

## Clustering :

### Introduction

The Human Development Index (HDI) is an index that measures key dimensions of human development that mainly depends on human life expectancy, literacy rate and standard of living. In this assignment we are basically using the life expectancy and literacy rate of all the districts in Nepal for the year 2013.

Objective:

Using the dataset : Life Expectancy Income of Nepal by District and Literacy rates of different district of Nepal and if needed any other data set figure out the similar districts forming a cluster to be able to determine HDI based on the observations

### Models to be used:

#### K-means

Categorizes data items into clusters using the kMeans algorithm; an unsupervised learning algorithm. KMeans is used to categorize the items into k groups of similarity using euclidean distance as measurement.

Flow:

1. First initialize k points, called means, randomly.
2. Then categorize each item to its closest mean and update the mean's coordinates, which are the averages of the items categorized in that mean so far.
3. Then repeat the process for a given number of iterations and at the end,
4. Then at the end when we get the same centroid in iteration our cluster is formed.

Initialize k means with random values

For a given number of iterations:

Iterate through items:

Find the mean closest to the item

Assign item to mean

Update mean

To determine the optimal number of clusters we can use

**Elbow Method** : It uses Distortion ( average of square distances from cluster center) and Inertia (sum of squared distances of samples to their closest cluster center) for each value of K in given range

**Silhouette Algorithm** : Assumes data to be clustered into k clusters by a clustering technique

and then for each data point, we define the following:-

- I. The cluster assigned to the  $i$ th data point
- II. The number of data points in the cluster assigned to the  $i$ th data point
- lii. It gives a measure of how well assigned the  $i$ th data point is to its cluster
- iv. It is defined as the average dissimilarity to the closest cluster which is not its cluster
- V. Determine silhouette coefficient

The average silhouette for each value of  $k$  and for the value of  $k$  which has the maximum value of *silhouette coefficient* is considered the optimal number of clusters for the unsupervised learning algorithm.

## DBSCAN

Clustering analysis basically an Unsupervised learning method that divides the data points into a number of specific groups, such that the data points in the same groups have similar properties and data points in different groups have different properties

Similar to K-Means which is distance between points, DBSCAN is distance between nearest points, i.e **Density-based spatial clustering of applications with noise**. The **DBSCAN algorithm** is based on the concept of “clusters” and “noise”. The key idea is that for each point of a cluster, the neighborhood of a given radius has to contain at least a minimum number of points.

DBSCAN algorithm requires two parameters –

*eps* : It defines the neighborhood around a data point i.e. if the distance between two points is lower or equal to ‘eps’ then they are considered as neighbors. If the eps value is chosen too small then large part of the data will be considered as outliers. If it is chosen very large then the clusters will merge and majority of the data points will be in the same clusters. One way to find the eps value is based on the *k-distance graph*.

*MinPts*: Minimum number of neighbors (data points) within eps radius. Larger the dataset, the larger value of MinPts must be chosen. As a general rule, the minimum MinPts can be derived from the number of dimensions  $D$  in the dataset as,  $\text{MinPts} \geq D+1$ . The minimum value of MinPts must be chosen at least 3.

DBSCAN Clusters have 3 types of data points.

Core Point: A point is a core point if it has more than MinPts points within eps.

Border Point: A point which has fewer than MinPts within eps but it is in the neighborhood of a core point.

Noise or outlier: A point which is not a core point or border point.

## Gaussian Mixture Model

A Gaussian mixture model finds a mixture of multi-dimensional Gaussian probability distributions that best model the input dataset.

## Pipeline

### 1. Data exploration and analysis

Life-expectancy-income(set1) and literacy-rate(set2) data sets were loaded and was found that Life expectancy-income dataset consists of:

	District	Life expectancy(In Years)	Per Capita Income(In USD)	Unnamed: 3
58	Achham	67.14	536	NaN
36	Arghakhanchi	68.56	909	NaN
34	Baglung	68.83	868	NaN
32	Baitadi	68.88	573	NaN
66	Bajhang	65.22	487	NaN

*Sample Data first 5 row*

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 75 entries, 58 to 28
Data columns (total 4 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   District                             75 non-null     object
1   Life expectancy(In Years)            75 non-null     float64
2   Per Capita Income(In USD)            75 non-null     object
3   Unnamed: 3                           0 non-null      float64
dtypes: float64(2), object(2)
memory usage: 2.9+ KB
```

*Data Information*

And literacy-rate dataset consists of:

	District	Total	Female	Male	Year
59	Achham	55.7	42.9	70.7	2013
15	Arghakhanchi	72.6	65.8	81.8	2013
18	Baglung	71.9	65.3	80.6	2013
44	Baitadi	63.0	49.2	79.0	2013
61	Bajhang	55.6	40.1	73.0	2013

*Sample Data first 5 row*

```

<class 'pandas.core.frame.DataFrame'>
Int64Index: 75 entries, 59 to 31
Data columns (total 5 columns):
#   Column      Non-Null Count  Dtype
---  ---
0   District    75 non-null    object
1   Total       75 non-null    float64
2   Female      75 non-null    float64
3   Male        75 non-null    float64
4   Year        75 non-null    int64
dtypes: float64(3), int64(1), object(1)
memory usage: 3.5+ KB

```

### Data information

Both dataset has data for 75 districts

On generating descriptive statistics for the sets:

Life expectancy(In Years) Unnamed: 3		
count	75.000000	0.0
mean	68.405333	NaN
std	2.251472	NaN
min	61.200000	NaN
25%	67.285000	NaN
50%	68.550000	NaN
75%	70.190000	NaN
max	72.900000	NaN

Stats for set1 data

	Total	Female	Male
count	75.0000	75.000000	75.000000
mean	65.1120	56.053333	75.016000
std	9.5331	11.199968	8.464617
min	41.7000	32.000000	50.900000
25%	57.0500	47.600000	69.300000
50%	66.2000	57.100000	76.200000
75%	71.9000	64.350000	80.900000
max	86.3000	79.800000	92.200000

Stats for set2 data

## 2. Feature Selection and preprocessing:

There are some features in both the datasets which is not required for further processing.

- From Set1 data 'Unnamed: 3' is not needed as all the Data in that column has NaN Value so it can be dropped. Similarly from Set2 all the data are of year 2013 so that too can be dropped

Now to merge the data frames using pandas merge feature

```

1 exp_income_literacy = pd.merge(expectency_income, literacy_rate, on='District', how='inner')
2 exp_income_literacy = exp_income_literacy.sort_values(by=['District'])
3 exp_income_literacy.shape
4

```

(150, 6)

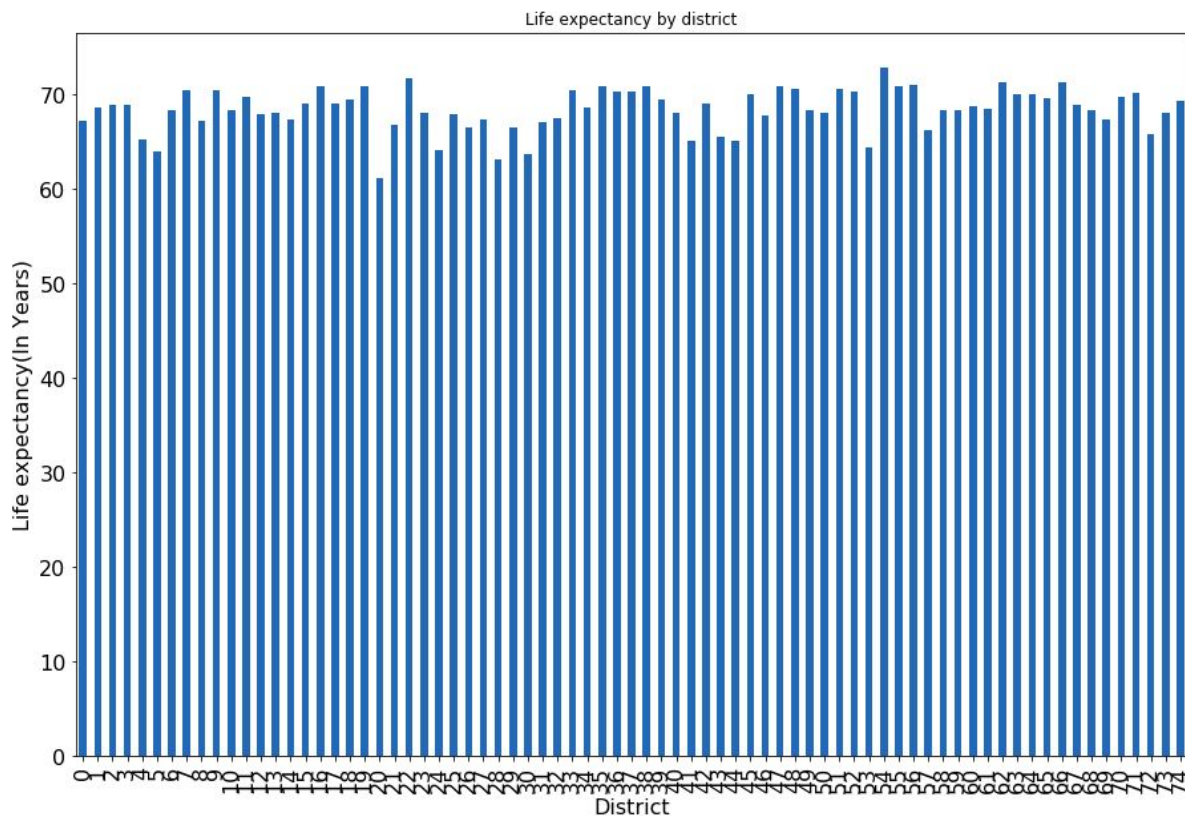
There seems to be something off with the data as we have 75 districts and the merged data is showing 150 districts so further data analysis and preprocessing is required

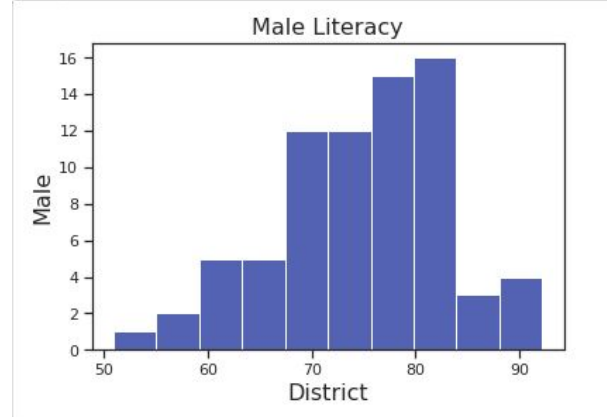
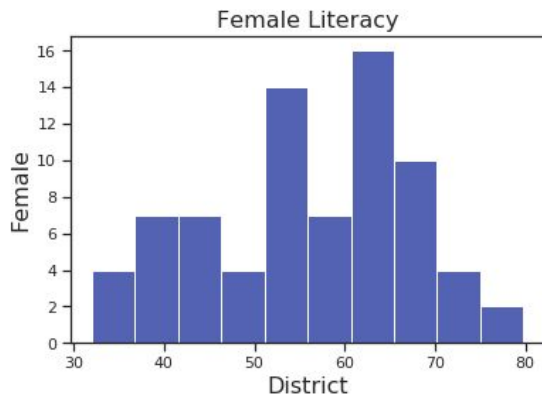
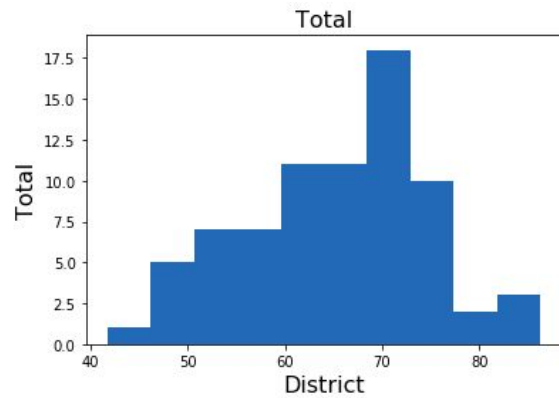
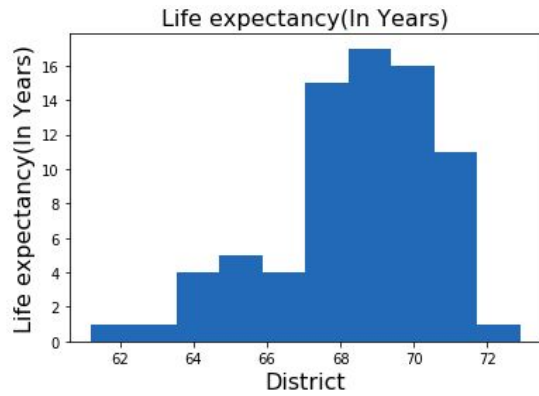
It is found that the District Name in both the sets do not match with each other as one of the set has space at the end of the name so after preprocessing the data the sets were merged and the final merged table obtained is as below:

	District	Life expectancy(In Years)	Per Capita Income(In USD)	Total	Female	Male
0	Achham	67.14	536	55.7	42.9	70.7
1	Arghakhanchi	68.56	909	72.6	65.8	81.8
2	Baglung	68.83	868	71.9	65.3	80.6
3	Baitadi	68.88	573	63.0	49.2	79.0
4	Bajhang	65.22	487	55.6	40.1	73.0

## Data Visualization

Simple Bar plot showing the Life Expectancy in years by District shows District 20 has lowest life span of around 60

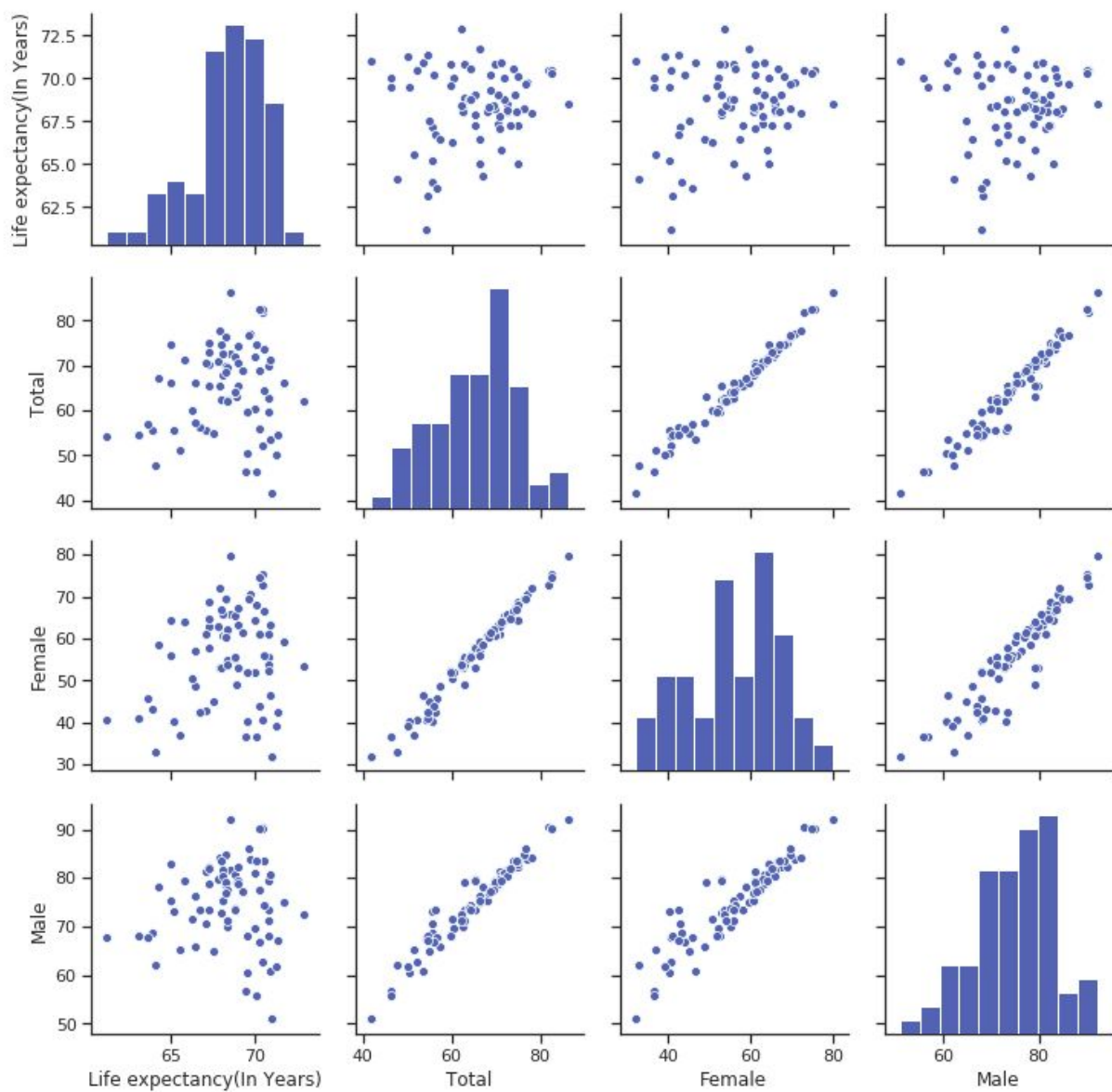




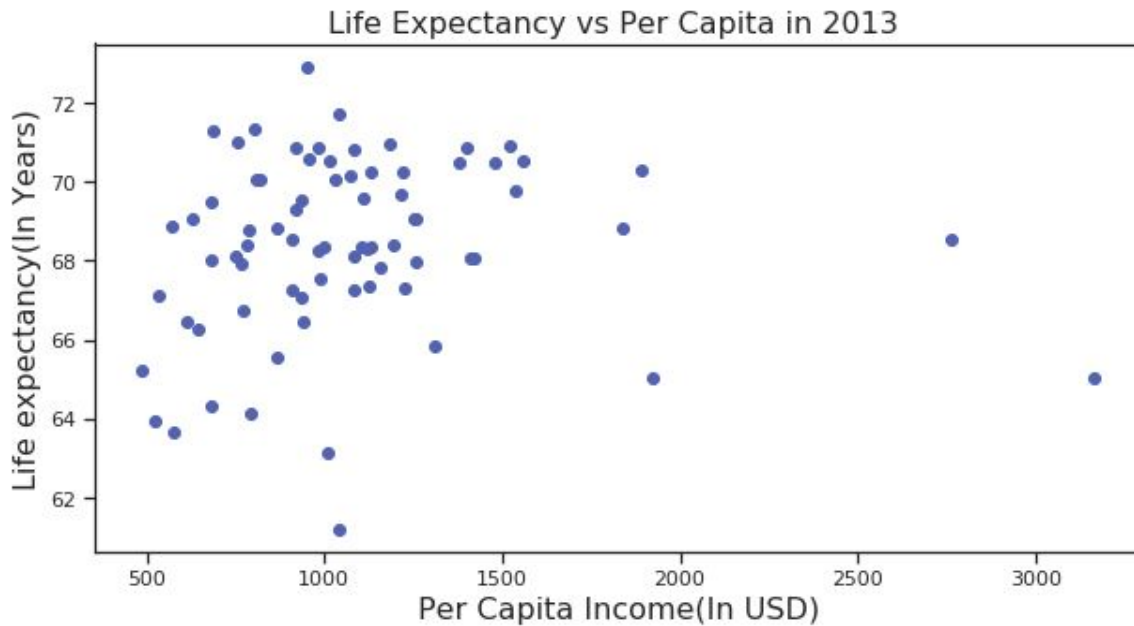
Histogram plot for each attributes based on district

	Life expectancy(In Years)	Per Capita Income(In USD)	Total	Female	Male
Life expectancy(In Years)	1.000000	0.070954	0.134884	0.209683	0.046707
Per Capita Income(In USD)	0.070954	1.000000	0.506040	0.498649	0.419661
Total	0.134884	0.506040	1.000000	0.986817	0.971105
Female	0.209683	0.498649	0.986817	1.000000	0.924431
Male	0.046707	0.419661	0.971105	0.924431	1.000000

Determining the correlation and visualizing it using pairplot







Scatter plot between Life Expectancy and Capita Income

### Standardization

Now Need to normalize the data to scale the data to 0 and 1 range so we gonna use standard scaler library to scale down the data

**from sklearn.preprocessing import StandardScaler**

```
array([[ -0.56578717, -1.23572093, -0.99394547, -1.1823167 , -0.51332084],
       [ 0.06915839, -0.38816428,  0.79076324,  0.87610087,  0.80685092],
       [ 0.18988747, -0.48132734,  0.7168404 ,  0.83115725,  0.66412965],
       [ 0.21224471, -1.15164695, -0.22303579, -0.61602715,  0.47383462],
       [-1.42430511, -1.34706215, -1.00450587, -1.43400094, -0.23977174]])
```

Above sample scaled data shows that the data has been normalized

### 3. K-Means | Approach taken and findings

Library used:

```
from sklearn.cluster import KMeans
from sklearn.decomposition import PCA
from sklearn.metrics import silhouette_score
```

To determine the optimal cluster:



```

: 1 Sum_of_squared_distances=[]
  2 Silhouette=[]
  3 for k in range(2,8):
  4     kmeans = KMeans(n_clusters=k, random_state=1)
  5     kmeans.fit(LitInc_scaled_data)
  6     Sum_of_squared_distances.append(kmeans.inertia_)
  7     b = silhouette_score(LitInc_scaled_data, kmeans.labels_)
  8     Silhouette.append(b)
  9     print(k,kmeans.inertia_)
10     print('Silhouette:', silhouette_score(LitInc_scaled_data, kmeans.labels_))

```

```

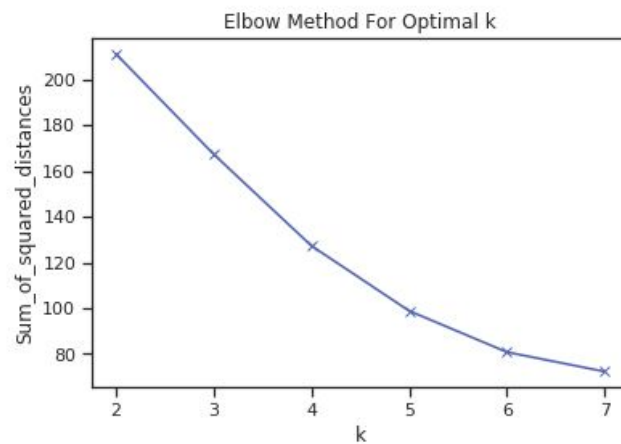
2 211.20721018478366
Silhouette: 0.38093100314853784
3 167.34711398154263
Silhouette: 0.3571966109546932
4 127.29913069221143
Silhouette: 0.3772347815392019
5 98.85166518122696
Silhouette: 0.3342231166632858
6 80.81867444869053
Silhouette: 0.3142932206212687
7 72.31968942752886
Silhouette: 0.28249121561822

```

```

1 plt.plot(range(2,8), Sum_of_squared_distances, 'bx-')
2 plt.xlabel('k')
3 plt.ylabel('Sum_of_squared_distances')
4 plt.title('Elbow Method For Optimal k')
5 plt.show()

```

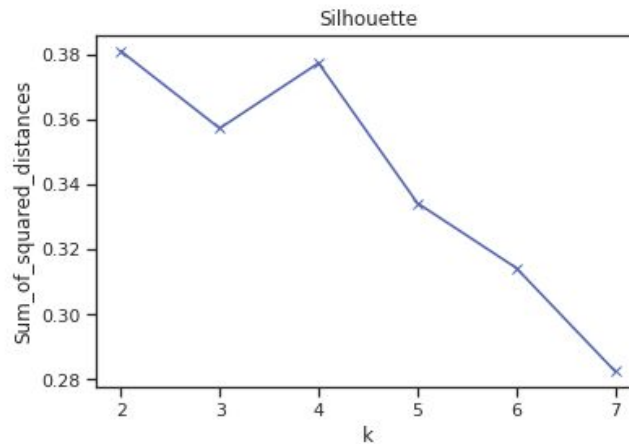


Elbow Method

```

1 plt.plot(range(2,8), Silhouette, 'bx-')
2 plt.xlabel('k')
3 plt.ylabel('Sum_of_squared_distances')
4 plt.title('Silhouette')
5 plt.show()

```



Silhouette Method

Analysing both the method we found that 6 is the optimal number of cluster

### Model 1: Simple KMeans

```

1 kmeans=KMeans(n_clusters=6, random_state=1)
2 kmeans.fit(LitInc_scaled_data)
3 kmeans_df_clustered = exp_income_literacy.copy()
4 kmeans_df_clustered['KMeans Label'] = kmeans.labels_

```

```

1 pd.set_option('max_rows', 100)
2 kmeans_df_clustered.sort_values(by='KMeans Label')

```

```

1 y_kmeans = kmeans.predict(LitInc_scaled_data)
2 y_kmeans

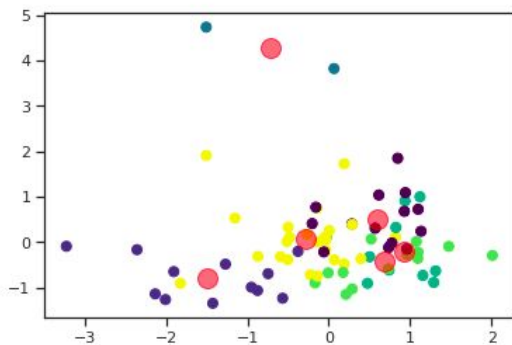
```

array([1, 5, 5, 4, 1, 1, 4, 3, 5, 0, 5, 0, 5, 4, 5, 4, 4, 0, 3, 4, 1, 1,  
4, 5, 1, 0, 1, 5, 1, 5, 1, 5, 1, 0, 2, 0, 5, 0, 0, 3, 5, 2, 5, 1,  
5, 0, 5, 4, 4, 0, 5, 0, 3, 5, 4, 3, 3, 1, 4, 5, 4, 5, 3, 3, 4, 4,  
3, 5, 5, 5, 0, 0, 5, 0, 5], dtype=int32)

```

1 plt.scatter(LitInc_scaled_data[:, 0], LitInc_scaled_data[:, 1], c=y_kmeans, s=50, cmap='viridis')
2
3 centers = kmeans.cluster_centers_
4 plt.scatter(centers[:, 0], centers[:, 1], c='red', s=200, alpha=0.5);
5

```



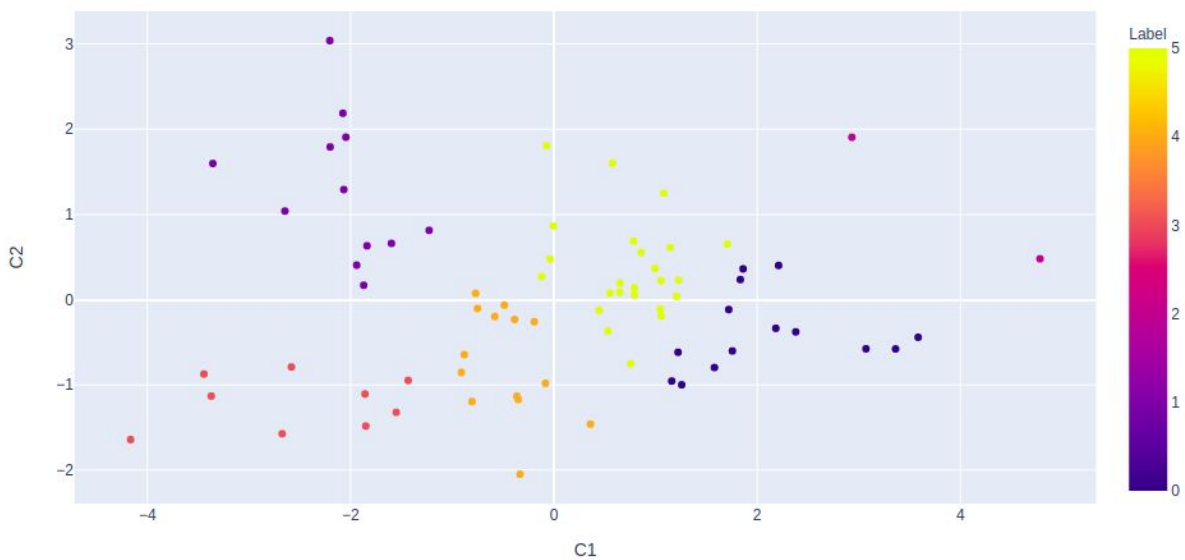
Simple scatter plot for the data with Centroid

Using PCA to plot the Kmeans plot

```

1 pca = PCA(n_components=2, random_state = 1)
2 pca = pca.fit(LitInc_scaled_data)
3 scaled_transformed = pca.transform(LitInc_scaled_data)
4 plot_data = pd.DataFrame(scaled_transformed, columns=['C1', 'C2'])
5 plot_data['Label'] = kmeans.labels_
6 plot_data['District'] = exp_income_literacy['District'].tolist()
7 fig = px.scatter(plot_data, x="C1", y="C2", color='Label', hover_data=['District'])
8 fig.show()

```



Model 2: KMeans after outlier detection

```

1 kmeans_outlier = KMeans(n_clusters=1)
2 kmeans_outlier.fit(LitInc_scaled_data)

```

```

KMeans(algorithm='auto', copy_x=True, init='k-means++', max_iter=300,
       n_clusters=1, n_init=10, n_jobs=None, precompute_distances='auto',
       random_state=None, tol=0.0001, verbose=0)

```

```

1 # identify the 5 closest points
2 distances = kmeans_outlier.transform(LitInc_scaled_data)

```

```

1 # argsort returns an array of indexes which will sort the array
2 # in ascending order. Reverse it with[::-1]
3 sorted_idx = np.argsort(distances.ravel()[::-1])[::-1][:5]

```

```

1 print(LitInc_scaled_data[sorted_idx][:, 0])
2 print(LitInc_scaled_data[sorted_idx][:, 1])

```

```

[ 0.06468694 -1.50479117  1.15572016 -3.22182705 -1.90275    ]
[ 3.82689626  4.74034873 -0.7335493  -0.09049693 -0.64947531]

```

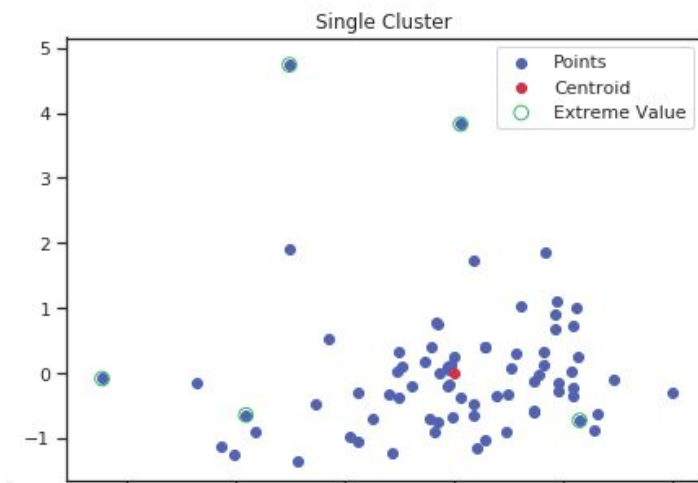
Found that 5 point can be taken as an outlier

```

1 f, ax = plt.subplots(figsize=(7,5))
2 ax.set_title('Single Cluster')
3 ax.scatter(LitInc_scaled_data[:, 0], LitInc_scaled_data[:, 1], label='Points')
4 ax.scatter(kmeans_outlier.cluster_centers[:, 0],
5           kmeans_outlier.cluster_centers[:, 1],
6           label='Centroid', color='r')
7 ax.scatter(LitInc_scaled_data[sorted_idx][:, 0],
8           LitInc_scaled_data[sorted_idx][:, 1],
9           label='Extreme Value', edgecolors='g',
10          facecolors='none', s=75)
11 ax.legend(loc='best')

```

<matplotlib.legend.Legend at 0x7f890870ed68>



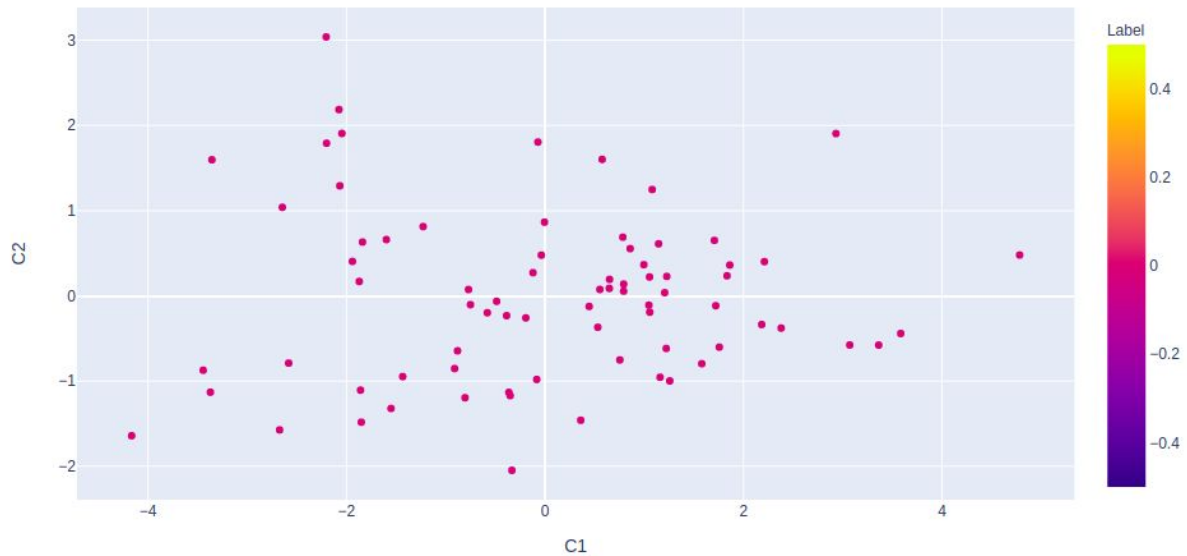
From the above plot got that there are 5 outliers denoted by Extreme Values and Point in red is centroid



```

1 plot_data = pd.DataFrame(scaled_transformed, columns=['C1', 'C2'])
2 plot_data['Label'] = kmeans_outlier.labels_
3 plot_data['District'] = exp_income_literacy['District'].tolist()
4 fig = px.scatter(plot_data, x="C1", y="C2", color='Label', hover_data=['District'])
5 fig.show()

```



### Model 3: Kmeans using PCA fit

```

1 n_samples, n_features = LitInc_scaled_data.shape
2 n_digits = len(np.unique(y_kmeans))
3 labels = y_kmeans
4 n_noise_k_means = kmeans.labels_
5 sample_size = 75

```

```

1 pca = PCA(n_components=2, random_state = 1)
2 pca = pca.fit(LitInc_scaled_data)
3 scaled_transformed = pca.transform(LitInc_scaled_data)
4 n_digits = 4
5 # labels = outlier_y_kmeans
6

```

```

1 # pca = PCA(n_components= 4).fit(LitInc_scaled_data)
2 print('-----PCA-ANALYSIS-----')
3 print("Noise:", n_noise_k_means)
4 evaluate_kmeans(KMeans(init='k-means++', n_clusters=n_digits, n_init=10),
5                 name="PCA-based", data=scaled_transformed)

```

On Evaluation of the model:

Model 1:

noise: 0  
Model Evaluation: k-means++  
\*\*\*\*\*  
Estimated number of clusters: 6  
Homogeneity: 0.944  
Completeness: 0.938  
V-measure: 0.941  
Adjusted Rand Index: 0.936  
Adjusted Mutual Information: 0.934  
silhouette\_score: 0.312037534584409

Model 2:

noise: 5  
Model Evaluation: k-means with Outlier detection  
\*\*\*\*\*  
Estimated number of clusters: 6  
Homogeneity: 0.897  
Completeness: 0.872  
V-measure: 0.884  
Adjusted Rand Index: 0.886  
Adjusted Mutual Information: 0.868  
silhouette\_score: 0.3228536072142073

Model 3:

-----PCA-ANALYSIS-----  
Model Evaluation: PCA-based  
\*\*\*\*\*  
Estimated number of clusters: 4  
Homogeneity: 0.611  
Completeness: 0.791  
V-measure: 0.689  
Adjusted Rand Index: 0.539  
Adjusted Mutual Information: 0.663  
silhouette\_score: 0.37752484349028187

#### 4. DBSCAN | Approach taken and findings

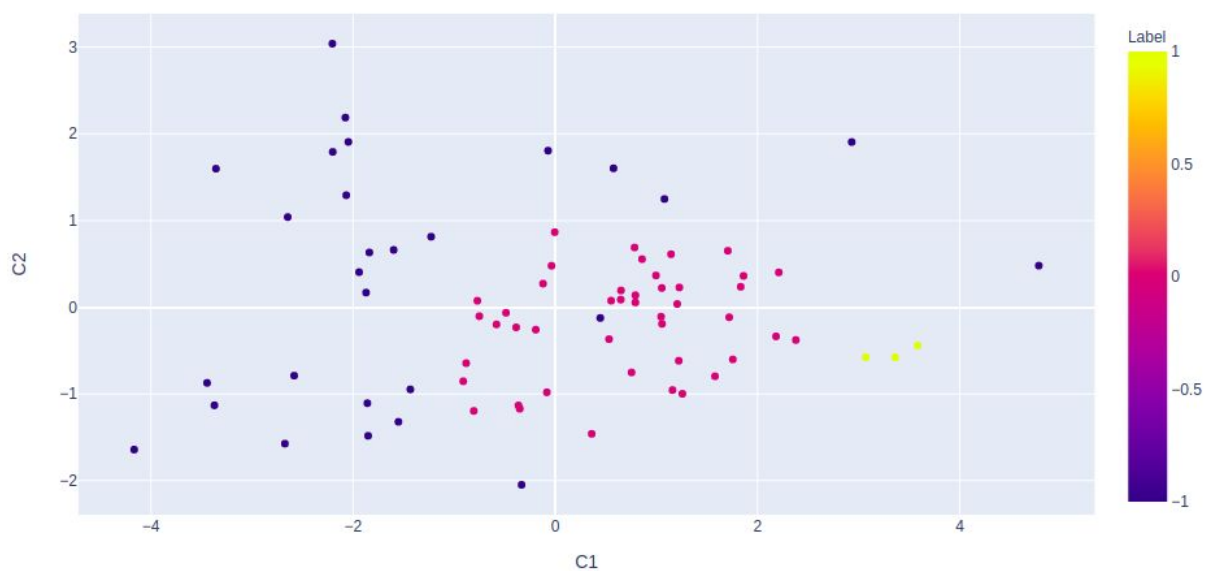
Using eps=0.8 and min\_sample =3

```

1 min_samples = 3
2 eps = 0.8
3
4 # dbscan = DBSCAN(eps=2., min_samples=8)
5 dbscan = DBSCAN(eps=0.8, min_samples=3)
6
7 y_dbscan = dbscan.fit_predict(LitInc_scaled_data)
8 exp_income_literacy['DBSCAN Label'] = dbscan.labels_

1 str(exp_income_literacy['DBSCAN Label'].tolist())
2
'[-1, 0, 0, 0, -1, -1, 0, -1, 0, 1, 0, 0, 0, 0, 0, 0, 0, -1, 0, -1, -1, 0, 0, -1, 0, -1, 0, -1, 0, -1,
1, -1, 0, 0, 1, 0, -1, 0, -1, 0, -1, -1, 0, 0, 0, 0, 0, 0, 0, -1, -1, -1, -1, -1, -1, 0, 0, 0, 0, -1, -1, 0, 0, -
1, -1, 0, 0, 0, 0, -1, 0, 0]'
```

No of clusters = 2 obtained



#### Model 4: DBSCAN

Estimated number of clusters: 2  
 Estimated number of noise points: 28  
 Homogeneity: 1.000  
 Completeness: 1.000  
 V-measure: 1.000  
 Adjusted Rand Index: 1.000  
 Adjusted Mutual Information: 1.000  
 Silhouette Coefficient: 0.095



## 5. GMM | Approach taken and findings

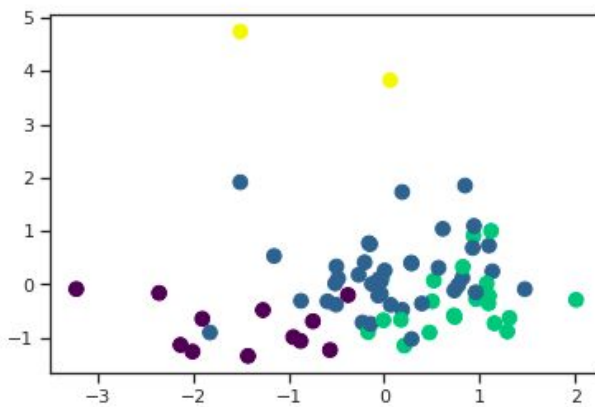
No of cluster = 4

```
1 from sklearn import mixture
```

```
1 mix_data = mixture.GaussianMixture(n_components=4, covariance_type='diag')  
2 mix_data.fit(X)
```

```
GaussianMixture(covariance_type='diag', init_params='kmeans', max_iter=100,  
                means_init=None, n_components=4, n_init=1, precisions_init=None,  
                random_state=None, reg_covar=1e-06, tol=0.001, verbose=0,  
                verbose_interval=10, warm_start=False, weights_init=None)
```

```
1 labels = mix_data.predict(X)  
2 plt.scatter(X[:, 0], X[:, 1], c=labels, s=75, cmap='viridis');
```



```

1 probs = mix_data.predict_proba(X)
2
3 print(probs[:5].round(3))

```

```

[[0.989 0.    0.011 0.   ]
 [0.    1.    0.    0.   ]
 [0.    1.    0.    0.   ]
 [0.    0.097 0.903 0.   ]
 [1.    0.    0.    0.   ]]

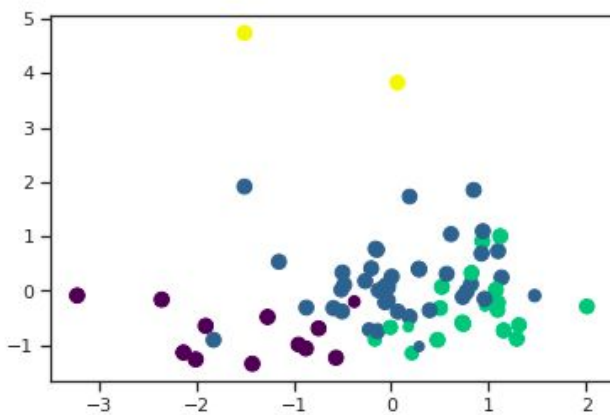
```

Visualizing the probabilistic cluster

```

1 size = 75 * probs.max(1) ** 2 # square emphasizes differences
2 plt.scatter(X[:, 0], X[:, 1], c=labels, cmap='viridis', s=size);

```



## Comparison

	Model 1	Model 2	Model 3	Model 4
Estimated number of clusters:	6	6	4	2
Estimated number of noise points:	0	5		28
Homogeneity Score	0.944	0.897	0.611	1.000
Completeness	0.938	0.872	0.791	1.000
V-measure:	0.941	0.884	0.689	1.000
Adjusted Rand Index:	0.936	0.886	0.539	1.000
Adjusted Mutual Information:	0.934	0.868	0.663	1.000
Silhouette Coefficient:	0.312	0.32285	0.3775	0.095

**URL**

<https://github.com/Kristiee/HDI-clustering>

**PCA Task:**

<https://github.com/Kristiee/HDI-clustering/blob/master/PCA.ipynb>

**Submitted by:**

**Shristi Maskay**