

IMAGE CLUSTERING

LINK TO COLAB FILE:-

https://colab.research.google.com/drive/1rFlxnNU_EOw_A56cH2HR_C9pgLPC4fmf?usp=sharing

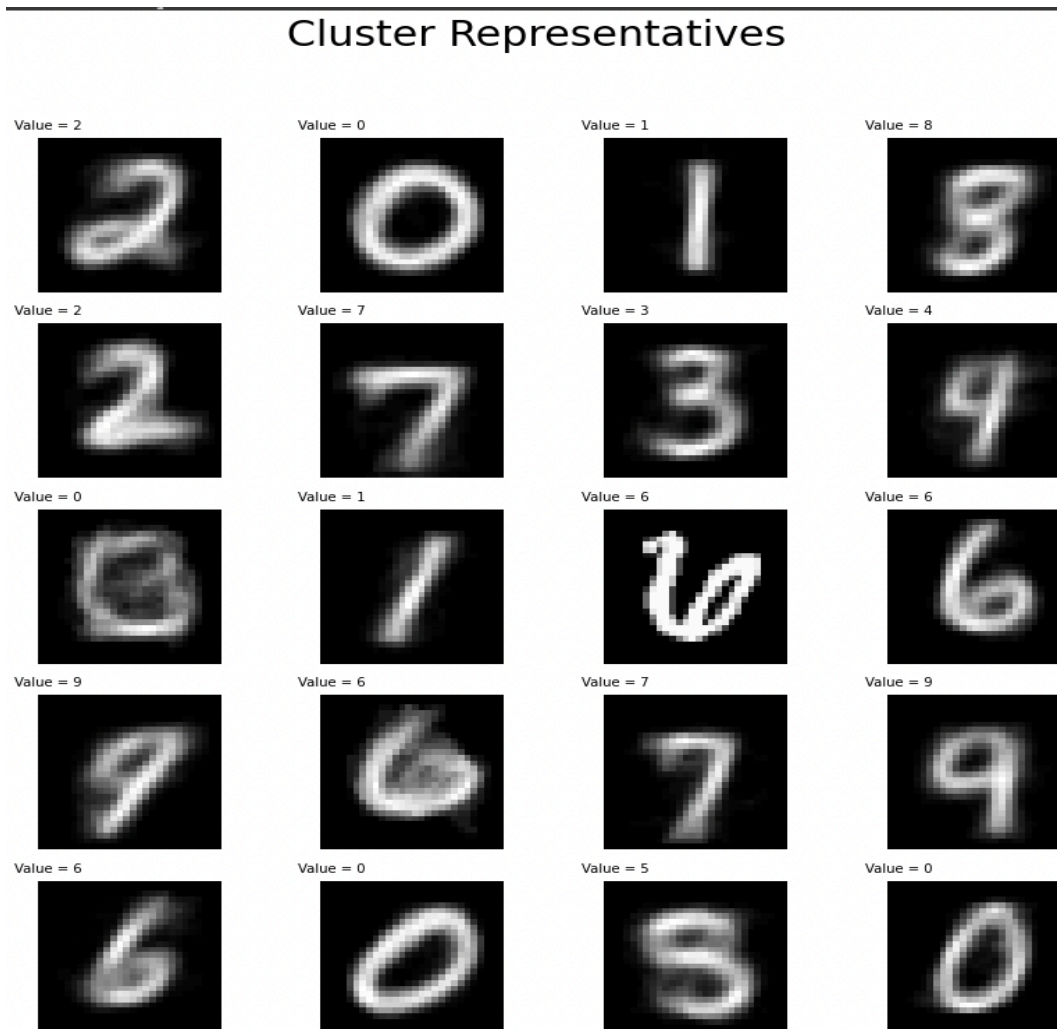
As we are selecting 100 images of each digit we have 1000 training samples in all, so $N=1000$. Each image is 28×28 pixels which is 784 features. Thus $n=784$. Convergence criteria is that successive J_Clust values should differ by less than $1e-6$.

$N = 1000$

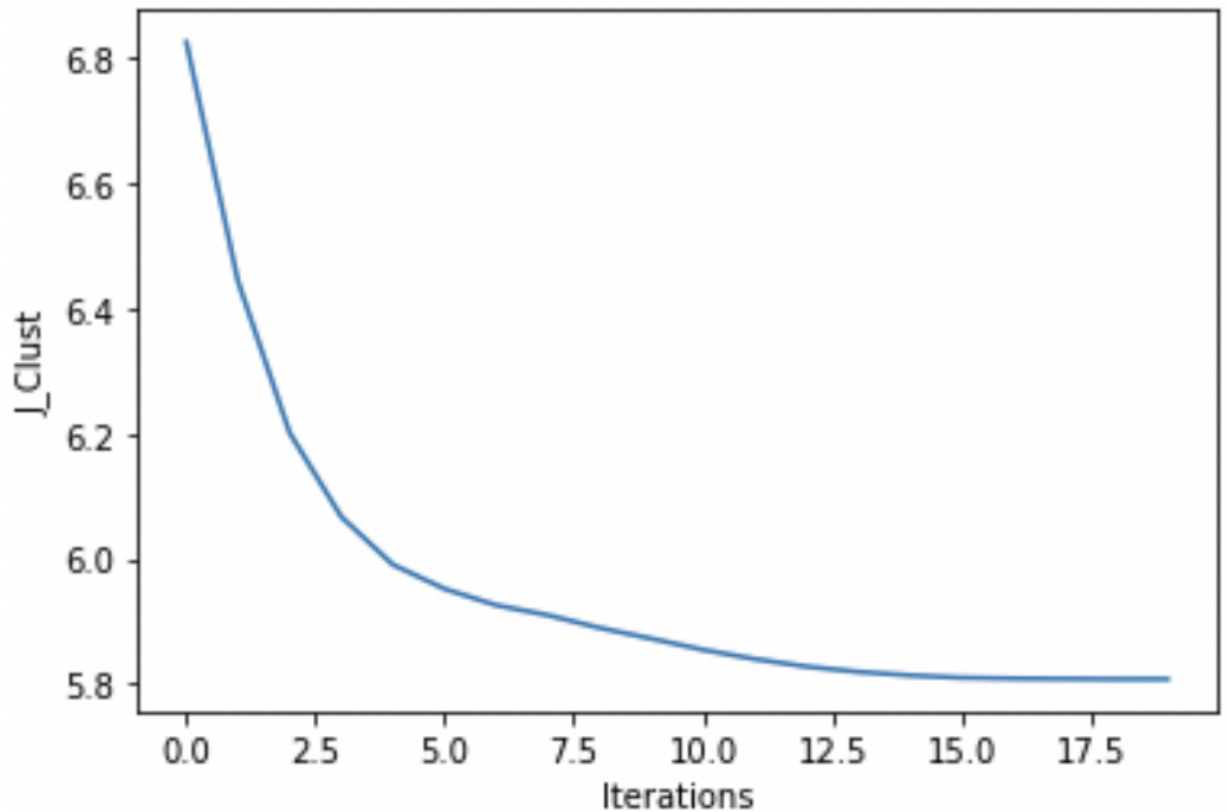
$n = 784$

Case (i) Random initialization of cluster representatives

Converged after 48 iterations. Cluster Representatives:-



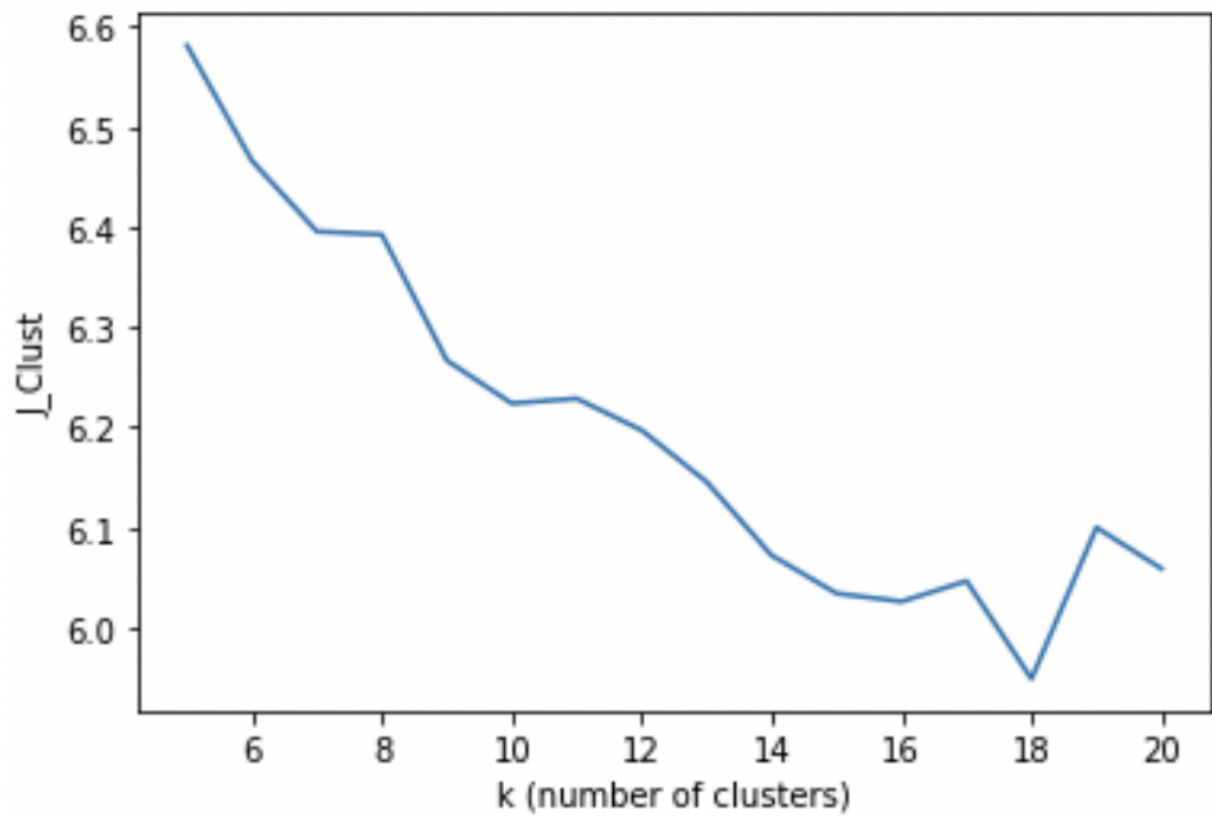
a) Variation of J_Clust with iterations



b) Test Accuracy is 68.0% and Train Accuracy is 64.40%

c) Variation of J_Clust with k (number of clusters)

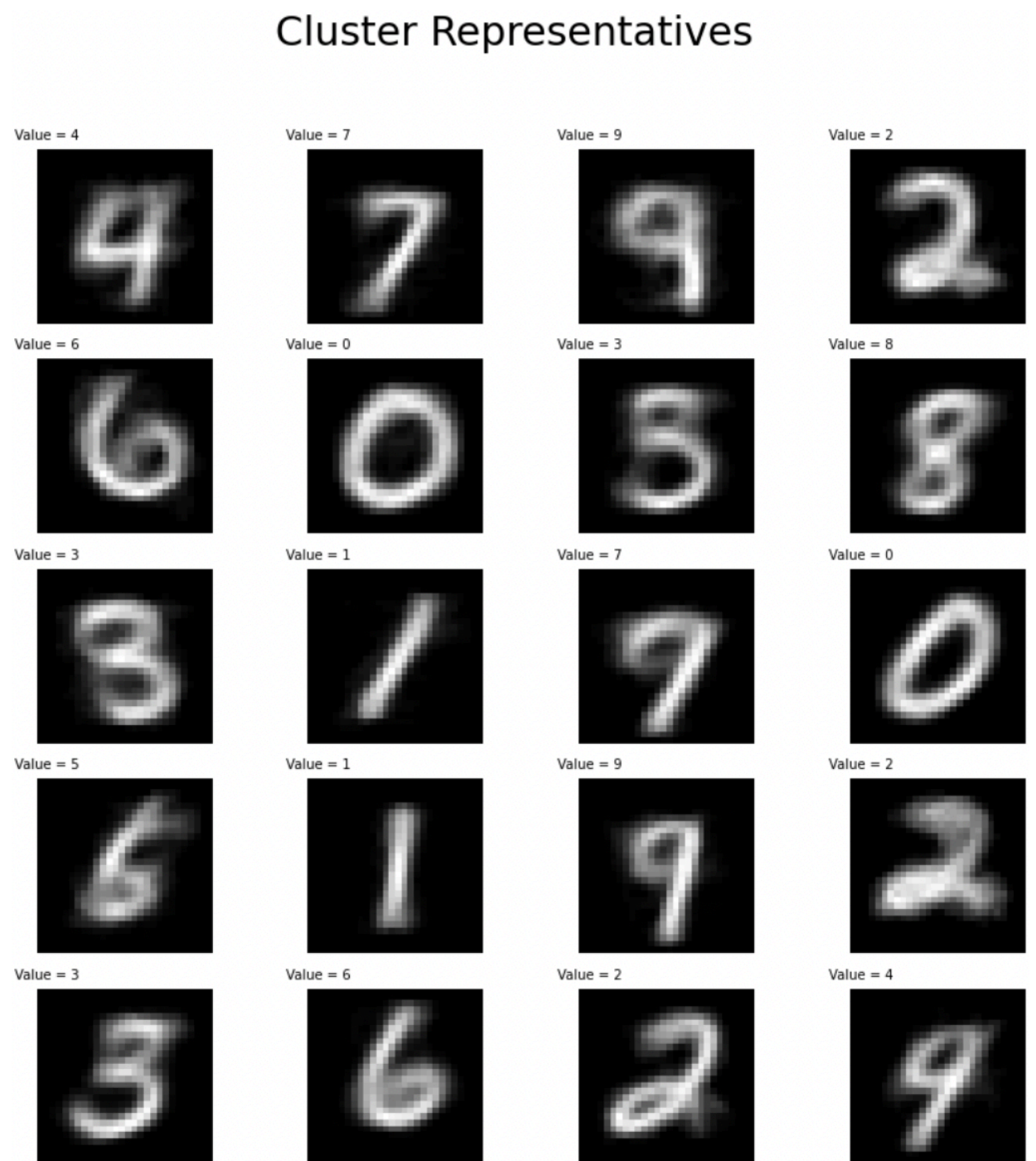
```
k = 5, J_clust = 6.581260896121875
k = 6, J_clust = 6.466127462591003
k = 7, J_clust = 6.395571088328437
k = 8, J_clust = 6.392430836900881
k = 9, J_clust = 6.26688595952364
k = 10, J_clust = 6.224017680188111
k = 11, J_clust = 6.228877296580053
k = 12, J_clust = 6.197223348250238
k = 13, J_clust = 6.145686704369392
k = 14, J_clust = 6.072384100309183
k = 15, J_clust = 6.034401741846443
k = 16, J_clust = 6.026332508460405
k = 17, J_clust = 6.046822889380301
k = 18, J_clust = 5.949277423926605
k = 19, J_clust = 6.1002297220266035
k = 20, J_clust = 6.059174088967577
```



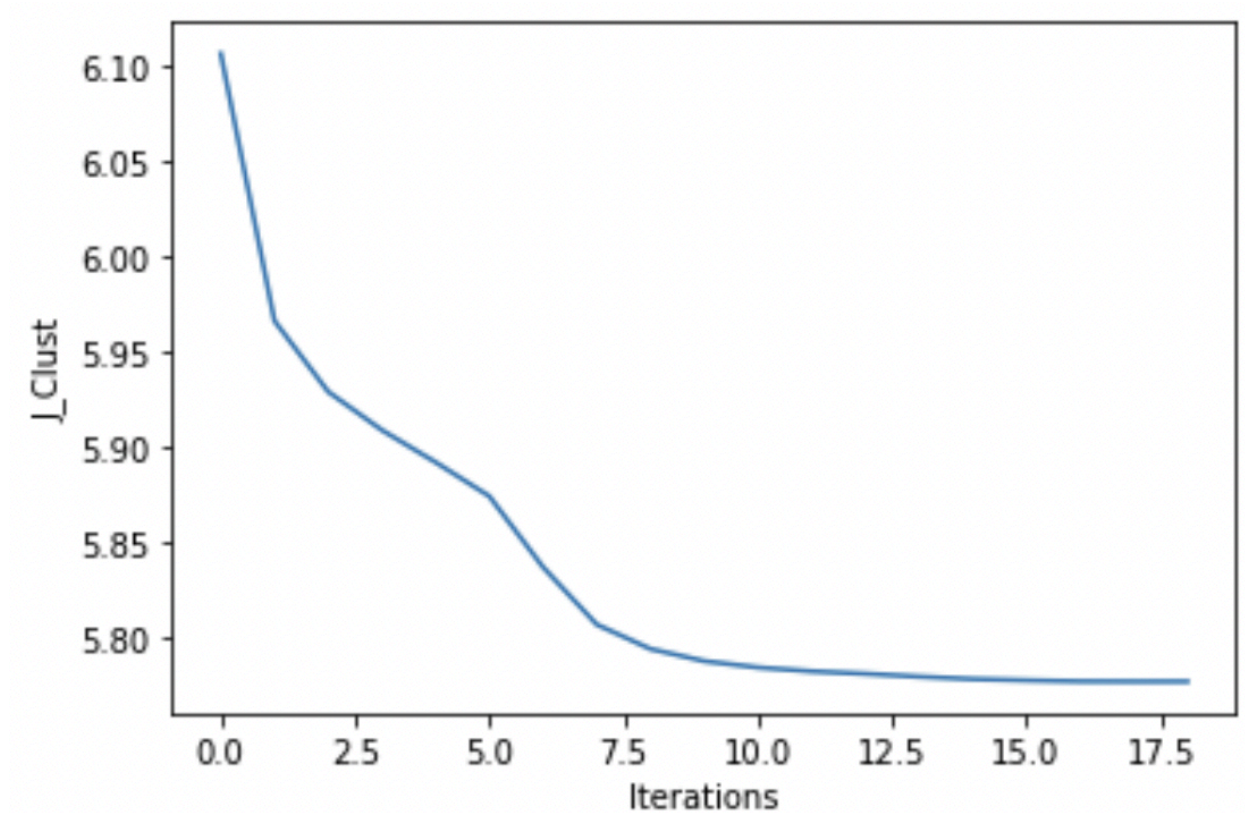
Min J_{clust} = 5.949277423926605 for k = 18

Case (ii) Choose cluster representatives from the given data set

Converged after 20 iterations. Cluster Representatives



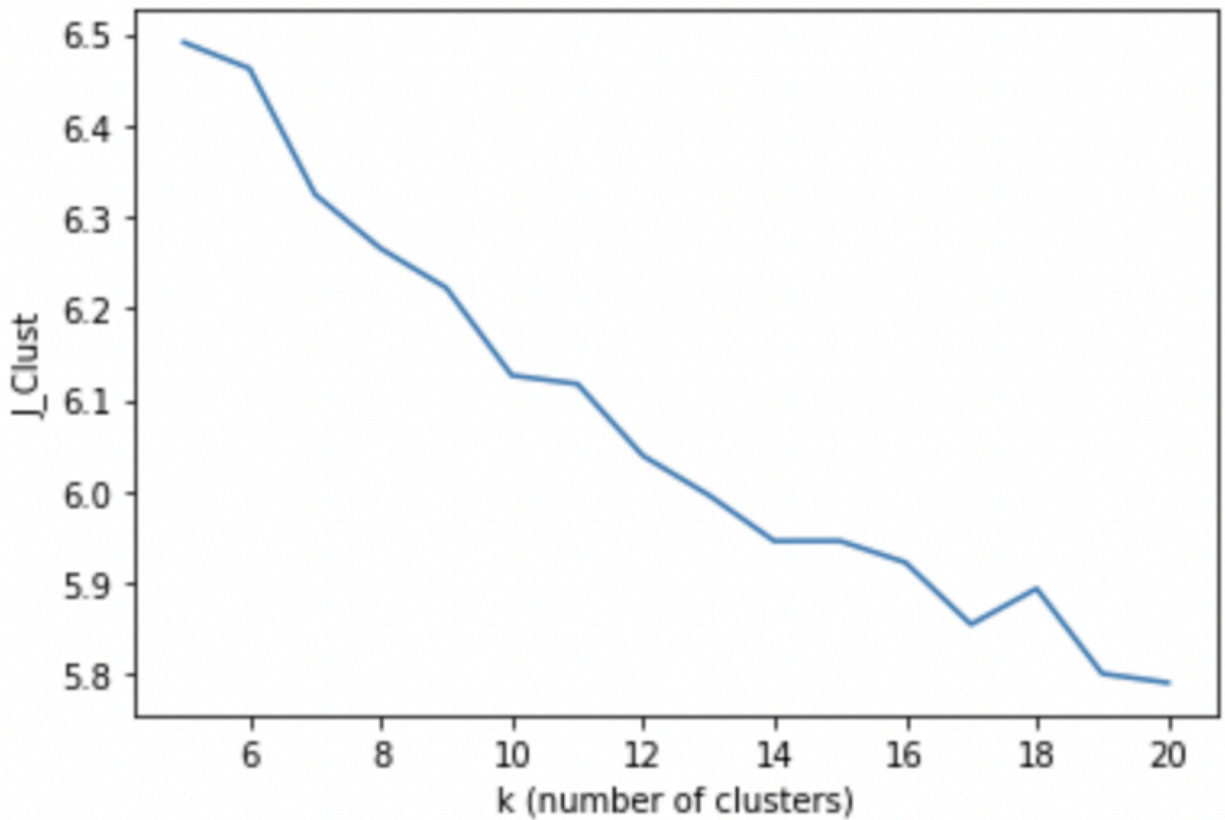
a. Variation of J_Clust with iterations



b. Accuracy is 74 %

c. Variation of J_Clust with k (number of clusters)

```
k = 5, J_clust = 6.490830181374058
k = 6, J_clust = 6.4622847327888016
k = 7, J_clust = 6.324576122306834
k = 8, J_clust = 6.265369819452733
k = 9, J_clust = 6.22225903939423
k = 10, J_clust = 6.126557642399606
k = 11, J_clust = 6.11723199074083
k = 12, J_clust = 6.0386687773903835
k = 13, J_clust = 5.995731399936434
k = 14, J_clust = 5.9454466082565505
k = 15, J_clust = 5.945566023102394
k = 16, J_clust = 5.921975325513757
k = 17, J_clust = 5.854224551668884
k = 18, J_clust = 5.893343848583136
k = 19, J_clust = 5.800335170526334
k = 20, J_clust = 5.790413544305315
```



Min $J_{\text{clust}} = 5.790413544305315$ for $k = 20$

In random initialisation the optimal value of k is 18 and for dataset initialisation the optimal value of k is 20.

Different styles of writing numbers leads to more clusters providing better accuracy. Each different style of writing a number goes into a different cluster. The choice of the initial condition primarily affects the number of iterations required for convergence.

With the initial cluster representatives from the dataset we observe a much faster convergence as compared to random initialisation. Moreover, we observe a better accuracy and slightly

lower J_Clust value in case of initialisation from the dataset. Also, in case of initialisation from the dataset the visualization shows the final cluster representatives are much smoother and closer to the actual numbers compared to the random initialisation.