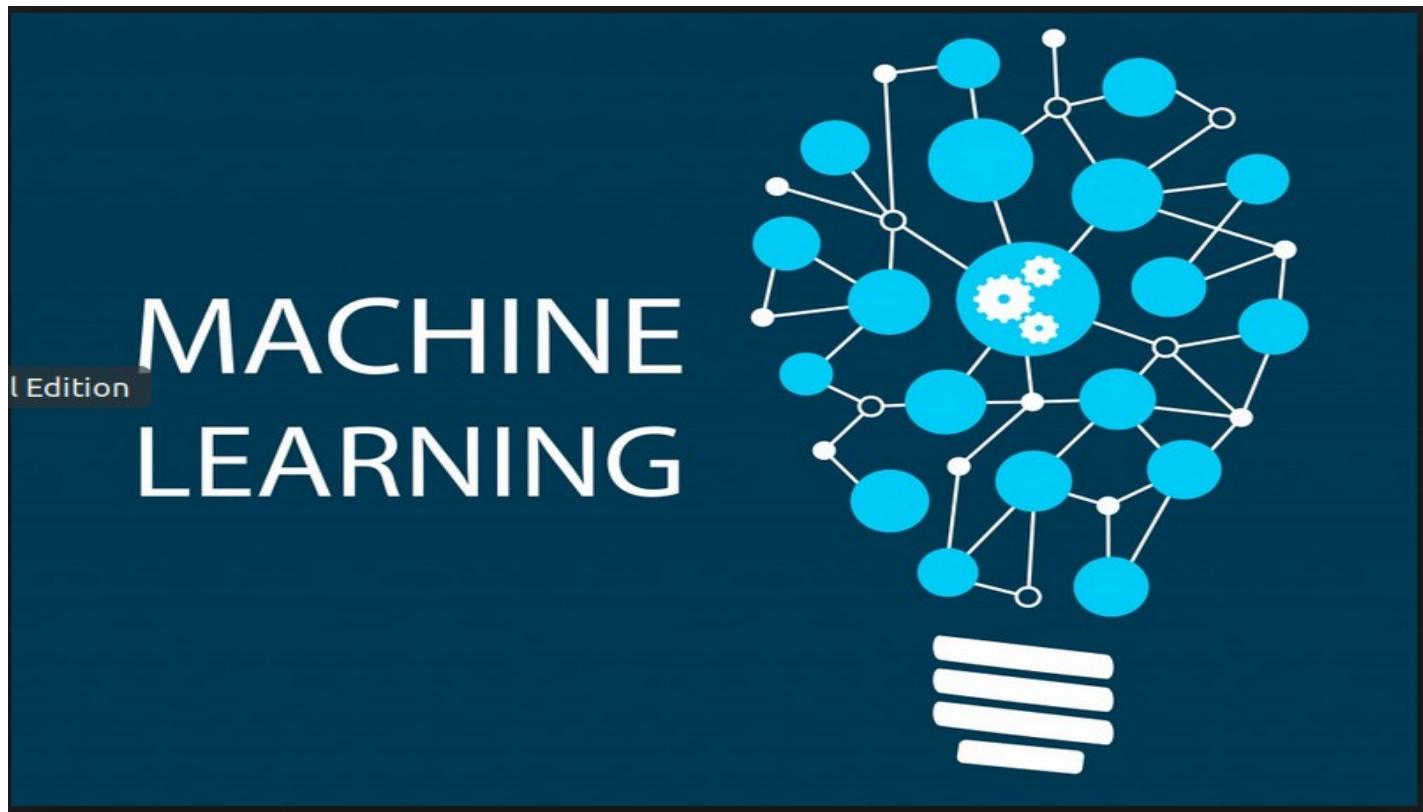


Machine Learning

Cluster



Name: Merna Zakaria

ID:4106

Cluster:

K-means clustering is one of the simplest and popular unsupervised machine learning algorithms.

DATASET USED:

CIFAR:

The CIFAR-10 dataset consists of 60000 32x32 colour images in 10 classes, with 6000 images per class. There are 50000 training images and 10000 test images.

The dataset is divided into five training batches and one test batch, each with 10000 images. The test batch contains exactly 1000 randomly-selected images from each class. The training batches contain the remaining images in random order, but some training batches may contain more images from one class than another. Between them, the training batches contain exactly 5000 images from each class.

Each of the batch files contains a dictionary with the following elements:

- **data** -- a 10000x3072 numpy array of uint8s. Each row of the array stores a 32x32 colour image. The first 1024 entries contain the red channel values, the next 1024 the green, and the final 1024 the blue. The image is stored in row-major order, so that the first 32 entries of the array are the red channel values of the first row of the image.
- **labels** -- a list of 10000 numbers in the range 0-9. The number at index i indicates the label of the i th image in the array data.

Screen Shot OF Cluster Code:

The screenshot shows the PyCharm IDE interface with the file 'Clustering.py' open. The code implements the K-Means clustering algorithm. It starts by initializing cluster centers and then iterates through the following steps: 1) Assigning each data point to its nearest cluster center. 2) Calculating the total distance from each point to its assigned cluster center. 3) Updating the cluster centers based on the mean of all points assigned to them. 4) Repeating the process until the total distortion (sum of squared distances) no longer decreases.

```
52 BelongToCluster=[0]*(len(features))*k
53 for it in range(0,iterations):
54     dist = 0
55     #get distance between labels and clusters
56     for labels in range(0,len(features)):
57         mnDis=1e16
58         for clusters in range(0,k):
59             ed=0
60             for dim in range(0,len(features[labels])):
61                 ed+=(features[labels][dim] - CentersOfCluster[clusters][dim])**2
62             if ed < mnDis:
63                 mnDis=ed
64             BelongToCluster[labels]=clusters
65
66     #get new centers
67     for centers in range(0,k):
68         total=0
69         points=[0]*len(features[0])
70         for labels in range(0,len(features)):
71             if BelongToCluster[labels] == centers:
72                 total+=1
73                 for dim in range(0,len(features[0])):
74                     points[dim]+=features[labels][dim]
75         if total>0:
76             for dim in range(0,len(features[0])):
77                 points[dim]/=total
78             CentersOfCluster[centers]=points
79
80     #get distortion after each iteration
81     for cluster in range(0,k):
82         for it in range(0,iterations) > for labels in range(0,len(featu...
```

This screenshot shows the same PyCharm project with an additional section of code added at the bottom to visualize the clustering results. It uses the Matplotlib library to plot the distortion value (sum of squared distances) against the number of iterations. The plot shows a decreasing trend, indicating that the algorithm is converging.

```
#get distortion after each iteration
for cluster in range(0,k):
    total = 0
    points = [0] * len(features[0])
    for labels in range(0, len(features)):
        if BelongToCluster[labels] == cluster:
            total += 1
            for dim in range(0, len(features[0])):
                points[dim] += features[labels][dim]
    if total > 0:
        for dim in range(0, len(features[0])):
            points[dim] /= total
        for label in range(0, len(features)):
            if BelongToCluster[label] == cluster:
                for dim in range(0, len(features[labels])):
                    dist += (CentersOfCluster[cluster][dim] - features[label][dim]) ** 2
    distortion.append(dist)

plt.plot(distortion)
plt.xlabel("iterations")
plt.ylabel("summation of diff. between each object and nearest cluster to it")
plt.show()
```

1) K-means clustering is one of the simplest and popular unsupervised machine learning algorithms.

Accuracy:

As we use clustering algorithm we can't calculate accuracy according to the labels of the image because cluster is unsupervised learning.

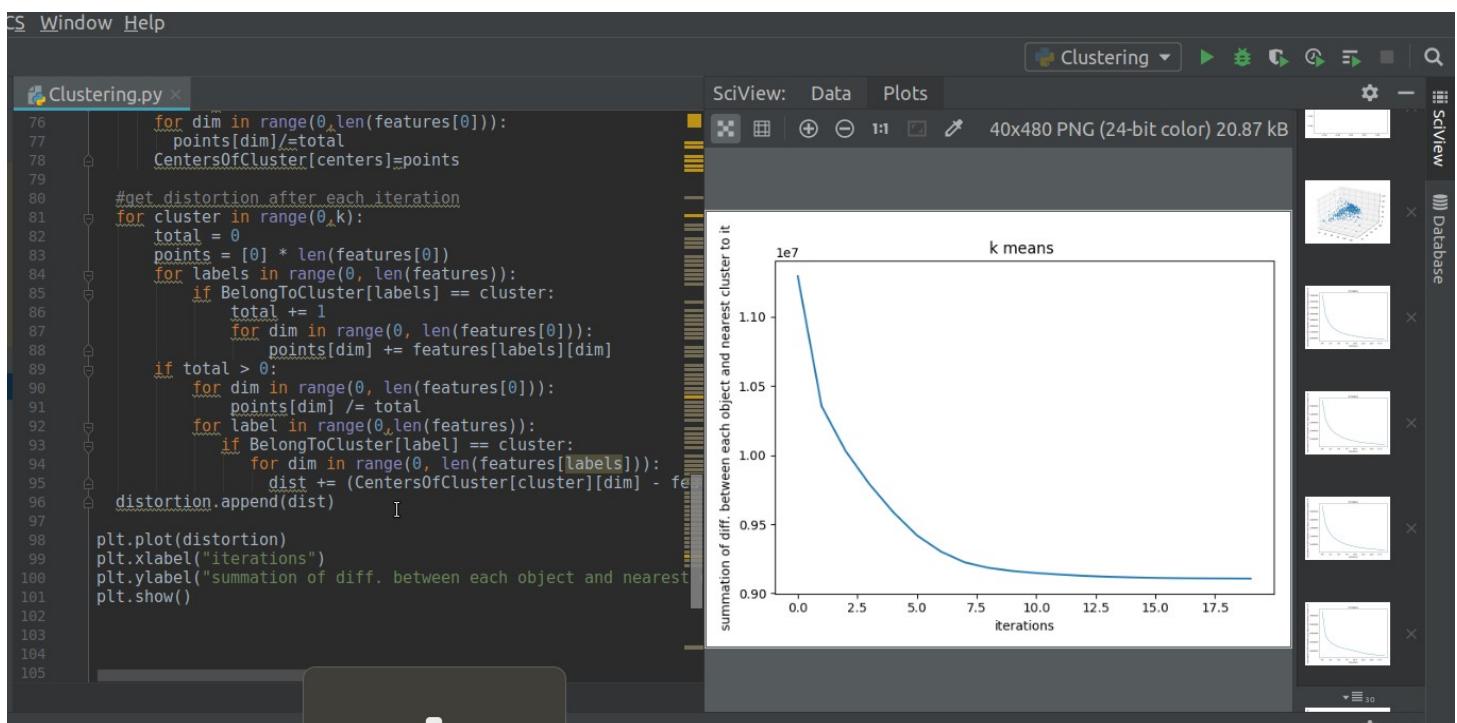
So we calculate our accuracy according to the distortion Function as we aim to make it tends to Zero.

What is Distortion Function?

It is how long our image far away from its nearest Cluster's Center. So we get summation for this Distance at each iteration.

It can be represented using variance Function.

At k=10:



- it's start to get near to x-axis at iteration:7 so it get to be stable and more accurate after the 7th iteration

At k=50:

Activities PyCharm Professional Edition 2:42 PM

untitled1 [-/PycharmProjects/untitled1] - .../Clustering.py [untitled1] - PyCharm

File Edit View Navigate Code Refactor Run Tools VCS Window Help

untitled1 Clustering.py

Project ~/PycharmProjects/untitled1

- untitled1
- cifar
- venv
- bin
- include
- lib
- lib64
- pip-selfcheck.json
- pyvenv.cfg
- cifar-10-python.tar.gz
- Clustering.py
- data_batch_1

External Libraries

Scratches and Consoles

Clustering.py

```
40     ⌄ np.array(features)
41     #clusters --> 10^3(10 clusters each have average of R.B.G
42     CentersOfClusters[]
43     #portion[1]
44     Ks50
45     #get Random positions For Clusters
46     for i in range(0,k):
47         CentersOfCluster.append(features[i])
48     #number of max iterations
49     iterations=20
50     #number of clusters
51
52     BelongToCluster=[0]*(len(features))*k
53     for it in range(0,iterations):
54         dist = 0
55         #get distance between labels and clusters
56         for labels in range(0,len(features)):
57             mnDis=16
58             for centers in range(0,k):
59                 ed=0
60                 for dim in range(0,len(features[labels])):
61                     ed+=(features[labels][dim] - CentersOfCluster[cluster])
62                 if ed < mnDis:
63                     mnDis=ed
64                     BelongToCluster[labels]=clusters
65
66     #get new centers
67     for centers in range(0,k):
68         totals=0
69         points=[0]*len(features[0])
```

SciView Data Plots

40x480 PNG (24-bit color) 24.74 kB

k means

Iterations	Summation of edit
0.0	2,250,000
2.5	1,950,000
5.0	1,850,000
7.5	1,800,000
10.0	1,750,000
12.5	1,700,000
15.0	1,680,000
17.5	1,660,000

summation of edit between each object and nearest cluster to it

iterations

Run: Clustering x

/home/merna/PycharmProjects/untitled1/venv/bin/python /home/merna/PycharmProjects/untitled1/Clustering.py

Process finished with exit code 0

Z-Structure

Z-Favorites

4: Run 6: TODO 7: Terminal 8: Python Console

58:32 LF: UTF-8 4 spaces Event Log

Looks like you're using NumPy Would you like to turn scientific mode on? Use scientific mode Keep current layout...

- it get closer to x-axis at iteration(13) so it get to be stable and more accurate after 13 iterations.

NOTE:

if we increase the number of clusters more than the normal as $k=50$ our function become less accuracy as each cluster has smaller amount of points belong to it and also if we make number of clusters so small as $k=2$ our function also become less accuracy as the distance increase.

2) Are the results wildly different for different restarts and/or different K?

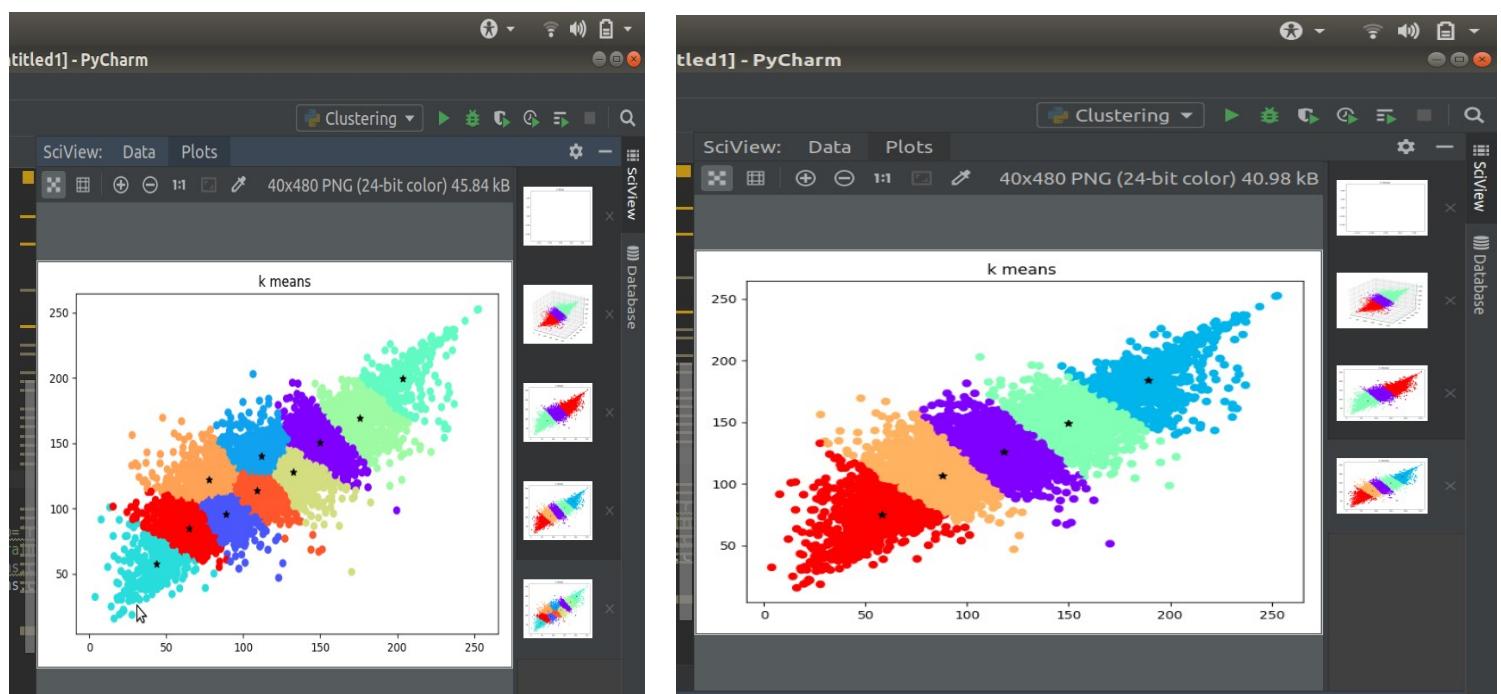
Yes, as we change our K accuracy Change.

Plotting Of clusters at difference Start and difference K .

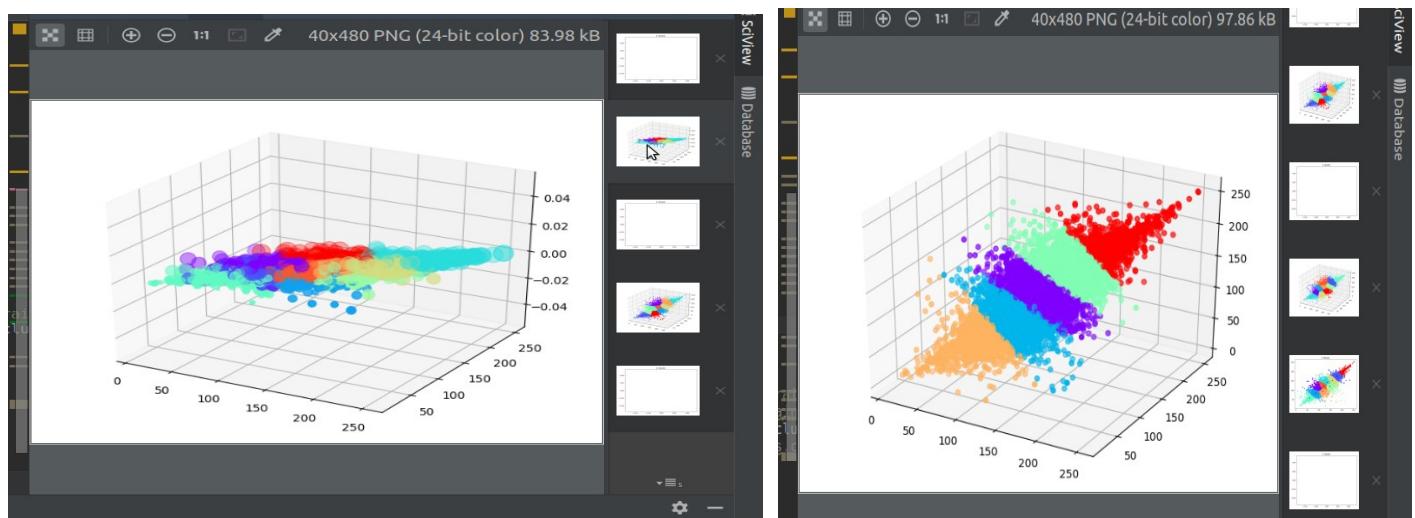
For Difference Restart: As K means get local maximum if we reach that every point belong to same cluster every iteration and nothing change, that doesn't mean that we reach the best case we have to start randomly and different restart to check for the best case.

plotting if we choose 2 features only (average of blue color, average of green color).

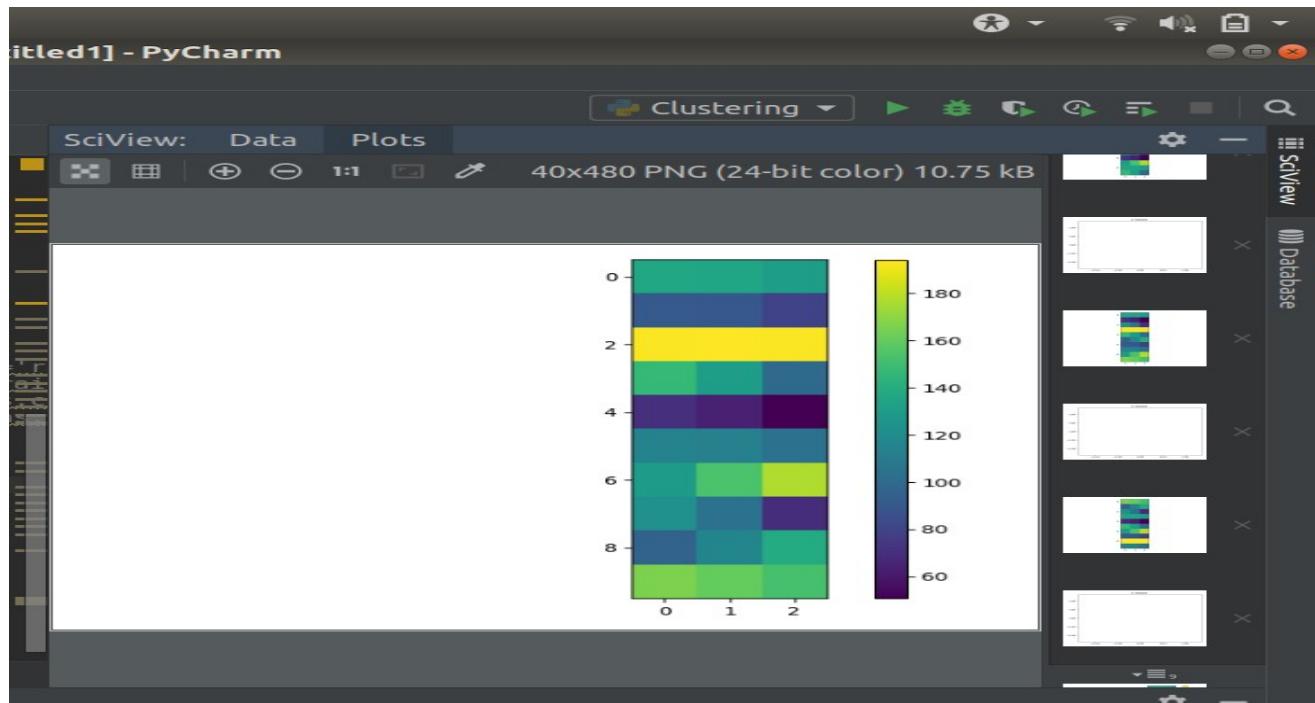
Black points are our cluster's Center.



Plotting if we choose 3 features (average of red color, average of blue color,average of green color).



3)What do the mean images look like? Provide screenshots with comments.

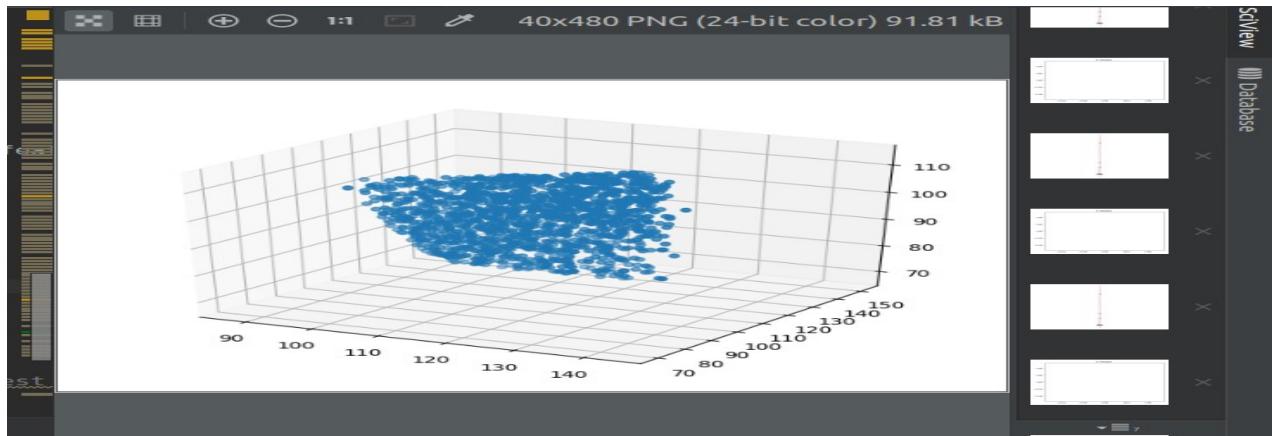


Note:

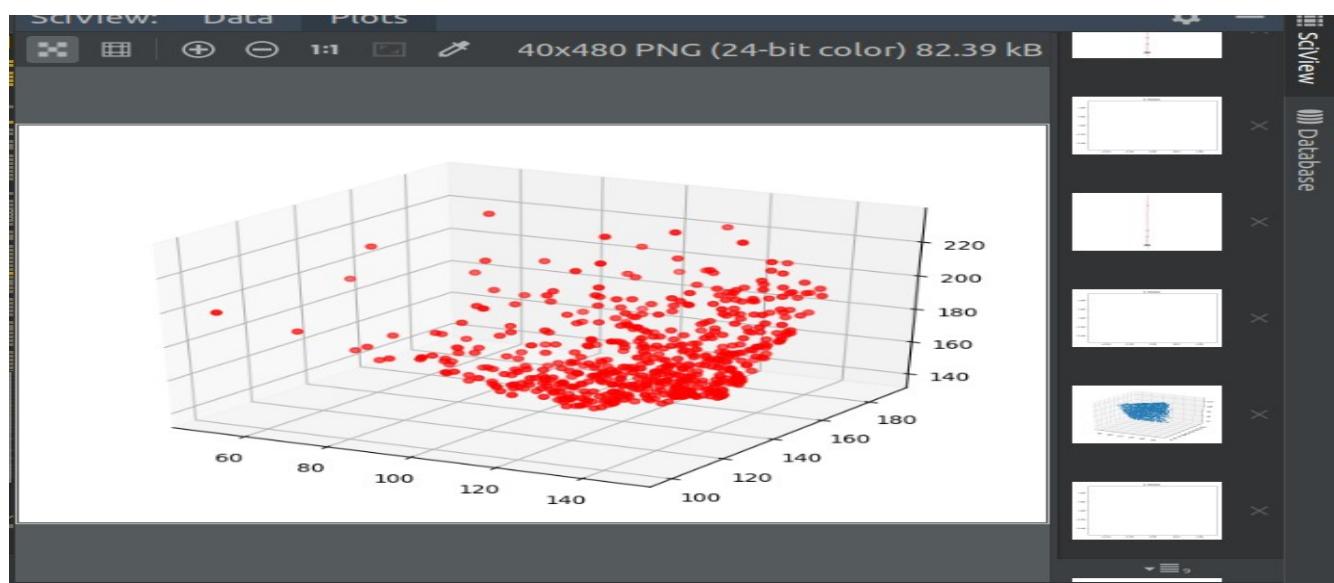
- K means get accuracy 40% with Cifar dataset as Cifar need deepLearning to get accuracy 90%.
- So it doesn't retrieve accurate images.

4) What are some representative images from each of the clusters?

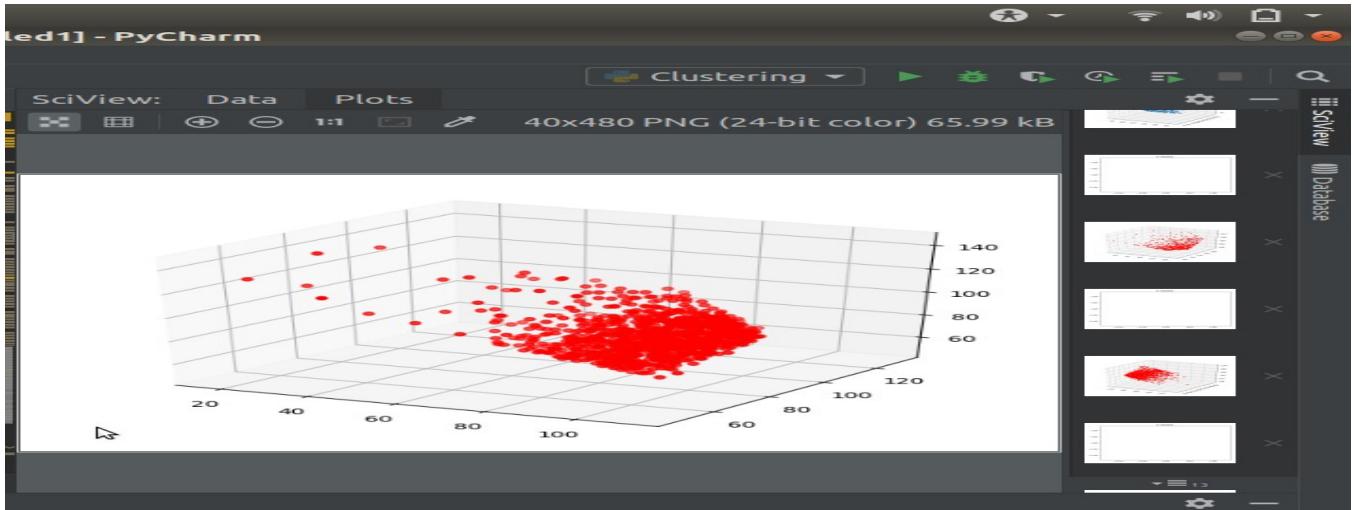
Cluster0:



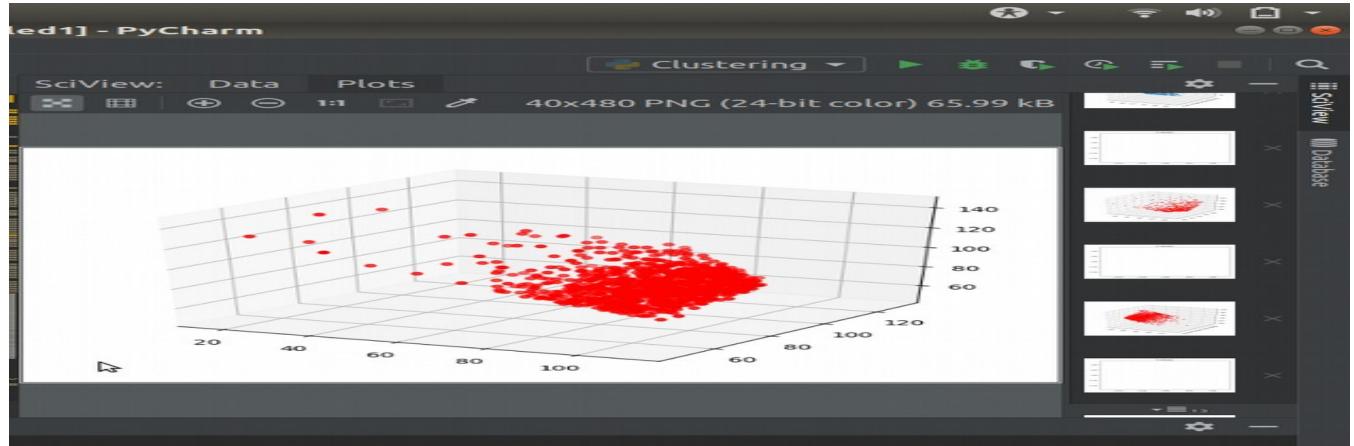
Cluster1:



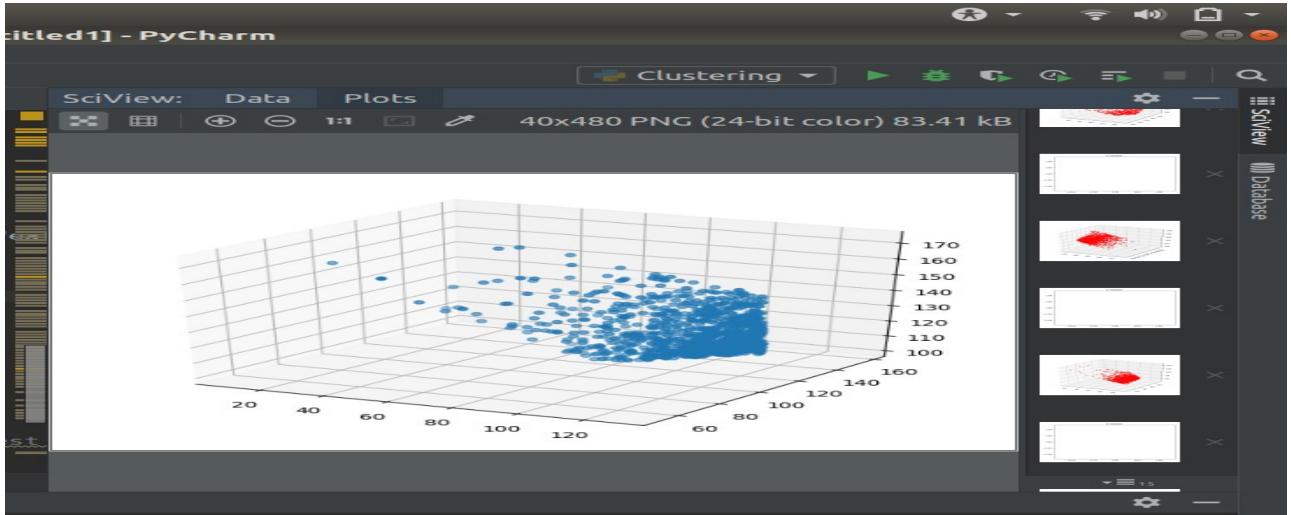
Cluster2:



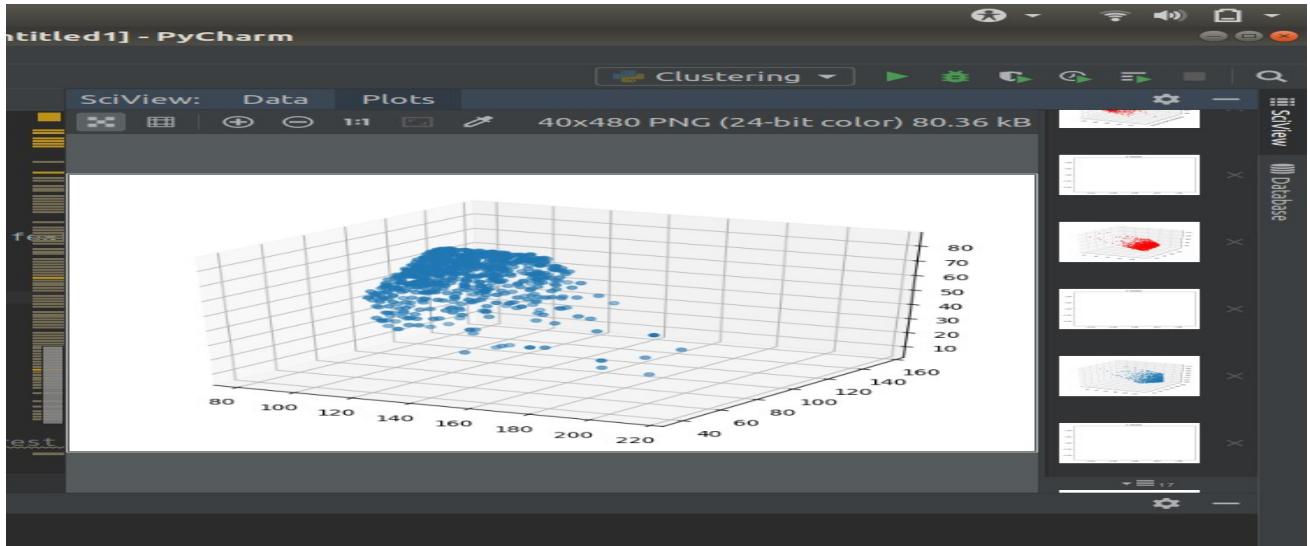
Cluster3:



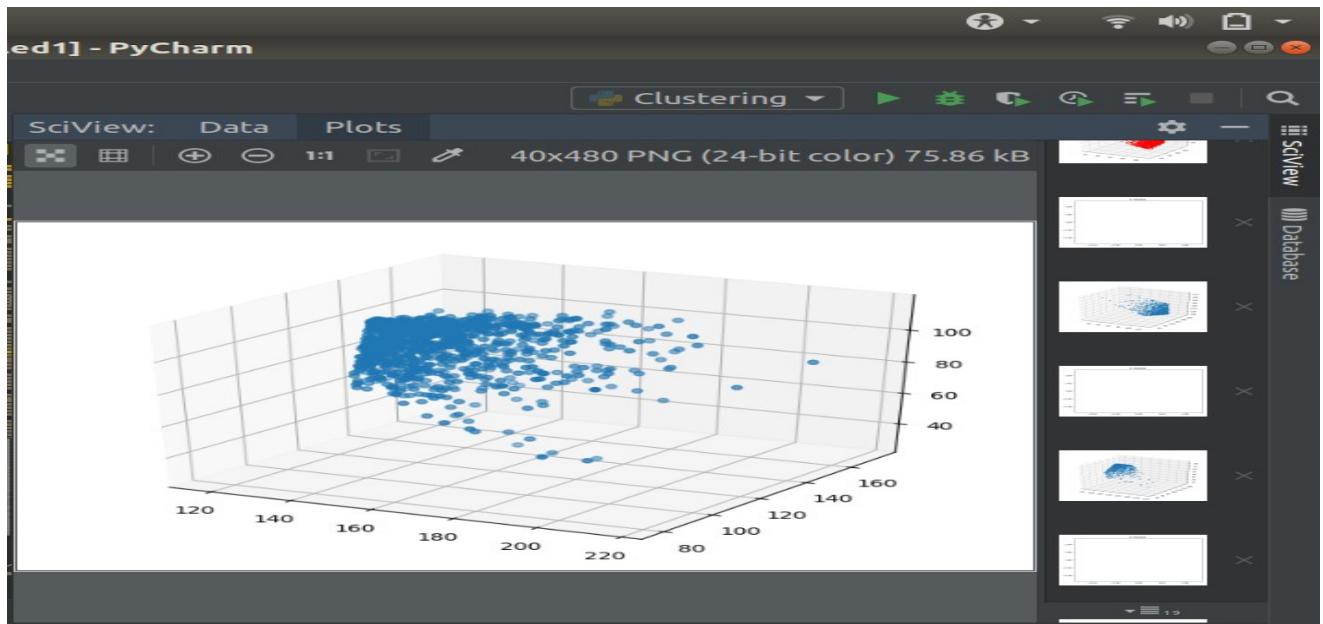
Cluster4:



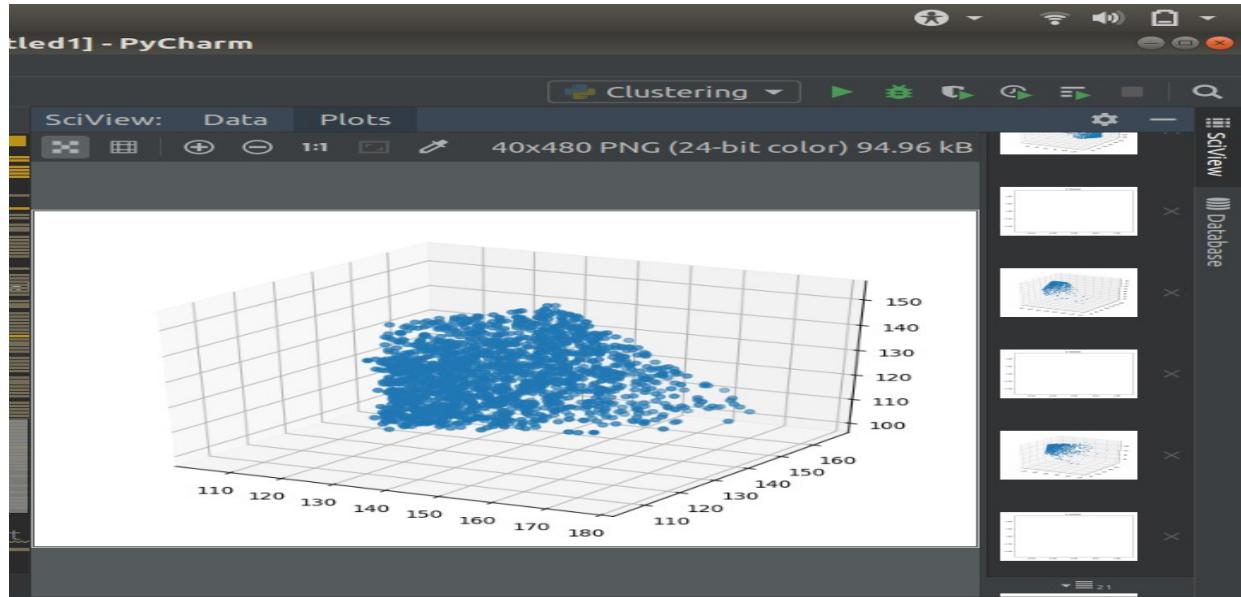
Cluster 5:



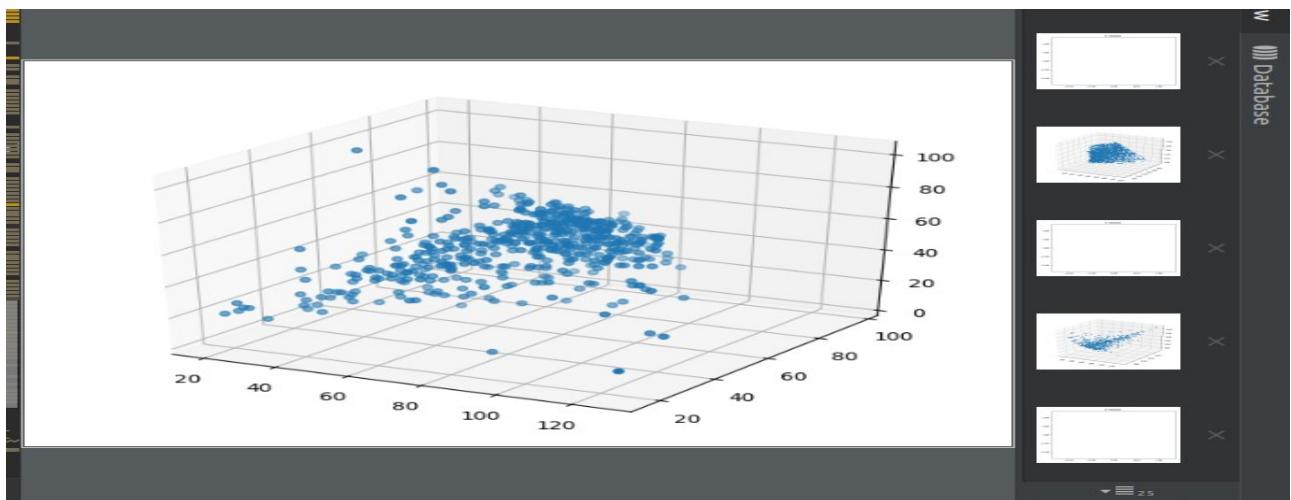
Cluster 6:



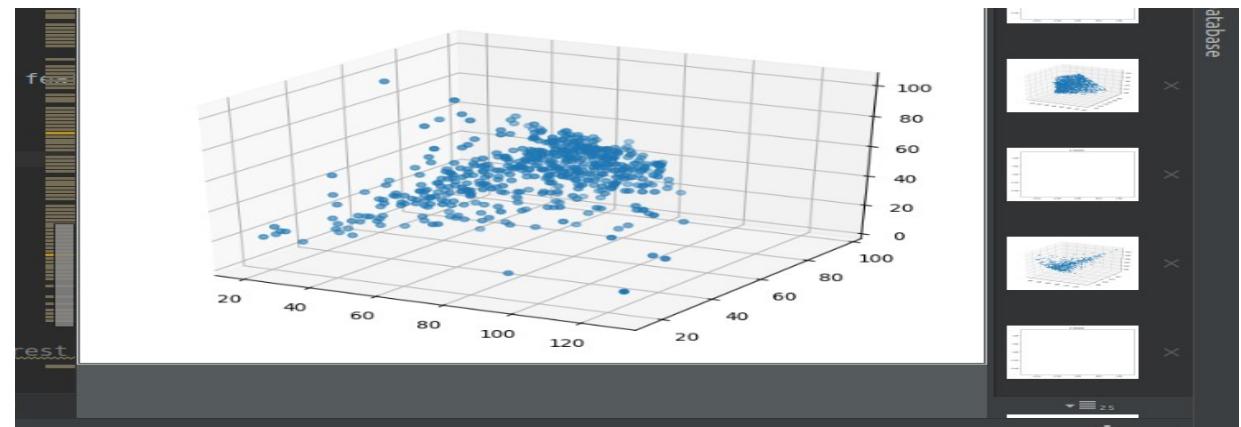
Cluster 7:



Cluster 8:

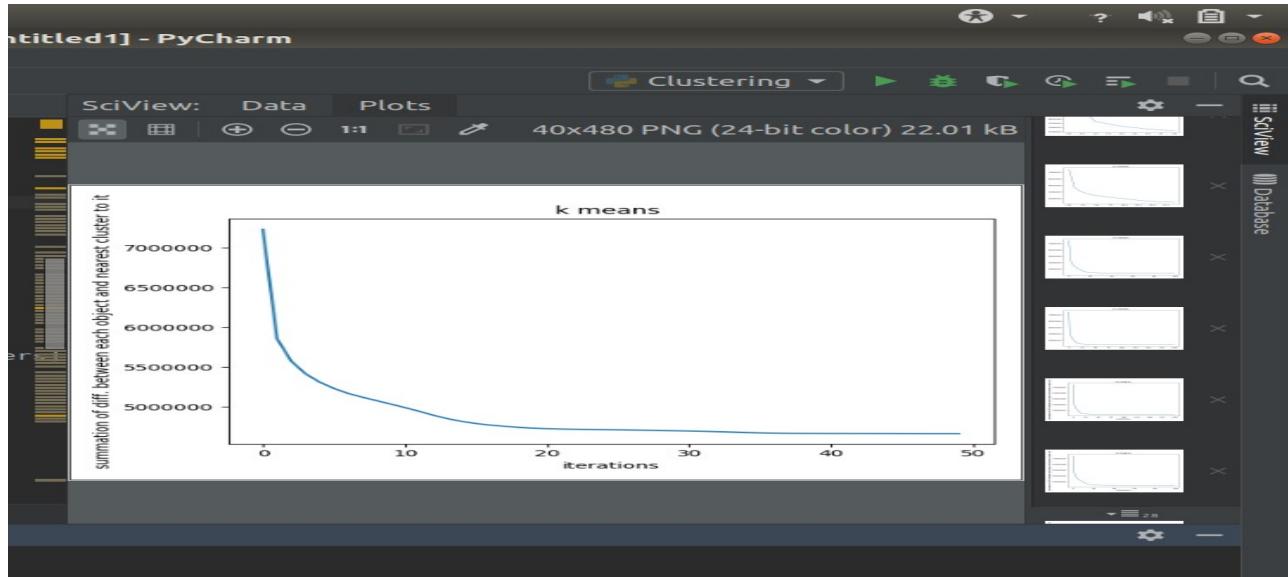


Cluster 9:

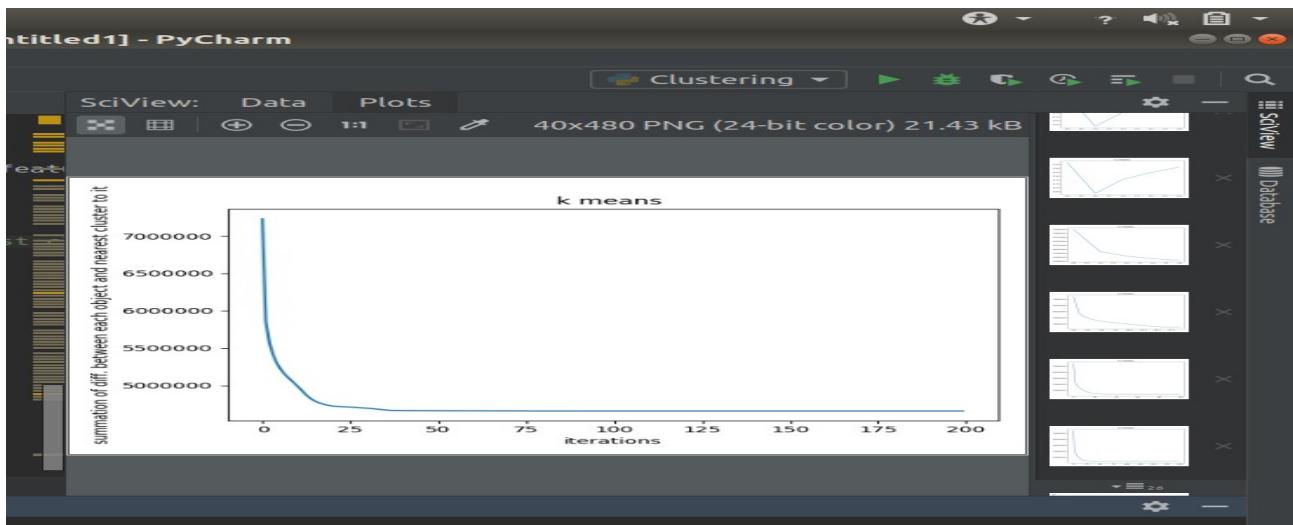


5) Plot the K-Means objective function (distortion measure) as a function of iteration and verify that it never increases.

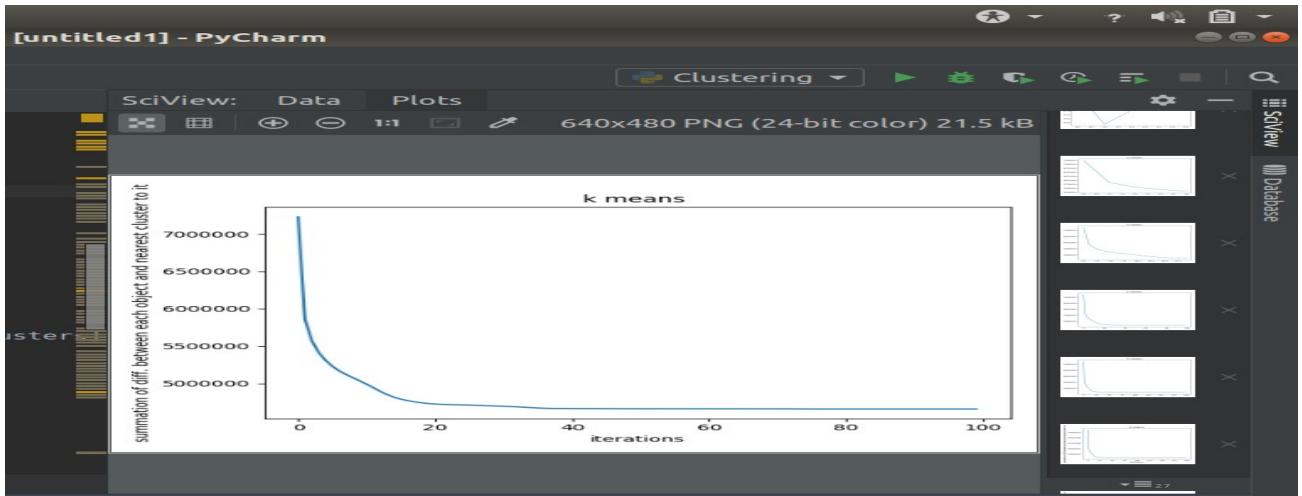
Iteration = 50 , k =10:



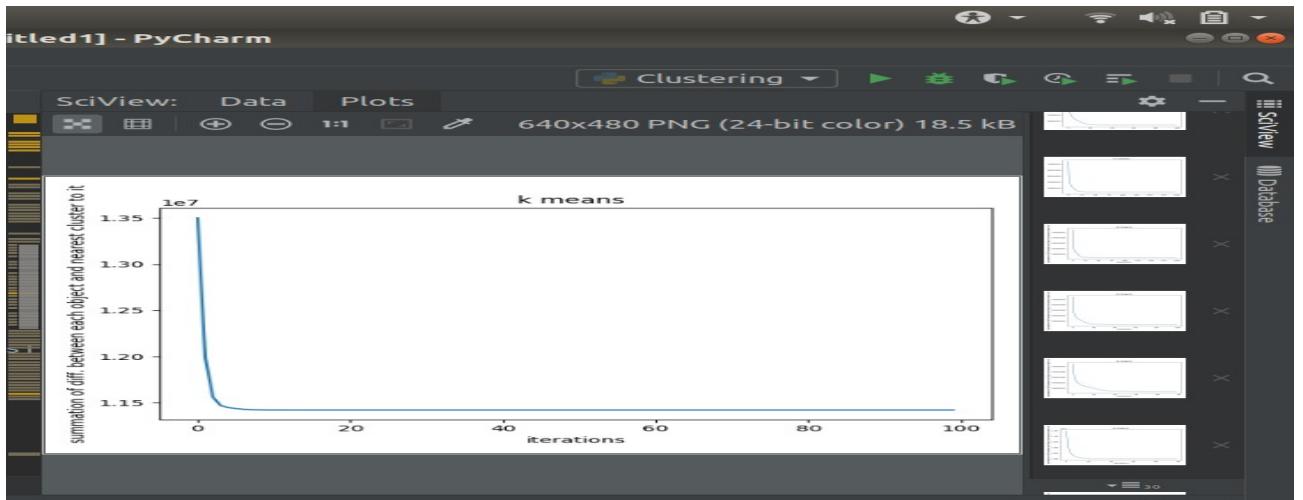
Iteration = 200, k=10:



Iteration = 100 , k= 10:



Iteration = 100, k= 3 :



Iteration = 50, k = 5:

