





```
In [63]: ms_lst = []
for b in bandwidths:
    ms = MeanShift(bandwidth=b)
    ms.fit(df_g_pca)

    try:
        sil_score = silhouette_score(df_g_pca, ms.labels_, metric='euclidean')
        ch_score = calinski_harabasz_score(df_g_pca, ms.labels_)
    except:
        sil_score = 0
        ch_score = 0
    catgeg = pd.Series(ms.labels_).nunique()
    ms_lst.append([b, catgeg, adjusted_rand_score(kingdoms, ms.labels_), sil_score, ch_score])

ms_df = pd.DataFrame(ms_lst, columns=['bandwidth', 'clust_no', 'adj_rand_score', 'sil_score', 'ch_score'])

In [64]: ms_df["quantile"] = quantiles
```

Again, different bandwidth values optimize different metrics. Adjusted Rand Score is the highest for the model with 92 clusters (0.14) - this model would be hard to use and interpret, so we will pay more attention to other metrics. Silhouette Coefficient is the highest for models with 2 clusters (0.25), while Calinski-Harabasz score is the highest for a 7-cluster model (686).

```
In [66]: ms_df

Out[66]:
```

	bandwidth	clust_no	adj_rand_score	sil_score	ch_score	quantile
0	4.038823	92	0.142449	0.004683	196.385013	0.02
1	4.576521	38	0.095270	0.014293	265.823592	0.04
2	4.977182	19	0.080759	0.047918	433.492867	0.06
3	5.285271	14	0.066663	0.068446	546.211831	0.08
4	5.548612	7	0.014155	0.093321	686.239630	0.10
5	5.785990	7	0.014175	0.094213	684.974095	0.12
6	6.004099	5	0.000849	0.045111	180.613098	0.14
7	6.203698	5	0.000832	0.043329	175.593573	0.16
8	6.388517	3	-0.002696	0.126270	200.481904	0.18
9	6.562058	2	0.000438	0.251856	25.768293	0.20
10	6.727673	2	0.000438	0.251856	25.768293	0.22
11	6.884945	1	0.000000	0.000000	0.000000	0.24
12	7.037239	1	0.000000	0.000000	0.000000	0.26
13	7.185963	1	0.000000	0.000000	0.000000	0.28
14	7.330660	1	0.000000	0.000000	0.000000	0.30

```
In [67]: max_rand_ms = ms_df[ms_df.adj_rand_score == np.max(ms_df.adj_rand_score)]

In [68]: max_sil_ms = ms_df[ms_df.sil_score == np.max(ms_df.sil_score)]

In [69]: max_ch_ms = ms_df[ms_df.ch_score == np.max(ms_df.ch_score)]

In [70]: max_ms = pd.concat([max_rand_ms, max_sil_ms, max_ch_ms])

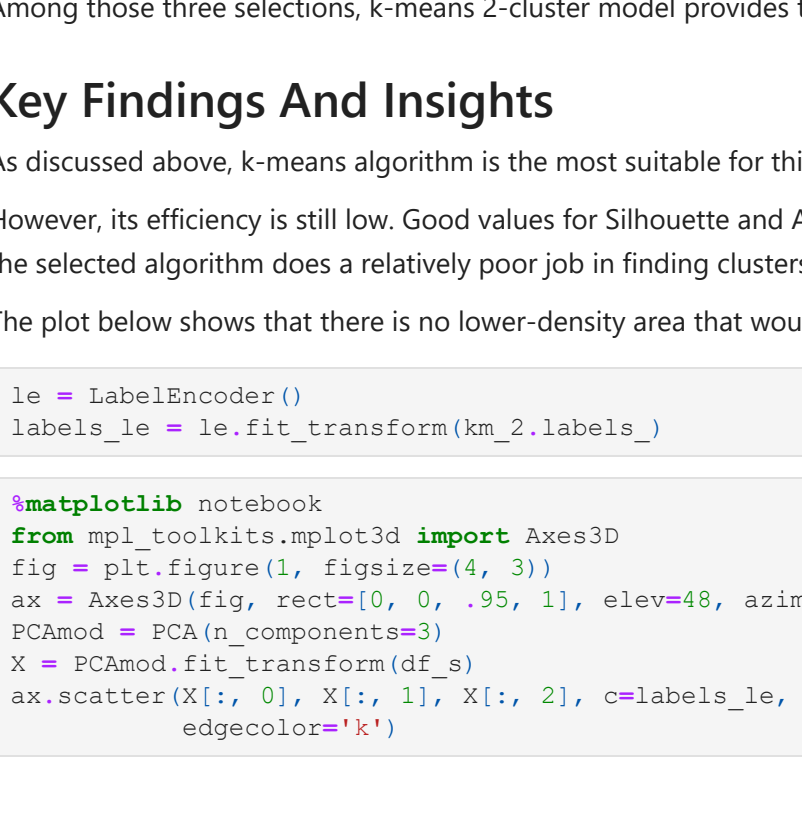
In [71]: max_ms

Out[71]:
```

	bandwidth	clust_no	adj_rand_score	sil_score	ch_score	quantile
0	4.038823	92	0.142449	0.004683	196.385013	0.02
9	6.562058	2	0.000438	0.251856	25.768293	0.20
10	6.727673	2	0.000438	0.251856	25.768293	0.22
4	5.548612	7	0.014155	0.093321	686.239630	0.10

We selected 2-cluster models for the previous algorithms, so we will choose it for Mean Shift as well for comparability. The plot below shows that while Adjusted Rand Score increases with the number of clusters (which shows improvement), Calinski-Harabasz and Silhouette scores decrease (a sign of worse performance) after the peak cluster number:

```
In [112]: ax = ms_df.plot(x="clust_no", y="ch_score", legend=False, color="g")
ax2 = ax.twinx()
ms_df.plot(ax=ax, x="clust_no", y=["adj_rand_score", "sil_score"], ax=ax2, legend=False)
ax.figure.legend()
plt.show()
```



Now let's analyze the 2-cluster model that gave the highest Silhouette Coefficient:

```
In [98]: ms_02 = MeanShift(bandwidth=bandwidths[9], bin_seeding=False)
ms_02 = ms_02.fit(df_g_pca)
```

The vast majority of samples fall under cluster 0, only 25 of them are clustered separately. Both clusters contain organisms from all kingdoms:

```
In [99]: pd.DataFrame(ms_02.labels_).value_counts()

Out[99]:
```

	0	1
dtype:	int64	

```
In [101]: pd.crosstab(kingdoms, ms_02.labels_)
```

	col_0	0	1
Kingdom			
arch	1625	5	
mon	3051	10	
pln	1515	8	
vri	3049	2	

## Model Selection

In the analysis above, we tried three different clustering algorithms - k-means, Agglomerative Clustering and Mean Shift. Results for different parameters were compared for each algorithm and the best one (or their combination) selected.

For k-means, 2-cluster algorithm was the best selection, with Silhouette Coefficient of 0.32, Calinski-Harabasz score of 5891 and Adjusted Rand Score of 0.08.

For Agglomerative Clustering, 2-cluster ward linkage model was the best, with Silhouette Coefficient of 0.25, Calinski-Harabasz score of 4509 and Adjusted Rand Score of 0.05.

For Mean Shift, 2-cluster model was again selected. Its Silhouette Coefficient is 0.25, Calinski-Harabasz score - around 26, and Adjusted Rand Score - close to 0.

Among those three selections, k-means 2-cluster model provides the best results and is thus selected as the final one.

## Key Findings And Insights

As discussed above, k-means algorithm is the most suitable for this analysis.

However, its efficiency is still low. Good values for Silhouette and Adjusted Rand scores would be close to 1. Low values demonstrate that the selected algorithm does a relatively poor job in finding clusters.

The plot below shows that there is no lower-density area that would separate the clusters obtained with k-means model:

```
In [75]: le = LabelEncoder()
labels_le = le.fit_transform(km_2.labels_)
```

```
In [77]: %matplotlib notebook
from mpl_toolkits.mplot3d import Axes3D
fig = plt.figure(1, figsize=(4, 3))
ax = Axes3D(fig, rect=[0, 0, .95, 1], elev=40, azim=134)
PCAnod = PCA(n_components=3)
X = PCAnod.fit_transform(df_g)
ax.scatter(X[:, 0], X[:, 1], X[:, 2], c=labels_le, cmap=plt.cm.nipy_spectral,
           edgecolor='k')
```



This is in line with our previous observation of overlapping organism kingdoms.

Although clustering power is low, our clusters still make some difference between organism kingdoms. For example, k-means clusters have a very different number of viruses (2601 vs 450), while the counts for other kingdoms are very similar. A similar pattern can be observed for the Agglomerative Clustering model, where one cluster contains much more virus samples (2029 vs 1022), while the other cluster of the two has more animals (1185 vs 445) and monomers (1900 vs 1161).

Another similarity between the results for those two algorithms is a similar cluster size. k-means algorithm split the dataset into clusters of 5742 and 3523, while Agglomerative Clustering provides counts of 4870 and 4395. Mean Shift clusters almost all samples into one cluster, keeping only 25 entries separately (so it has done the worst job in clustering).

## Suggestions

Relatively poor results of our clustering models might encourage us to say that codon information is not suitable for organism clustering. However, it is worth considering some possible model improvements, before making the final conclusion.

Firstly, we used only genomic DNA data, skipping other types. The analysis might be repeated with DNA data of all types. Also, it might be worth adding some other features, that might cluster organisms well together with codon frequencies.

Ground truth labels (in our case - kingdoms) is a result of a somewhat arbitrary choice. They were remade to keep taxonomic consistency, but we might have used the initial labels. After the initial analysis, we might regroup the samples into e.g. viruses and all the rest.

Also, we selected 12 PCA components (given 64 codon features). Changing the number of components might yield different results. Reducing the component number will make us lose some information, and using more components will make more susceptible to "the curse of dimensionality", so experimenting with several numbers might be useful.

Furthermore, using other clustering algorithms is also worth trying, depending on their properties. For example, Gaussian mixture models are not scalable, but given the size of our dataset (about 10000 entries - not too large), might be interesting to experiment with.