

01. Cluster Analysis

A Machine Learning technique

02. Applications of Clustering

Social Network Analysis, Costumer Segmentation, ...

03. Clustering Methods

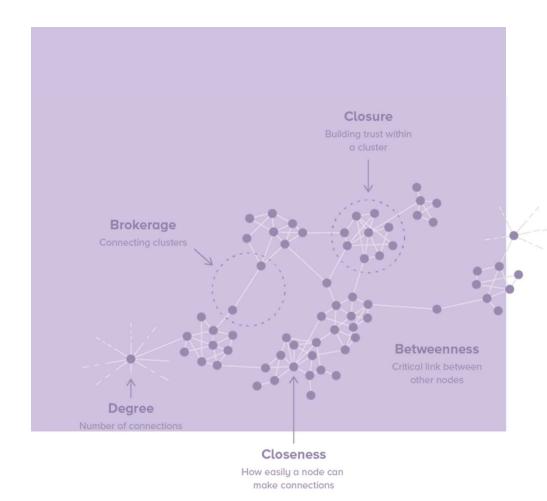
Based on Partition, Density, Hierarchy, ...

04. Tuning Clustering

Sillhoutte method, AIC, BIC, ...

05. Clustering Metrics

Separation Index – FD index



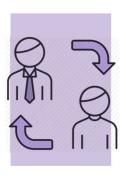
01 Cluster Analysis

Clustering is a Machine Learning technique that involves the grouping of data points.

Given a set of data points, we can use a clustering algorithm to classify each data point into a specific group.

O2 Applications of Clustering

Applications of Clustering







Social Network Analysis

The clustering in social network requires grouping objects into classes based on their links as well as their attributes.

Costumer Segmentation

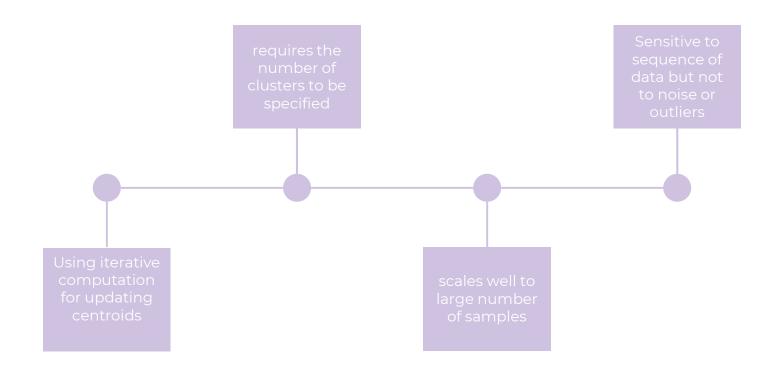
Customer Segmentation is the subdivision of a market into discrete customer groups that share similar characteristics.

Identifying fraudulent activity

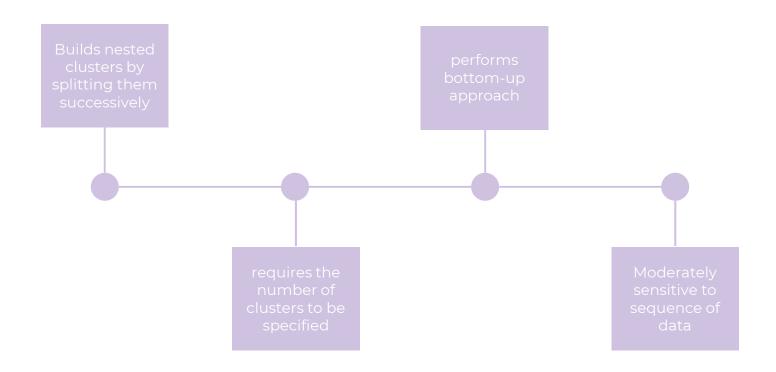
Based on the characteristics of the groups you are then able to classify them into those that are real and which are fraudulent.

03 Clustering Methods

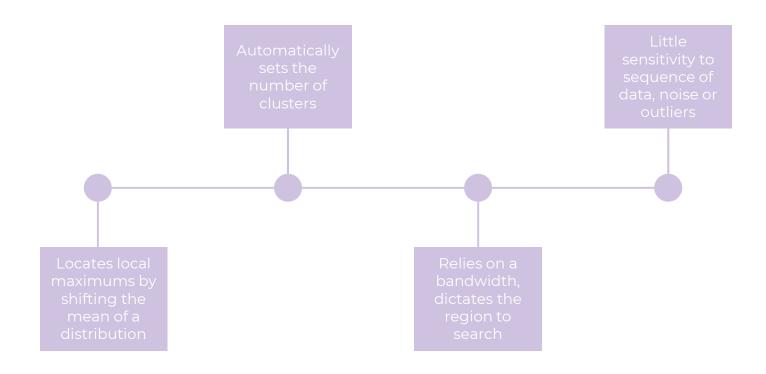
K-means – Based on Partiton



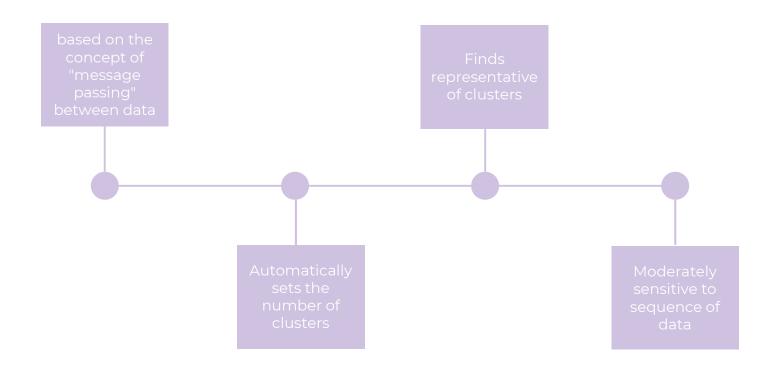
Agglomerative Clustering – based on Hierarchy



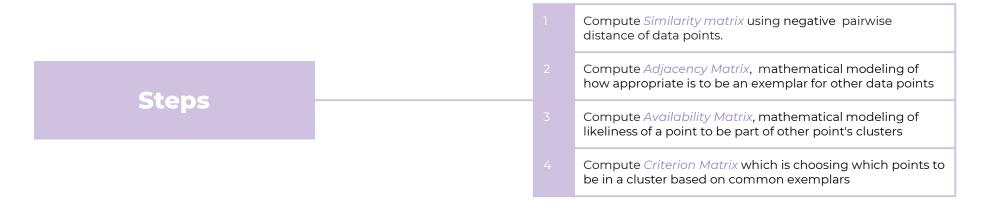
Mean shift - Based on Density



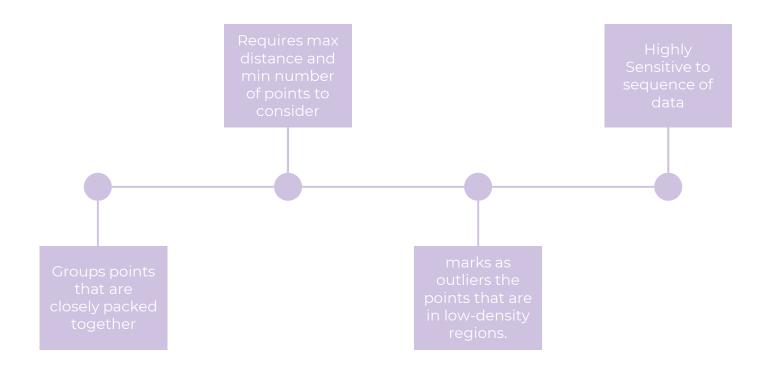
Affinity Propagation



Affinity Propagation



DBSCAN – Based on Density

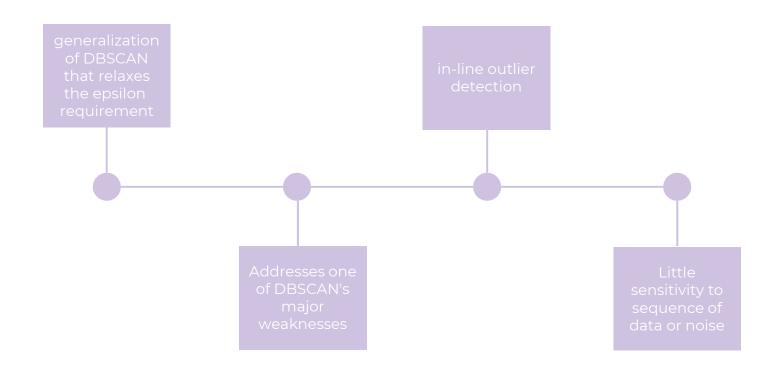


DBSCAN

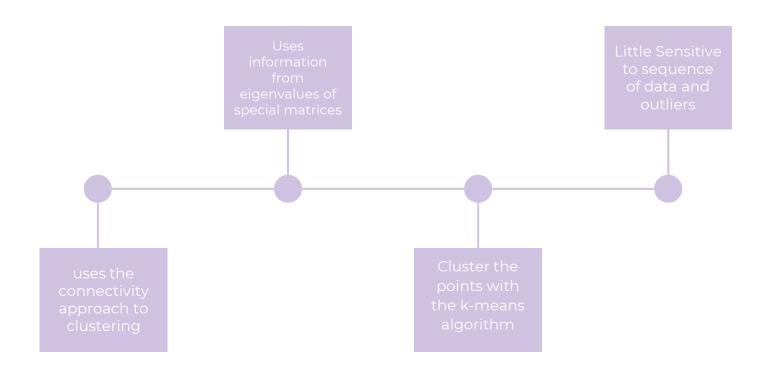
Checks the possibility of an unvisited point being a sample of estimated PDF or being noise, by measuring density of vicinity points.

for all unvisited vicinity numbers until reaching the border of estimated PDF, then checking other points to estimate other potential PDFs.

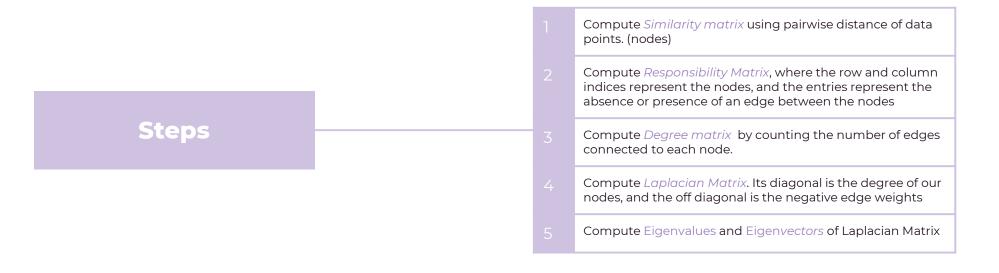
OPTICS – Based on Density



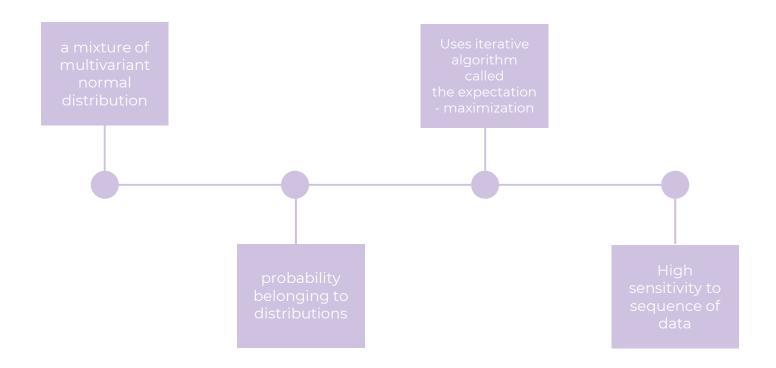
Spectral Clustering



Spectral Clustering



GMM - Based on Density



GMM

E – step: use current best knowledge of the center and shapes of each cluster to calculate the probability of belonging to cluster

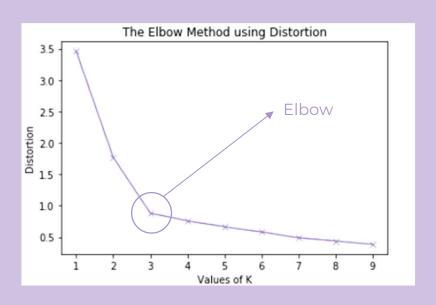
M – step: use current best knowledge of which class each point belongs to update and improve our estimates for the center and shape of each cluster

04 Tuning Cluster ng Methods

Elbow Method

In cluster analysis, the elbow method is a heuristic used in determining the number of clusters in a data set.

Silhouette coefficient and BIC score are better alternatives to the elbow method for visually discerning the optimal number of clusters.



Silhouette Method

The silhouette method computes silhouette coefficients of each point that measure how much a point is similar to its own cluster compared to other clusters.

The value of the silhouette ranges between [1, -1], where a high value indicates that the object is well matched to its own cluster and poorly matched to neighboring clusters

Information Criterion

a method for scoring and selecting a model

The Bayesian Information Criterion:

 $BIC: -2 \times LL + \log(N) \times k$

BIC penalizes the model more for its complexity, meaning that more complex models have worse scores

The Akaike Information Criterion:

$$AIC: -\frac{2}{N} \times LL + \log(N) \times \frac{k}{N}$$

AIC penalizes complex model less, meaning that it may put more emphasis on model performance O5
Clustering Metrics

Separation Index

Separation index is based on the ratio of the inter cluster distance to intra cluster distance, and then minimizing it across all the possible pairs.

$$d(S_i, S_j) = \min_{x,y} \{d(x, y | x \in S_i, y \in S_j)\}$$

$$d(S_l, S_l) = \max_{x,y} \{d(x, y | x, y \in S_l)\}$$

$$SI = \min_{j} \left\{ \min_{i(i \neq j)} \left\{ \frac{d(S_i, S_j)}{\max_{l} d(S_l, S_l)} \right\} \right\}$$

Fisher's Discrimination Index

like Separation, it maximizes a function by using covariance matrices of within and between cluster data points.

Between variability is calculated among centroids

Within variability is simply the sum of covariance matrices of the data points in a cluster.

$$S_{W} = \sum S_{i}$$

$$S_{B} = \sum Q_{i}(\widehat{\mu_{i}} - \hat{\mu})(\widehat{\mu_{i}} - \hat{\mu})^{T}$$

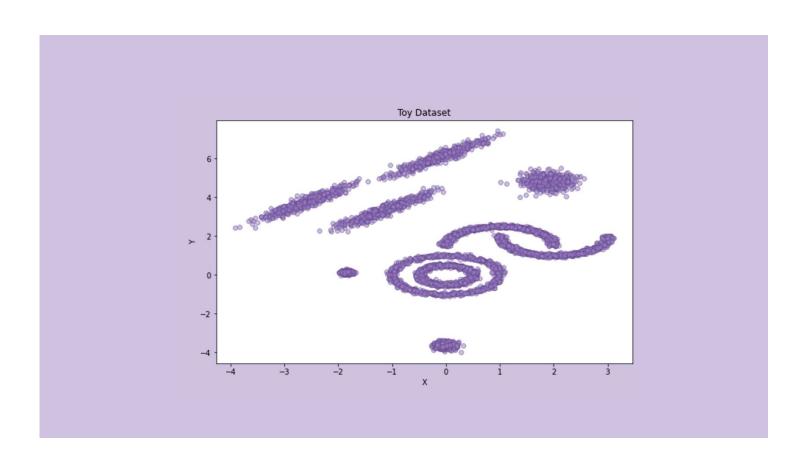
$$S_{i} = \sum (\widehat{\mu_{i}} - \hat{\mu})(\widehat{\mu_{i}} - \hat{\mu})^{T}$$

$$FDI = trace(S_{W}^{-1}S_{B})$$

06

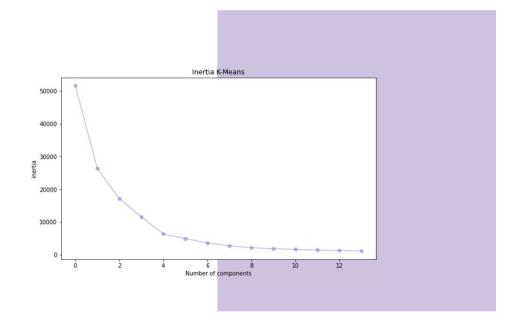
Study

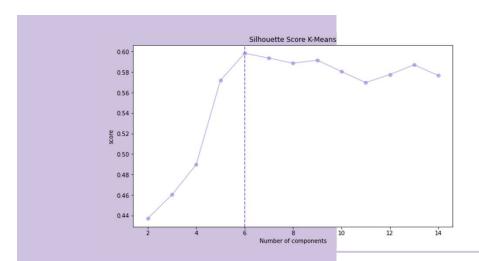
TOY dataset



K-Means

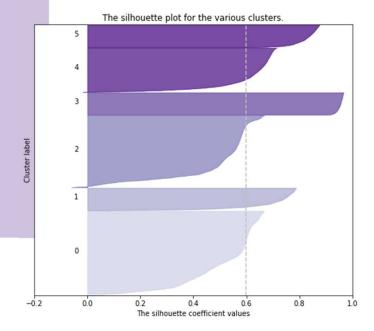
Elbow Method : Doesn't help with the clusters

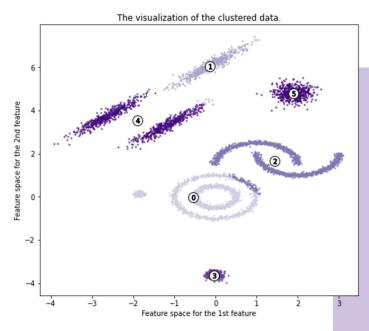




Sillhoutte Score : 6 clusters gives out the best score

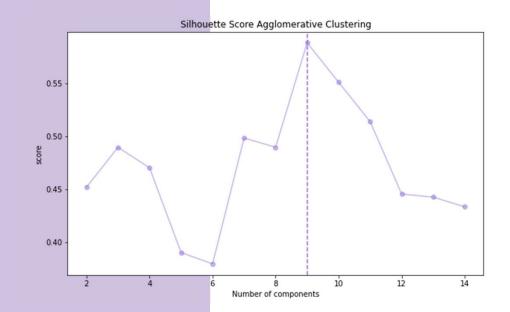
Silhouette analysis for K Means clustering on sample data with 6 clusters



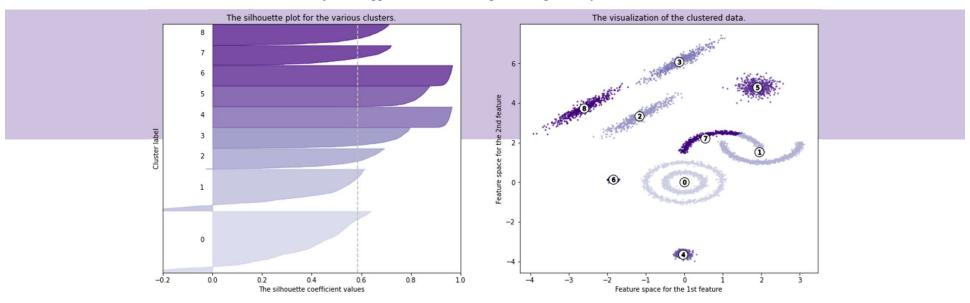


Agglomerative Clustering

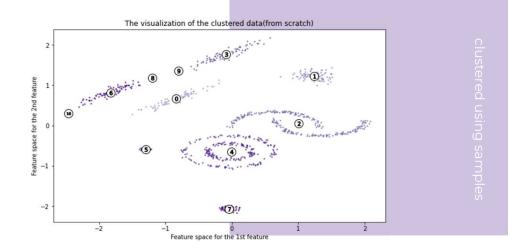
According to Sillhouette score, 9 clusters best represent our dataset.

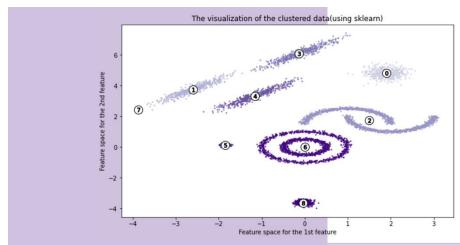


Silhouette analysis for Agglomerative Clustering clustering on sample data with 9 clusters

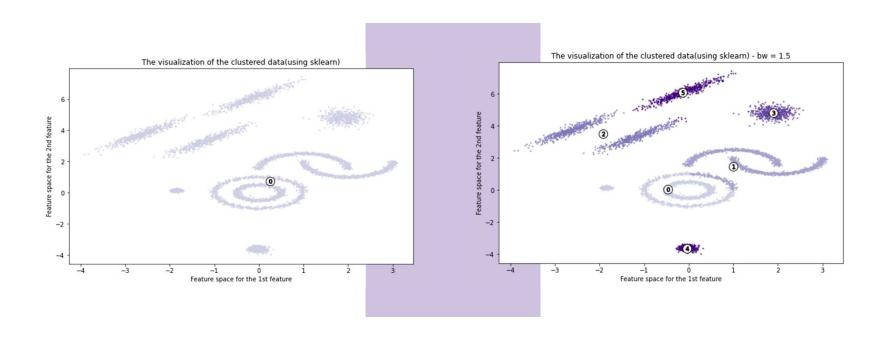


Agglomerative Clustering



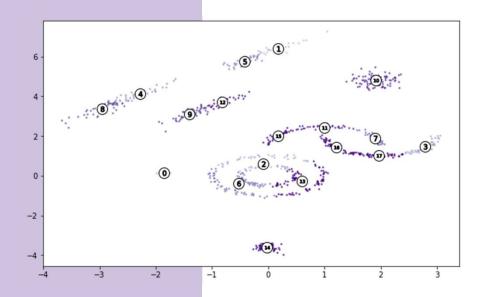


Mean Shift

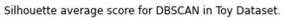


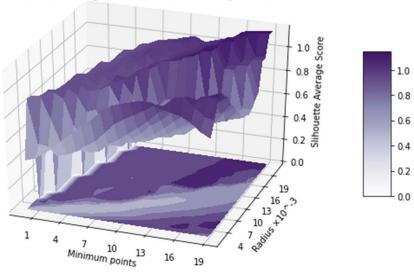
Affinity Clustering

clustered using samples and 0.8 damping factor!



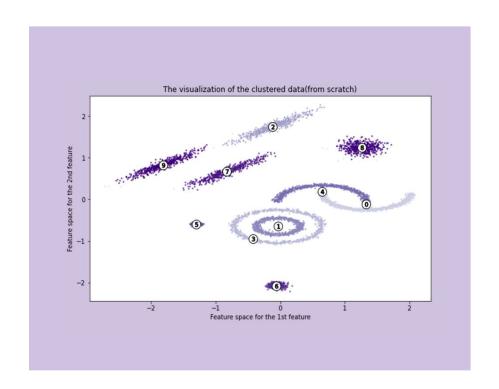
DBSCAN

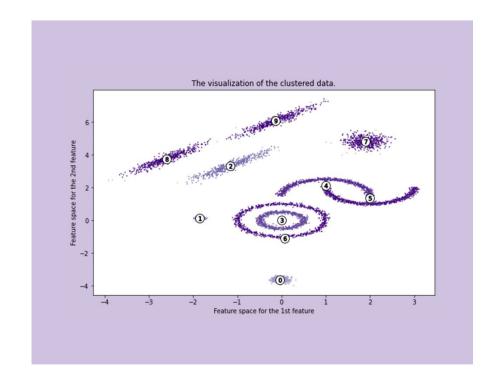




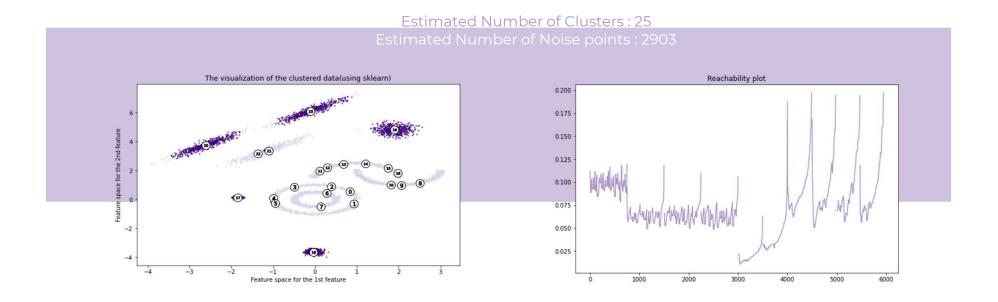
		Silhoue	ette Avera	age in 100	00 sample	s of Toy-	dataset		
2	0.168727	1.5579						7.79231	- 8
m -	0.0847504	0.578659	2.41361		7.08751		7.98646	8.09857	- 7
4	0.105631	0.230491	1 33167		6.80209			8.28545	- 6
r.	0.201421	0.168817	0.516635	2.62787	5.80232			7.40389	
9 .	0.337398	0.177046	0.323584	124008	3.77068		6.49768	7.09596	- 5
7	0.967308	0.34743	0.188807	0.72294	2.69039			7.08044	-4
ω -	110788	0.507871	0.254936	0.391491	1.58572			7.23568	- 3
σ.	1.04488	0.53736	0.313863	0.271644	0.972307	2.5579		6.40173	- 2
10	1 40401	1.9453	0.646137	0.239079	0.479133	1.79511		5.44881	- 1
п	1.40401	181603	0.646137	0.245068	0.347889	1.1271	2.85546		
	0.02	0.03	0.04	0.05	0.06	0.07	0.08	0.09	

DBSCAN



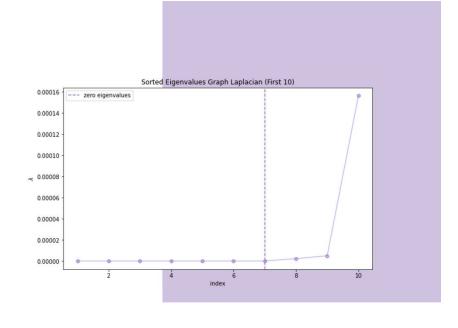


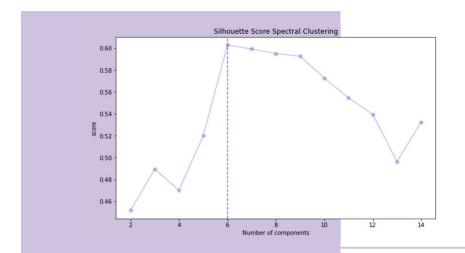
OPTICS



Spectral Clustering

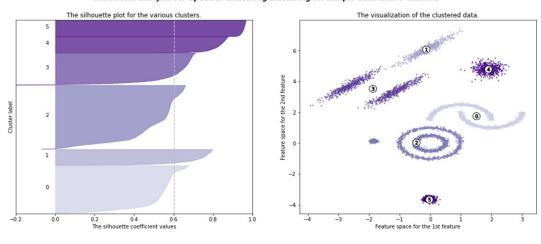
Optimal Number of Components turned out to be 7 according to spectral gap



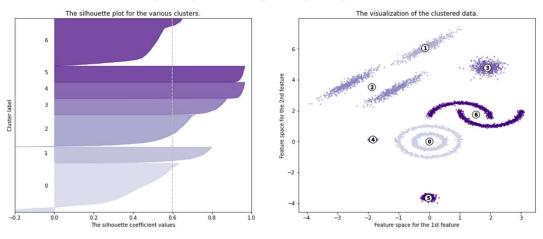


Optimal Number of Components turned out to be 7 according to Sillhoutte score

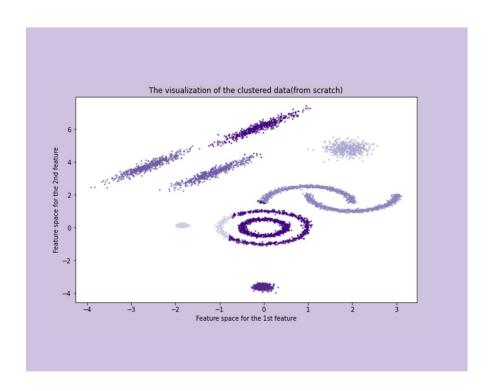
Silhouette analysis for Spectral Clustering clustering on sample data with 6 clusters

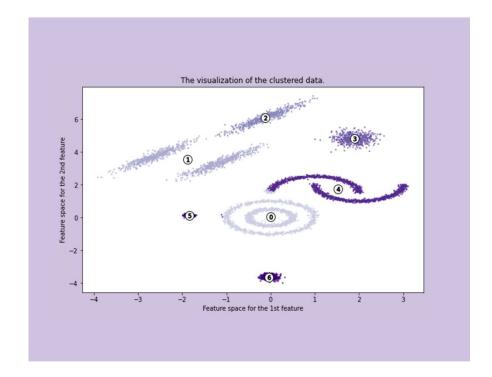


Silhouette analysis for Spectral Clustering clustering on sample data with 7 clusters

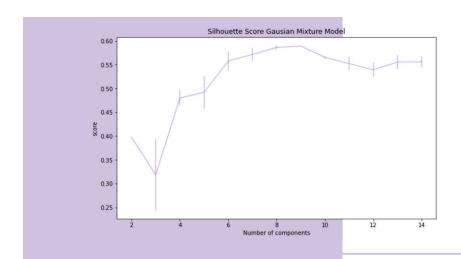


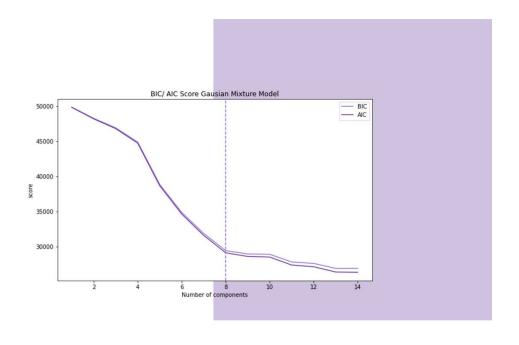
Spectral Clustering





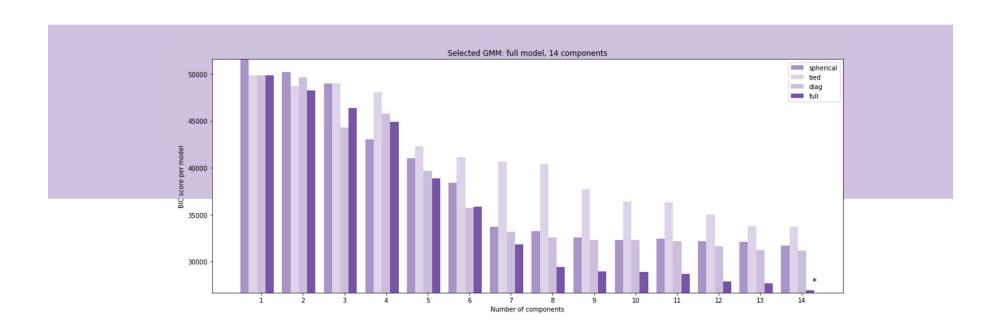
BIC/AIC score : can't really tell ...

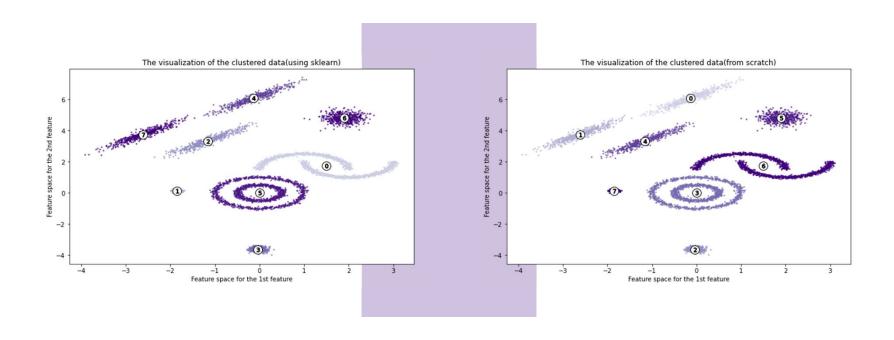


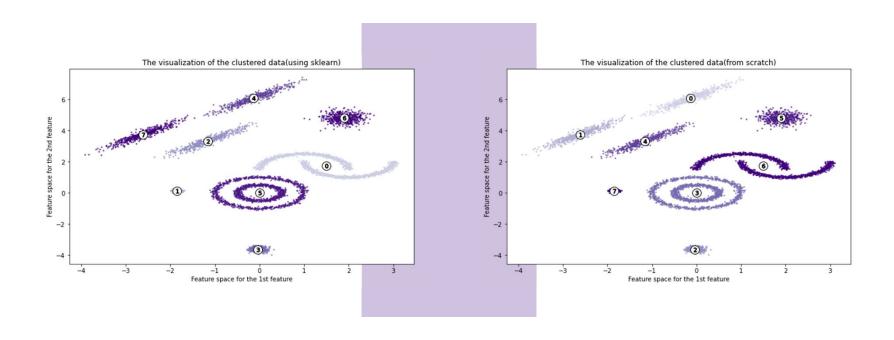


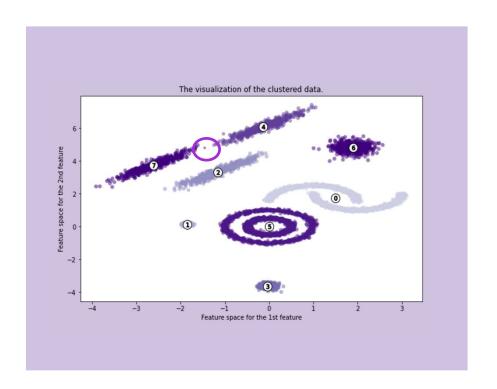
Sillhoutte Score : 8 clusters gives out the best score and lease amount of stochasticity

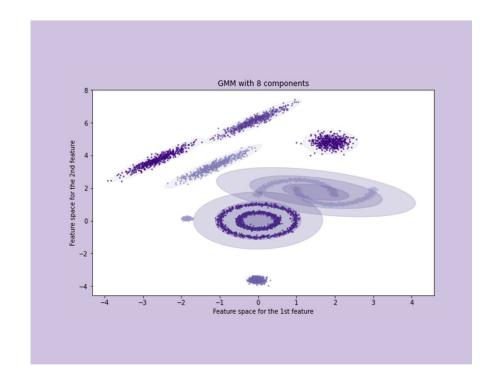
GMM - BIC score - Different covariance types



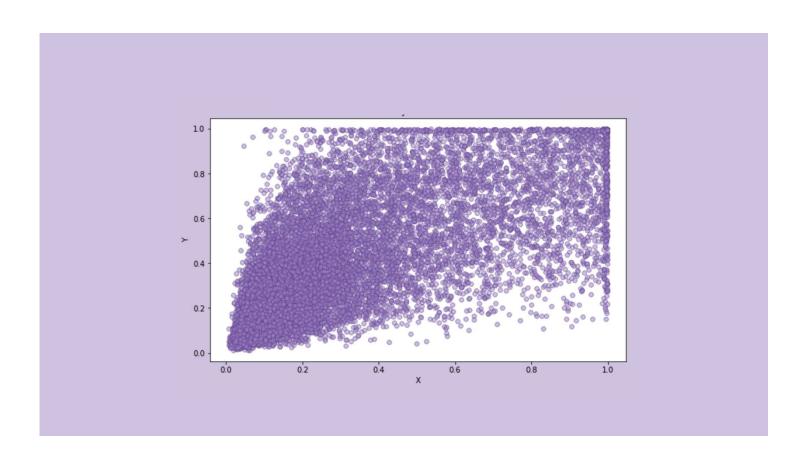






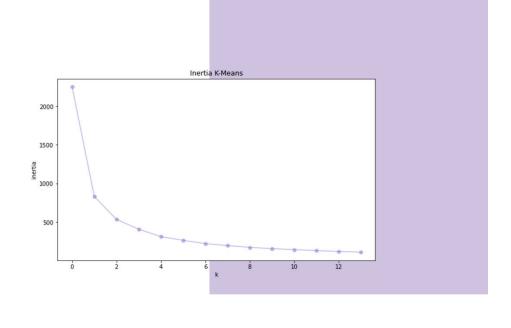


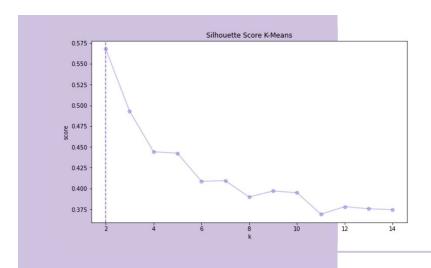
BBOX dataset



K-Means

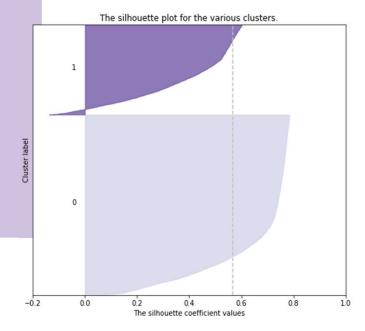
Elbow Method : Doesn't help with the clusters

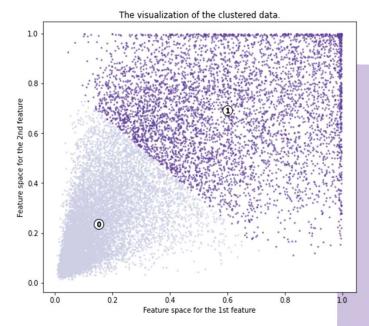




Sillhoutte Score: 2 clusters gives out the best score

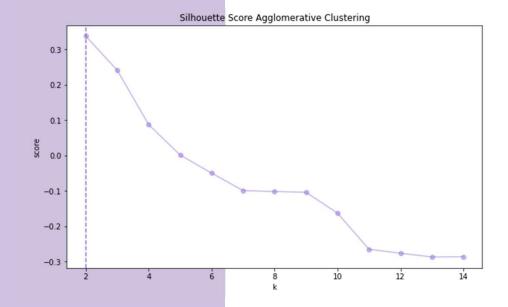
Silhouette analysis for K Means clustering on sample data with 2 clusters

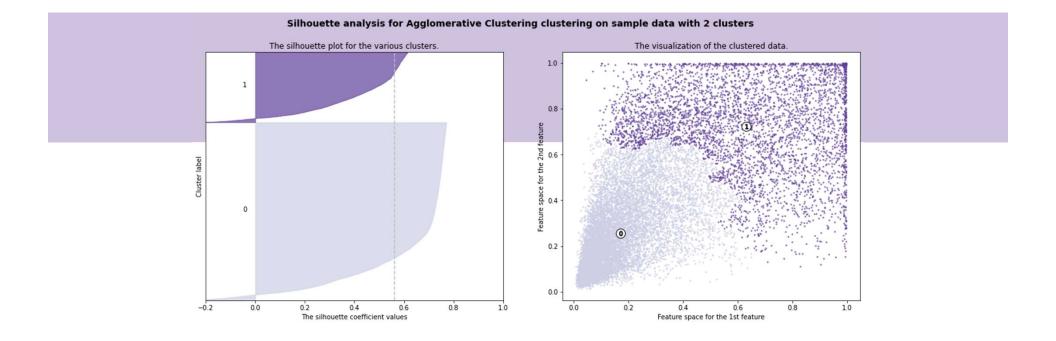


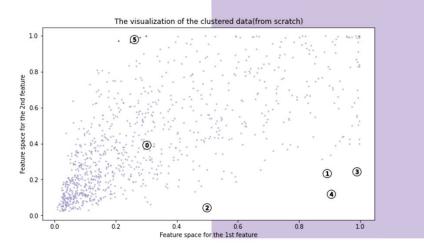


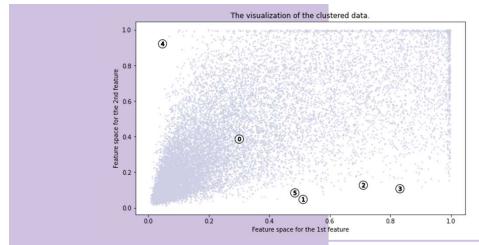
Agglomerative Clustering

According to Sillhouette score, 2 clusters best represent our dataset.

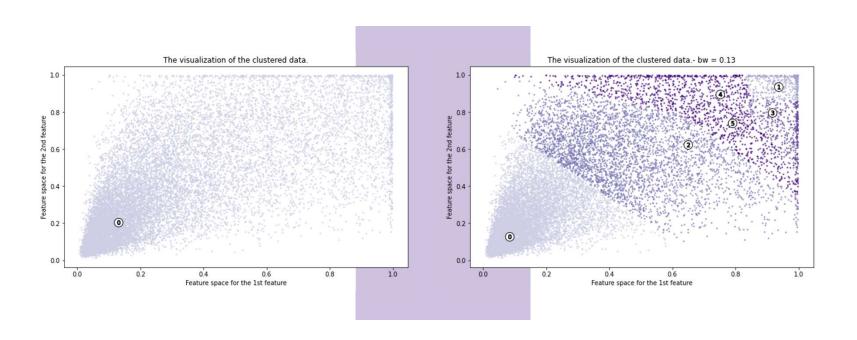






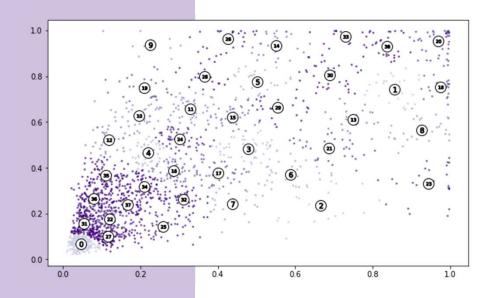


Mean Shift

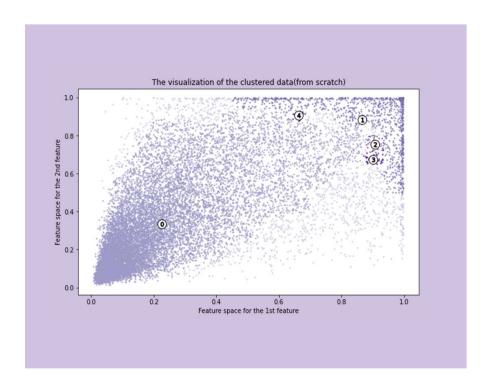


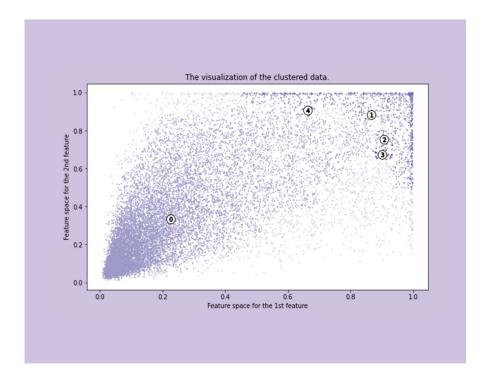
Affinity Clustering

clustered using samples and 0.8 damping factor!

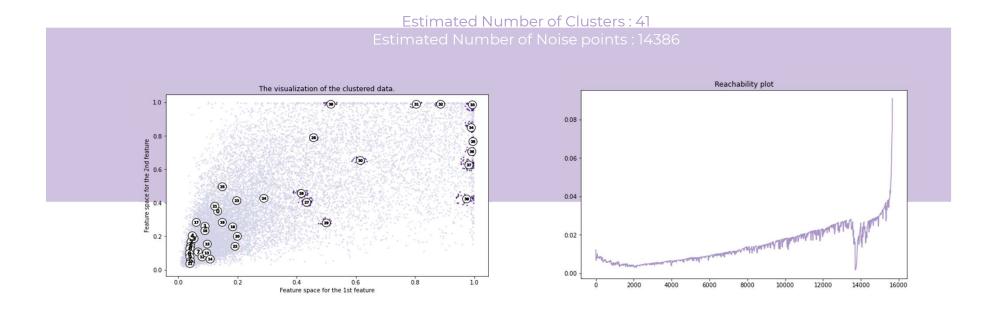


DBSCAN



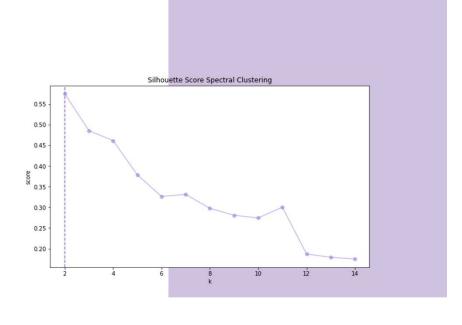


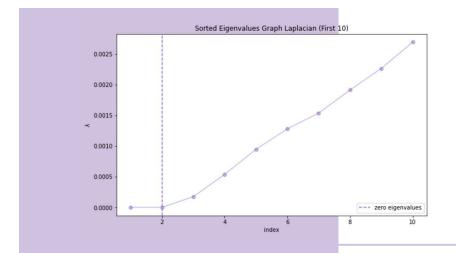
OPTICS



Spectral Clustering

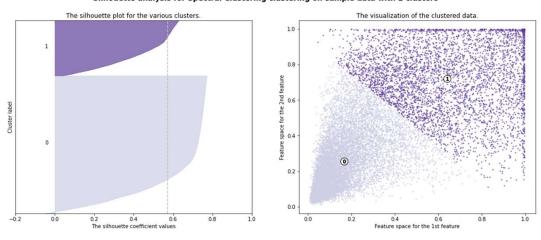
Optimal Number of Components turned out to be 2 according to spectral gap



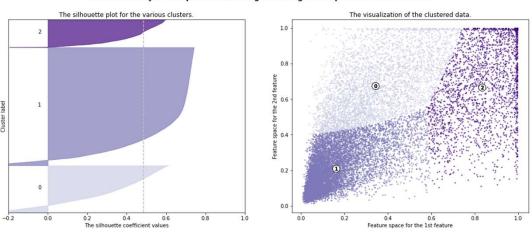


Optimal Number of Components turned out to be 2 according to Sillhoutte score

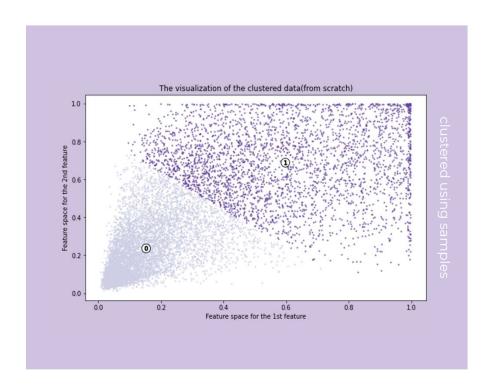
Silhouette analysis for Spectral Clustering clustering on sample data with 2 clusters

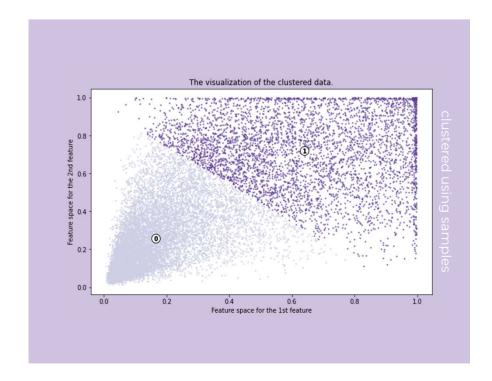


Silhouette analysis for Spectral Clustering clustering on sample data with 3 clusters

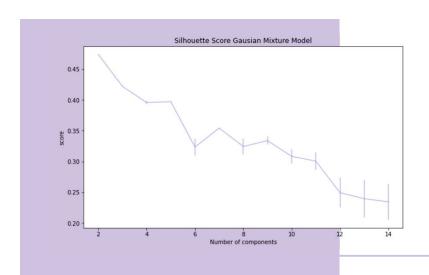


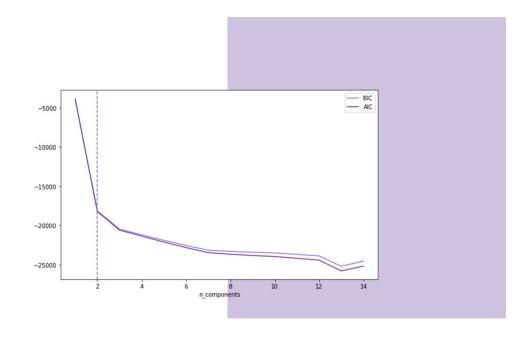
Spectral Clustering





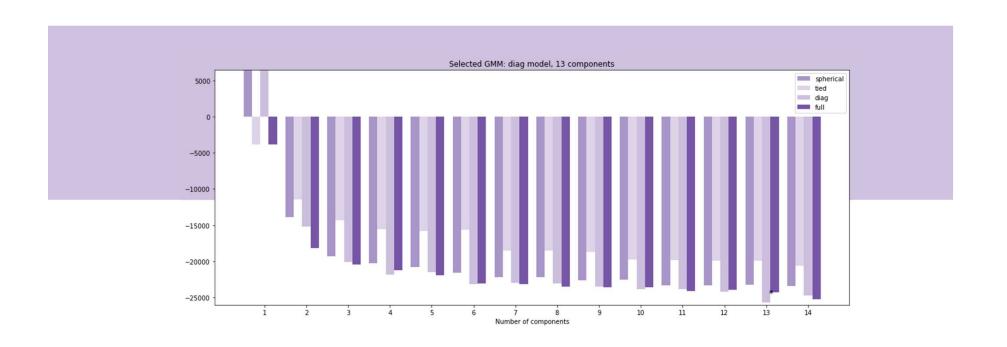
BIC/AIC score: 2 or 12!

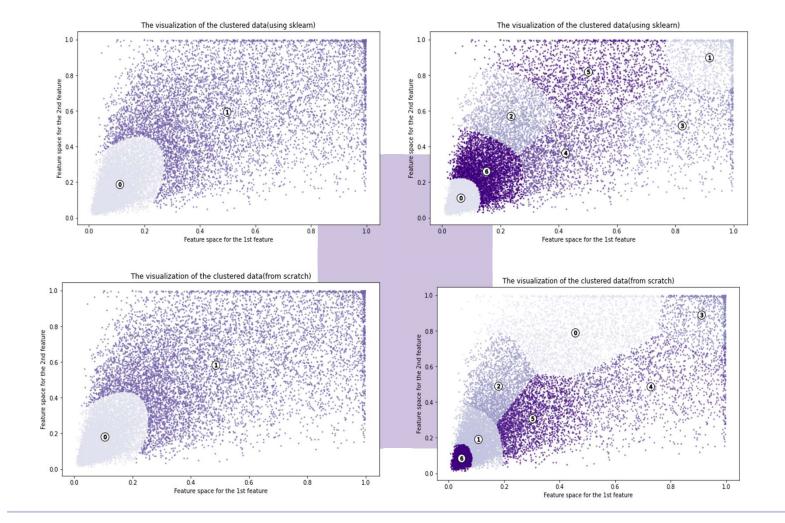




Sillhoutte Score: 2, 3, 5, 7 clusters all give out the very good scores and lease amount of stochasticity

GMM - BIC score - Different covariance types

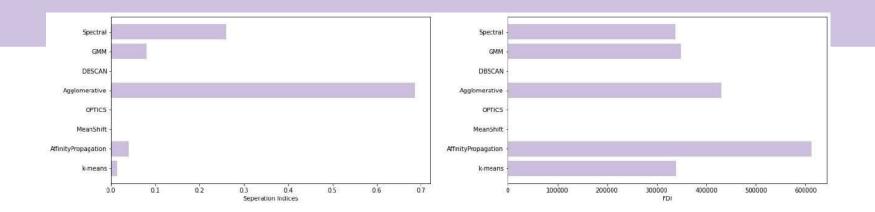




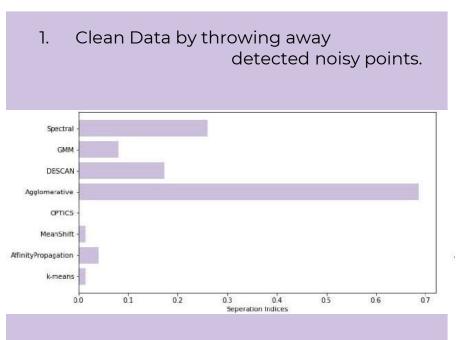
Comparison and Evaluation - Problems (sk-learn, Toy DataSet)

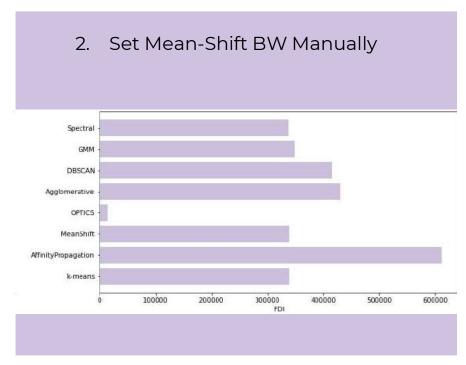
OPTICS, DBSCAN fail in evaluation due to Auto-Noise Detection

Data was not appropriate for mean-shift's estimated BW. Single Cluster



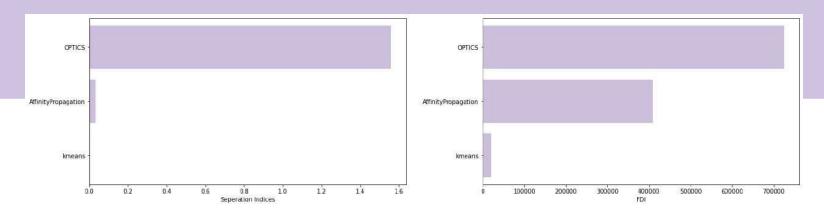
Comparison and Evaluation - Solutions (sk-learn, Toy DataSet)





Comparison and Evaluation (sk-learn, bboxes DataSet)

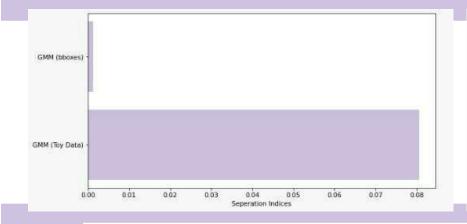
Due to high time and computation Complexity, 3 Algorithms that had fast convergence were evaluated on this dataset

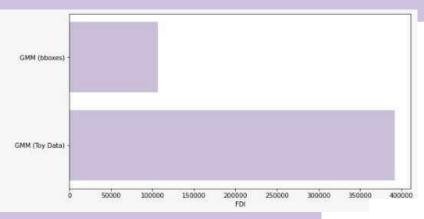


GMM - BIC score - Different covariance types

Scratch Algorithms we implemented usually performed better on Toy Dataset

One of the reasons can be the high dimensional structure and data complexity of 'bboxes'





THANKS!

CREDITS: This presentation template was created by Slidesgo, including icons by Flaticon, and infographics & images by Freepik.