

Assignment Project Report

Gaussian Mixture Models: Bag of Words Representation

Name: Pavan Tiwari

Course: AI and ML
(Batch 4)

- **Problem Statement**

Using a gaussian mixture model, perform a simple clustering on the given 2D Dataset. Try to find the optimal number of clusters using python (you may use any module to implement this). Now implement the same from scratch using python and a dummy dataset generated using scikit learn dataset generating functions such as make blob.

Software:

- Python 3 (Use anaconda as your python distributor as well)

– Tools:

- Pandas
- Numpy
- Matplotlib
- Seaborn
- OpenCv
- Sklearn

– Dataset Link:

Clustering_GMM

https://cdn.analyticsvidhya.com/wpcontent/uploads/2019/10/Clustering_gmm.csv

- **Method Used**

A Gaussian mixture model is a probabilistic model that assumes all the data points are generated from a mixture of a finite number of Gaussian distributions with unknown parameters. It attempts to find a mixture of multi-dimensional Gaussian probability distributions that best model any input dataset allowing the model to learn automatically, i.e. in an unsupervised manner. The bag-of-words model is a way of representing text data when modelling text with machine learning algorithms which can be combined with GMM to get a useful model representation.

- ### 1. Load all required libraries

1. Load all required libraries

[illegible]

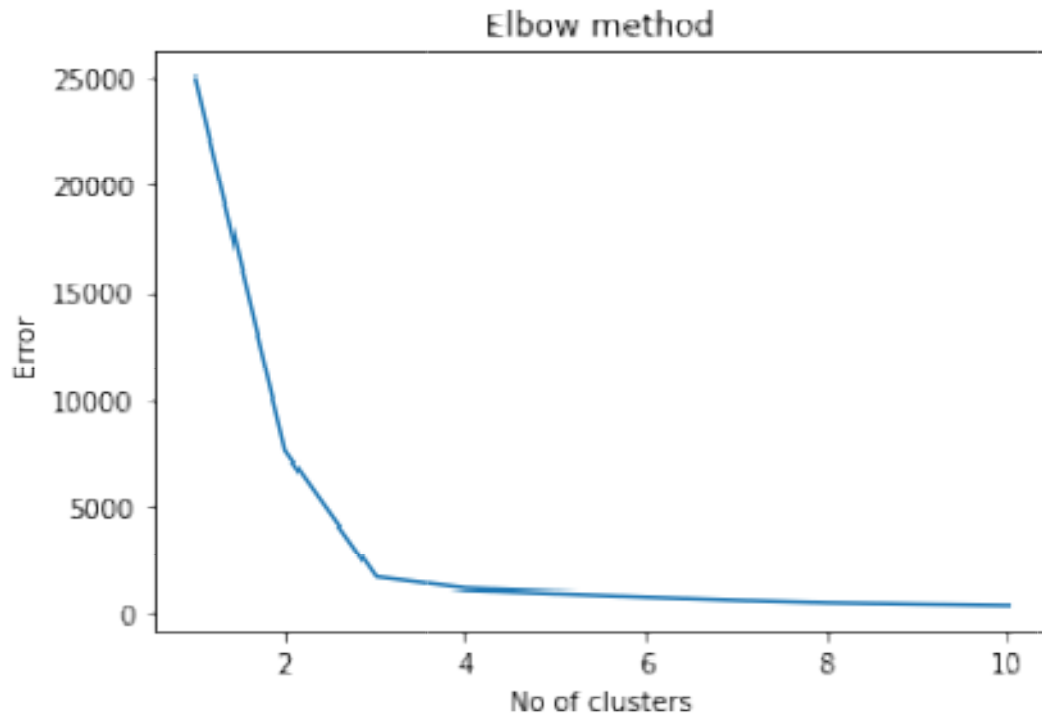
2. Deciding Number of Cluster

Elbow Method

```

1 Error = []
2 for i in range(1, 11):
3     kmeans = KMeans(n_clusters = i).fit(data)
4     kmeans.fit(data)
5     Error.append(kmeans.inertia_)
6 import matplotlib.pyplot as plt
7 plt.plot(range(1, 11), Error)
8 plt.title('Elbow method')
9 plt.xlabel('No of clusters')
10 plt.ylabel('Error')
11 plt.show()

```

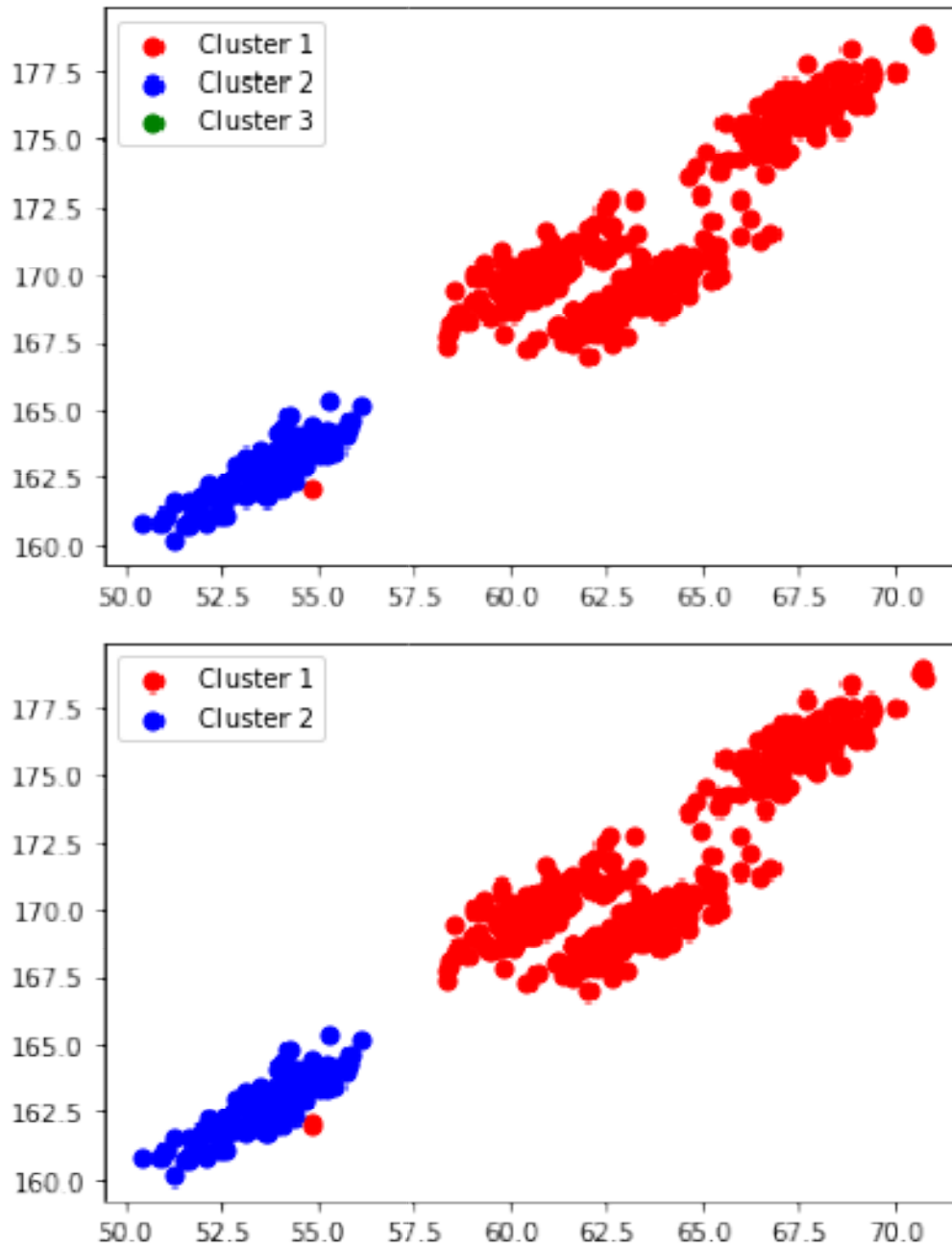


3. Model Building Using GMM On Clustering_GMM Dataset

Gaussian Mixture Models - Bag of Words Representation Last Checkpoint: Last Friday at 00:52 (

View Insert Cell Kernel Widgets Help

1



4. Importing Libraries and Making Dummy dataset

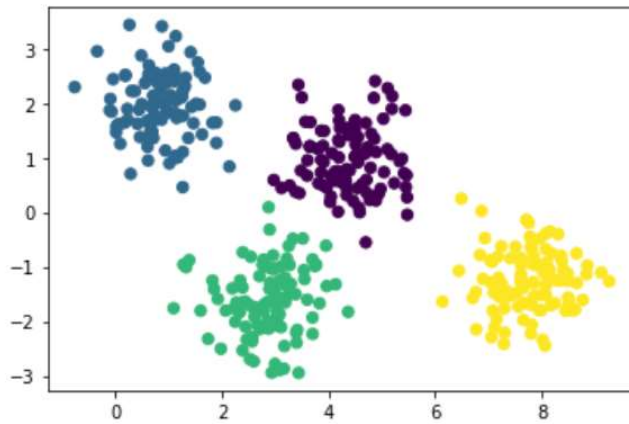
```

1 from sklearn.datasets import make_blobs

1 X, y_true = make_blobs(n_samples=400, centers=4,
2                        cluster_std=0.60, random_state=0)
3 X = X[:, ::-1] # flip axes for better plotting

1 plt.scatter(X[:, 0], X[:, 1], c=y_true, s=40, cmap='viridis')
2 plt.show()

```



5. Model Building

```

[18]: 1 from sklearn.mixture import GaussianMixture as GMM
2      gmm = GMM(n_components=4).fit(X)
3      labels = gmm.predict(X)
4      plt.scatter(X[:, 0], X[:, 1], c=labels, s=40, cmap='viridis')
5      plt.show()

```

