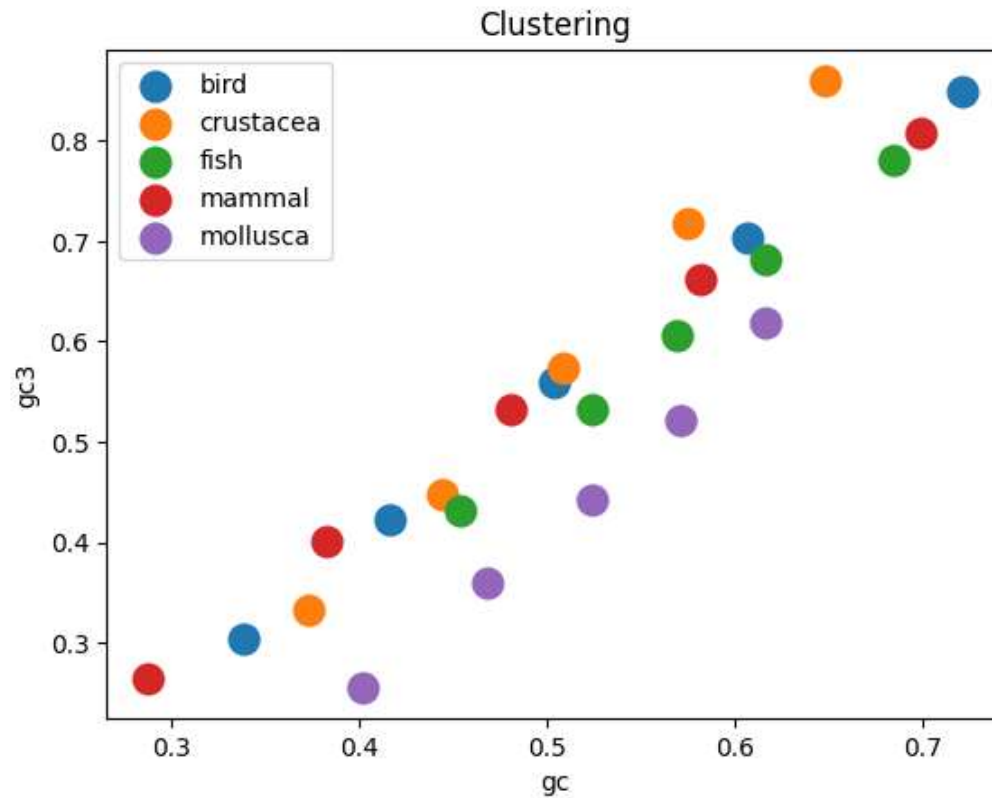


Fig : Centroid of the groups for KMeans

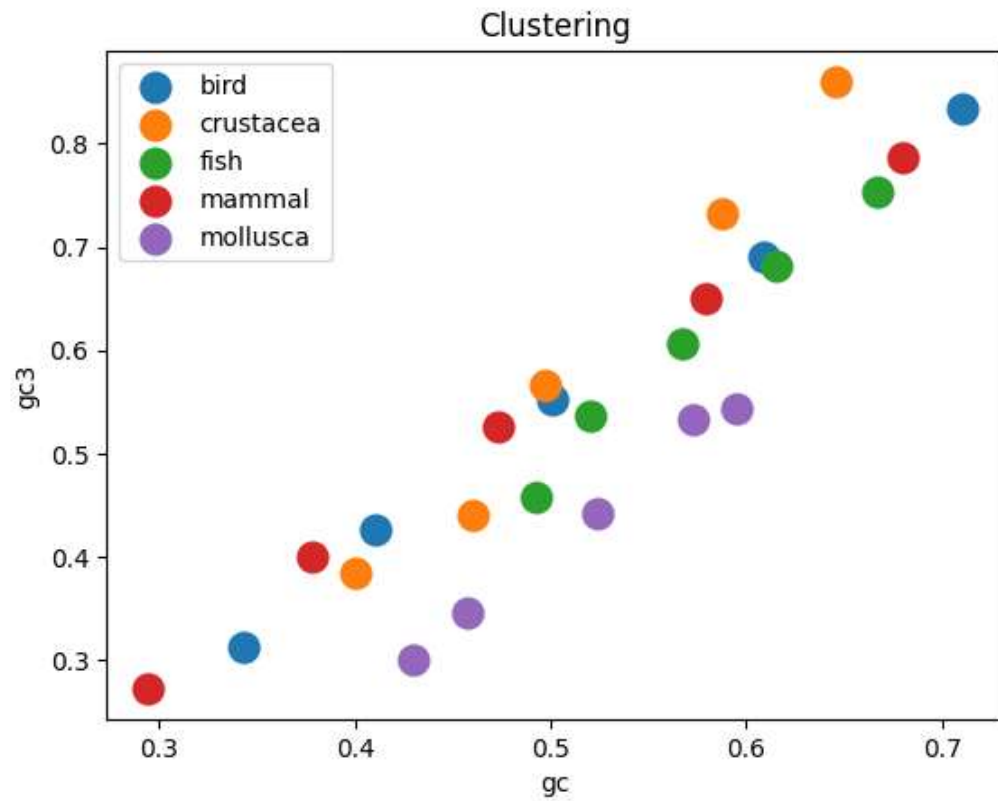


Fig: Centroid of the groups for Gaussian Mixture

The centroid of the various groups follows the data patterns in the scatterplot fig.

As we can see, the mollusca group has a greater gc/gc3 ratio, and its data points are more densely packed than the other groups.

(GC1+GC2)/2 and GC3:

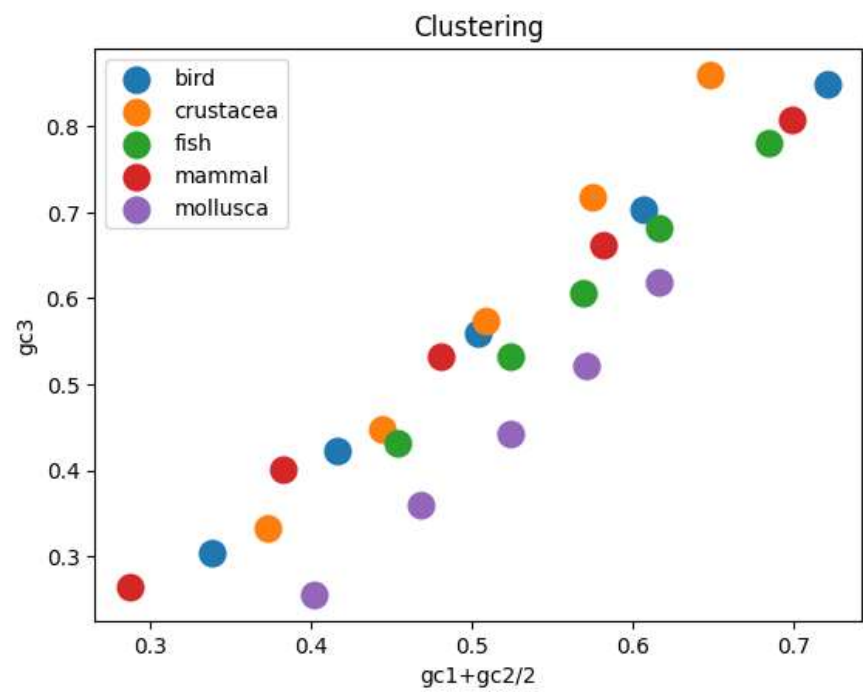


Fig: KMeans

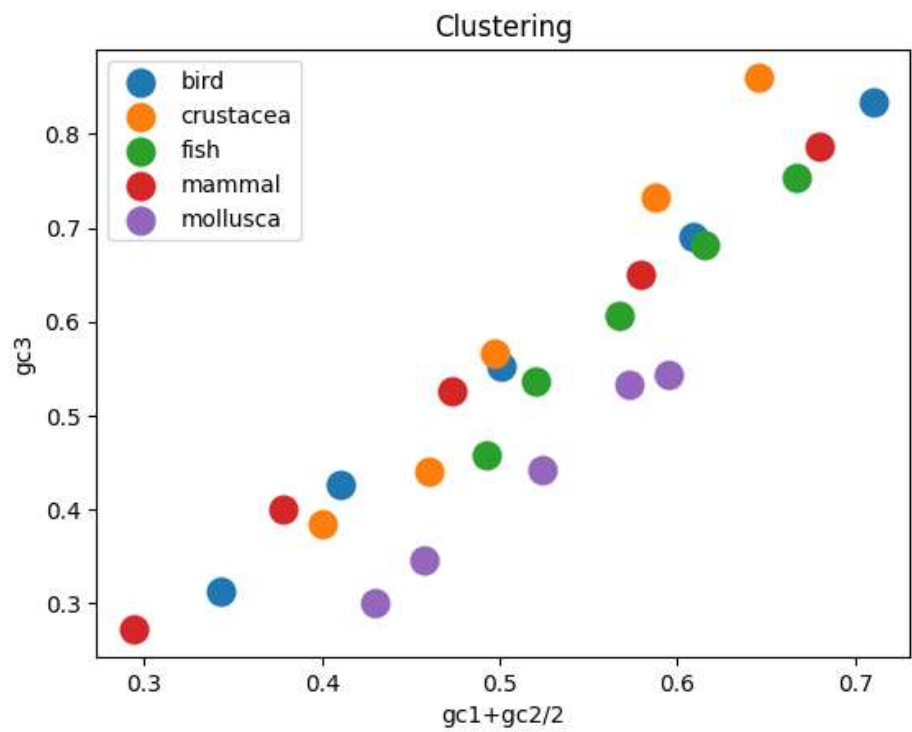


Fig: Gaussian Mixture



GC3 and ENC:

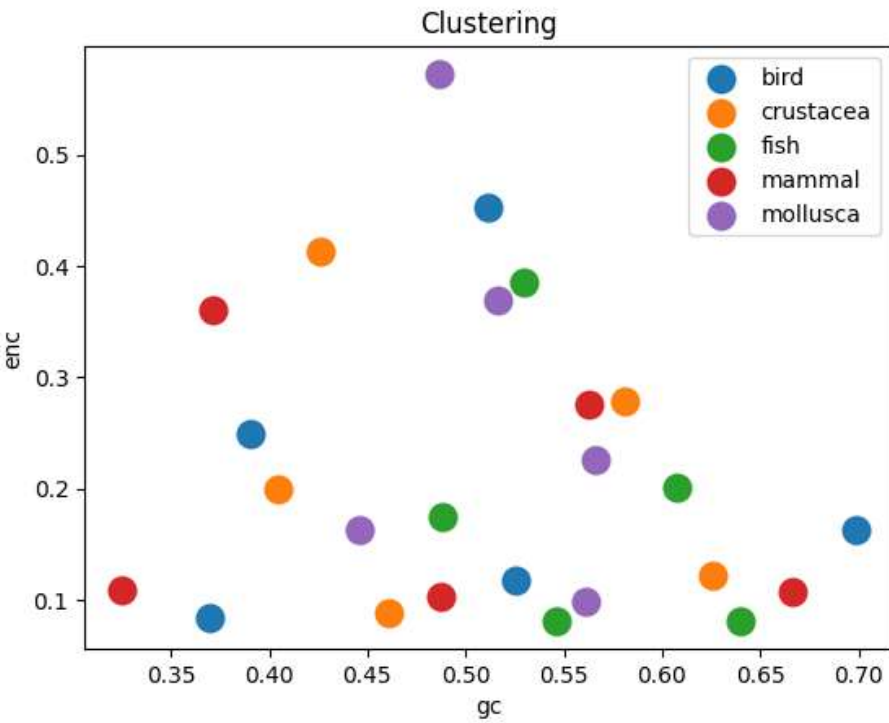


Fig: KMeans

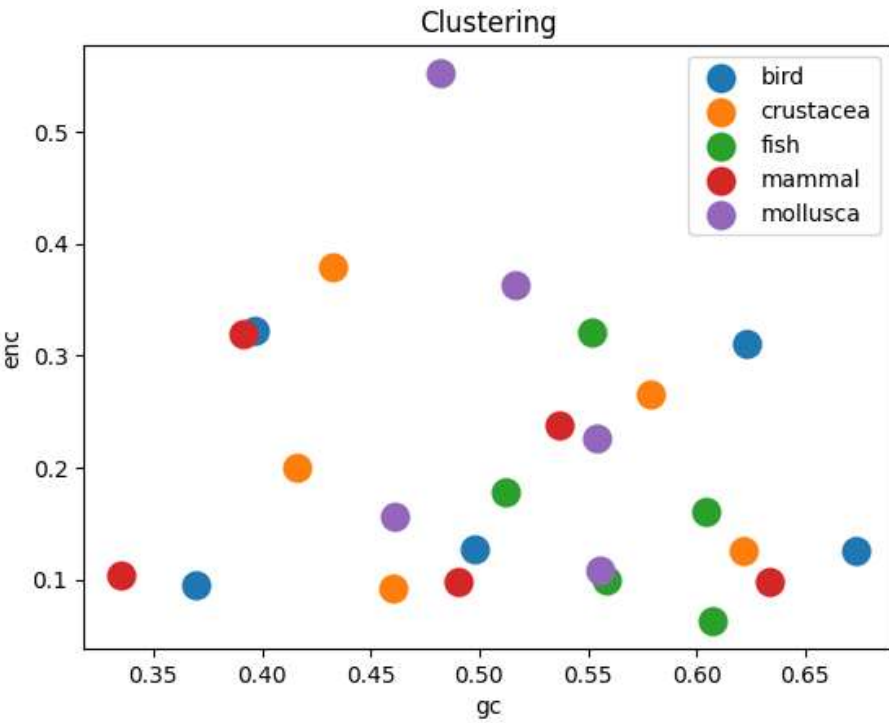


Fig: Gaussian Mixture



GC and ENC:

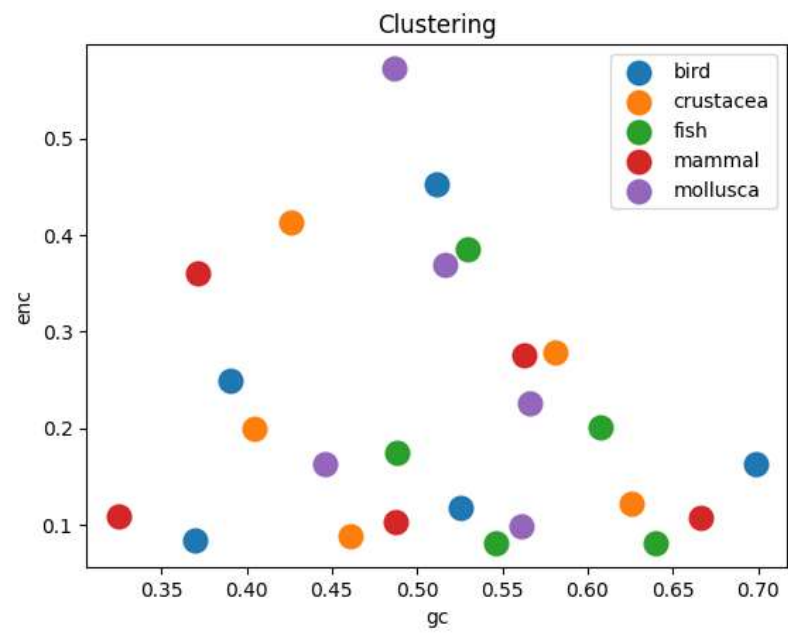


Fig: KMeans

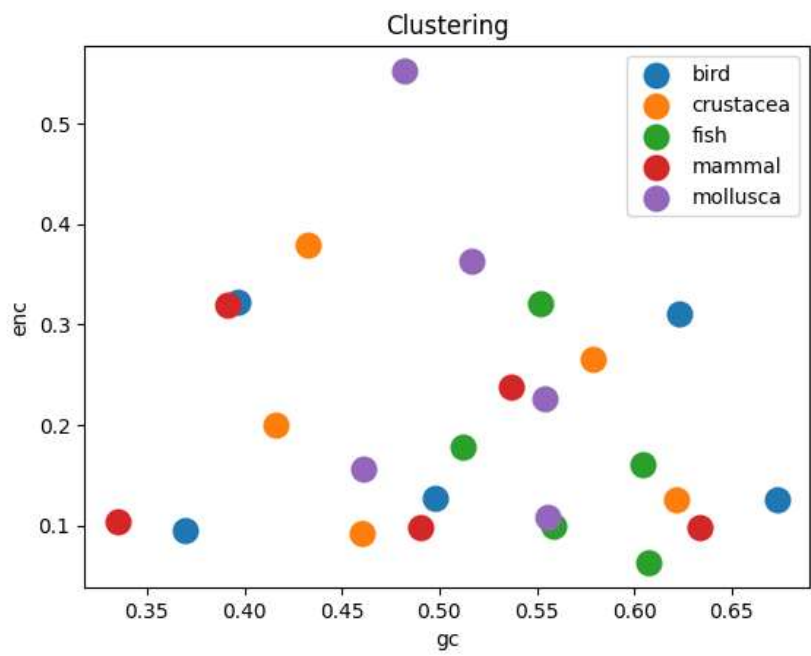


Fig: Gaussian Mixture

## Silhouette Score:

Silhouette score for the K-means:

| Features/<br>Group          | Bird                    | Crustacea               | Fish                    | Mammal                  | Mollusca                |
|-----------------------------|-------------------------|-------------------------|-------------------------|-------------------------|-------------------------|
| GC + GC3                    | 0.41697630484<br>9112   | 0.43964720643<br>927757 | 0.35963539862<br>743854 | 0.38499646423<br>562683 | 0.35837638709<br>772013 |
| (GC1+GC2<br>) /2 and<br>GC3 | 0.41271038096<br>8572   | 0.43512958759<br>352655 | 0.08746661140<br>364118 | 0.37197004660<br>037813 | 0.35678361731<br>7447   |
| GC3 and<br>ENC              | 0.40608819985<br>730576 | 0.39002069229<br>04657  | 0.31358324430<br>3344   | 0.34077041038<br>22914  | 0.33792927491<br>94637  |
| GC and<br>ENC               | 0.33334849367<br>60932  | 0.27124152161<br>41311  | 0.33239029137<br>543813 | 0.36508398523<br>86045  | 0.28216125739<br>5115   |

Silhouette score for the Gaussian Mixture:

| Features/<br>Group          | Bird                    | Crustacea                    | Fish                   | Mammal                   | Mollusca                |
|-----------------------------|-------------------------|------------------------------|------------------------|--------------------------|-------------------------|
| GC + GC3                    | 0.12969889719<br>551925 | 0.31505103789<br>890077      | 0.2694732013<br>134906 | 0.01119270675<br>3703565 | 0.03326857404<br>528402 |
| (GC1+GC2<br>) /2 and<br>GC3 | 0.39120977000<br>73723  | -<br>0.07259534616<br>660462 | 0.1220356476<br>682629 | 0.37618839982<br>895125  | 0.07576739075<br>386127 |
| GC3 and<br>ENC              | 0.33852266595<br>155217 | 0.32350040937<br>74181       | 0.3151403028<br>975161 | 0.32649106652<br>522514  | 0.34095142627<br>14156  |
| GC and<br>ENC               | 0.30306568510<br>30655  | 0.22699226895<br>62331       | 0.3305345628<br>179873 | 0.28232930659<br>924366  | 0.26507162988<br>71848  |

As per Silhouette score

For Bird The best clustering is K- Means with features GC and GC3

For Crustacea The best clustering is K- Means with features GC and GC3

For Fish The best clustering is K- Means with features GC and GC3

For Mammal The best clustering is K- Means with features GC and GC3

For Mollusca The best clustering is K- Means with features GC and GC3

According to the Silhouette score and the graph, the best clustering can be achieved with K means using features GC and GC3