# HW-2

## Mahesh Desai

**1.Assignment**

The goal of this assignment is to implement K Means with CISFR-10 dataset and find its Silhouette score and DUNN score. Here I have implemented K Means Clustering Algorithm from basic by using various methods to improve the accuracy of the model .

**2.Dataset**

We have the CISFR-10 dataset . 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 two parts as said train and test. Both the training set and the test set contains images from the 10 class equally. Here I while executing the program I have taken only 8000 images from the test set due to limitation of RAM.

The dataset contains the following classes.

1) Airplane
2) Automobile
3) Bird
4) Cat
5) Deer
6) Dog
7) Frog
8) Horse
9) Ship
10) Truck

Outcome: There will be a cluster of images belonging to each group.

**3. Editor**

I have use VS code to write my code.

**4. Description**

Here I have imported all the images from the keras . I have  selected few images from the test array to test my code.

a. Accuracy-This accuracy check will be used to check hoe accurectly our model is able to make cluster of the images based on their features. Here we are using two accuracy checking algorithms one is silhouette score and the second one is Dunn's score.

1) Silhouette score- The Silhouette Coefficient is calculated using the mean intra-cluster distance (a) and the mean nearest-cluster distance (b) for each sample. The Silhouette Coefficient for a sample is (b - a) / max(a, b). To clarify, b is the distance between a sample and the nearest cluster that the sample is not a part of. Note that Silhouette Coefficient is only defined if number of labels is 2 <= n_labels <= n_samples - 1.

   The best value is 1 and the worst value is -1. Values near 0 indicate overlapping clusters. Negative values generally indicate that a sample has been assigned to the wrong cluster, as a different cluster is more similar.

2) Dunn's score-The Dunn index (DI) (introduced by J. C. Dunn in 1974), a metric for evaluating clustering algorithms, is an internal evaluation scheme, where the result is based on the clustered data itself. Like all other such indices, the aim of this Dunn index to identify sets of clusters that are compact, with a small variance between members of the cluster, and well separated, where the means of different clusters are sufficiently far apart, as compared to the within cluster variance. Higher the Dunn's index better cluster.

b. Procedure- Here we have used CISFR-10 dataset. Now the first step of the algorithm is to extract features from the images. I have passed all the image list to the model to extract features .I have used Keras VGG16, VGG19, Xception, InceptionV3 model to extract features from all the images . We get the output as the feature List . Then we have to predict the cluster. For that we have to choose random clusters points, number of classes and number of iteration that we want to complete the entire in. I have taken classes as 10 as we have total 10 class. Then I have chosen random 10 centroids Then calculated the distance between the centroid and all the data points. I have used Euclidean distance to measure the distance between the points. And iterated the entire procedure till the centroids of the points don't change. After calculating the centroids and total label points we have calculate the Silhouette score and Dunn's score.

**5.Concusion**-

| MODELS | ITERATIONS | SILHOUETTE SCORE | DUNN'S SCORE |
|---|---|---|---|
| VGG16 | 2000 | 0.03275114 | 0.060244985 |
| VGG19 | 2000 | 0.037474357 | 0.04040027 |
| XCEPTION | 2000 | 0.25955316 | 0.03484359 |
| INCEPTIONV3 | 2000 | 0.0038415287 | 0.19932023 |

Here I have submitted the code for XCEPTION model as it performed better than other model.

# PART 2

## 1.Assignment

The goal of this assignment is to implement AutoEncoder with CISFR-10 dataset and find its Silhouette score on Encoded images. Here I have implemented Autoencoder Algorithm from basic by using Convolution Network model to improve the accuracy of the model and give proper generated images .

## 2.Dataset

We have the CISFR-10 dataset. I have imported the dataset from the keras . 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 two parts as said train and test. Both the training set and the test set contains images from the 10 class equally. Here I while executing the program I have taken only test images due to limitation of RAM. I have taken 10000. While calculation Silhouette score I have taken 10000 images due to limitation of RAM.

The dataset contains the following classes.

1) Airplane
2) Automobile
3) Bird
4) Cat
5) Deer
6) Dog
7) Frog
8) Horse
9) Ship
10) Truck

Outcome: The outcome will be regenerated images from the Autoencoder model.

## 3. Editor

I have use VS code to write my code.

## 4. Description

Here I have imported all the images from the keras . I have  selected few images from the test array to test my code.
Autoencoder is a type of neural network where the output layer has the same dimensionality as the input layer. Here is have kept the dimensions of the input image and the output image same which is 32x32x3. Autoencoder tries to replicate the input images unsupportively. It trains of the input images and tries to replicate them. There are two parts of Autoencoder a) Encoder b) Decoder.
   a) Encoder- An encoder is a feedforward, fully connected neural network that compresses the input into a latent space representation and encodes the input image as a compressed representation in a reduced dimension. The compressed image is the distorted version of the original image.
      Here I have used Convolution 2D 3 layers to do encoding with 256 neurons respectively. To compress the image I have used MaxPooling layers. The input image size is 32x32x3 and when the image is encoded the size is reduced to 4x4 .This encoded image array is used to

calculate the Silhouette score. We calculate the distance between the feature points sand the labels and pass it to the Silhouette function to calculate the score.

b) Decoder- Decoder is also a feedforward network like the encoder and has a similar structure to the encoder. This network is responsible for reconstructing the input back to the original dimensions from the code. This decoder is used to reconstruct the images. The decoder used the compressed images from the encoder to generate similar image. I have used 3 layers of Convolution 2D with 256 neurons each. To upscale the image I have used UpSampling2D , it upscale the image while generating it . The shape of the decoded image is 32x32x3.

So the shape of the input image and output image is same . Then the decoded images are put in the autoencoder model to train . I have used optimizer as "adam" and loss as " mean square error", with batch size of 1000 and iterations 15.
Then I have plotted 10 images to show the original and the generated images.

**5) Procedure**- Here we have used CISFR-10 dataset.I have imported the dataset from keras. I have used only 10000 images due to Ram limitations. Then I am creating the encoder and decoder layers , to create my model. The I am passing the images to the model to encode it compress it and then back upscale ot and regenerate it. I have put 15 iterations and 1000 batch size. The after the encoding is done the labels are generated and then it is passed to Silhouette function to calculate the score and see the performance of the model which is 0.0854. Then the images are displayed along with the original images.

**6) Conclusion-**