# ICP 6

**1. Apply K means clustering in this data set provided below:**

**https://umkc.box.com/s/s15r7m0gnxu7b1s2kaobvc5w7da2nc1c**

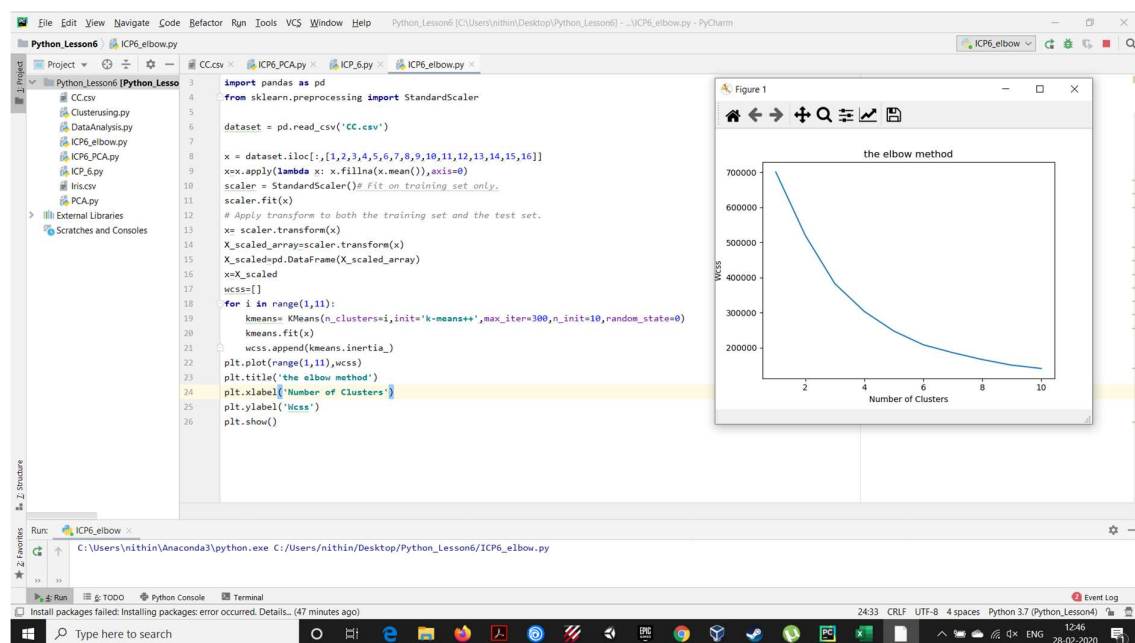●**Remove any null values by the mean.**

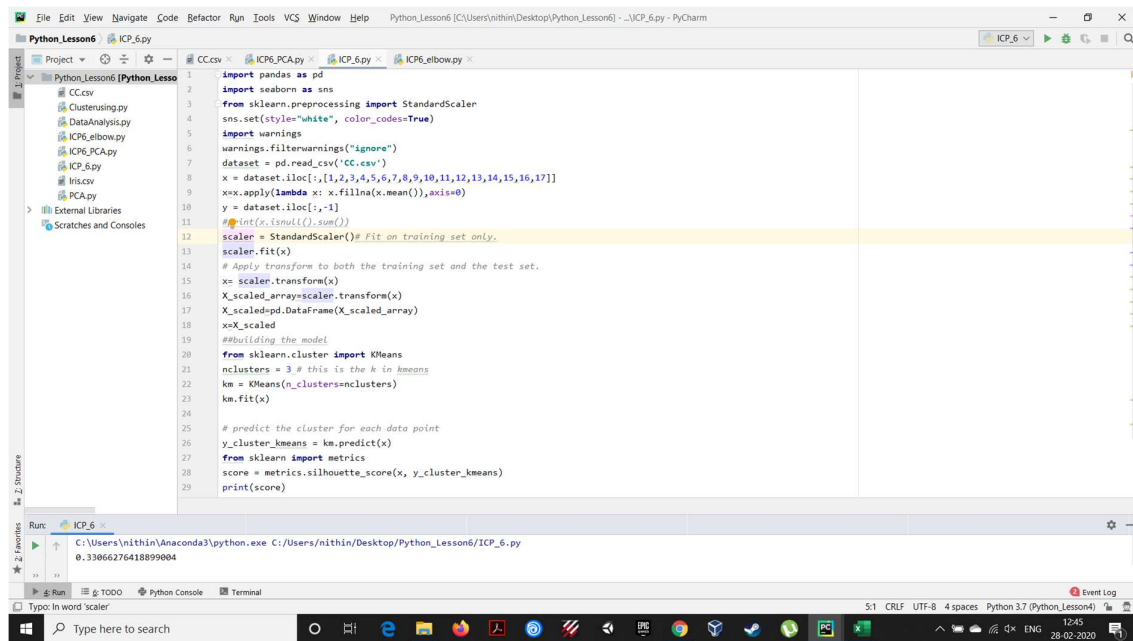●**Use the elbow method to find a good number of clusters with the KMeans algorithm**

**2.Calculate the silhouette score for the above clustering**

**3.Try feature scaling to see if it will improve the Silhouette score**

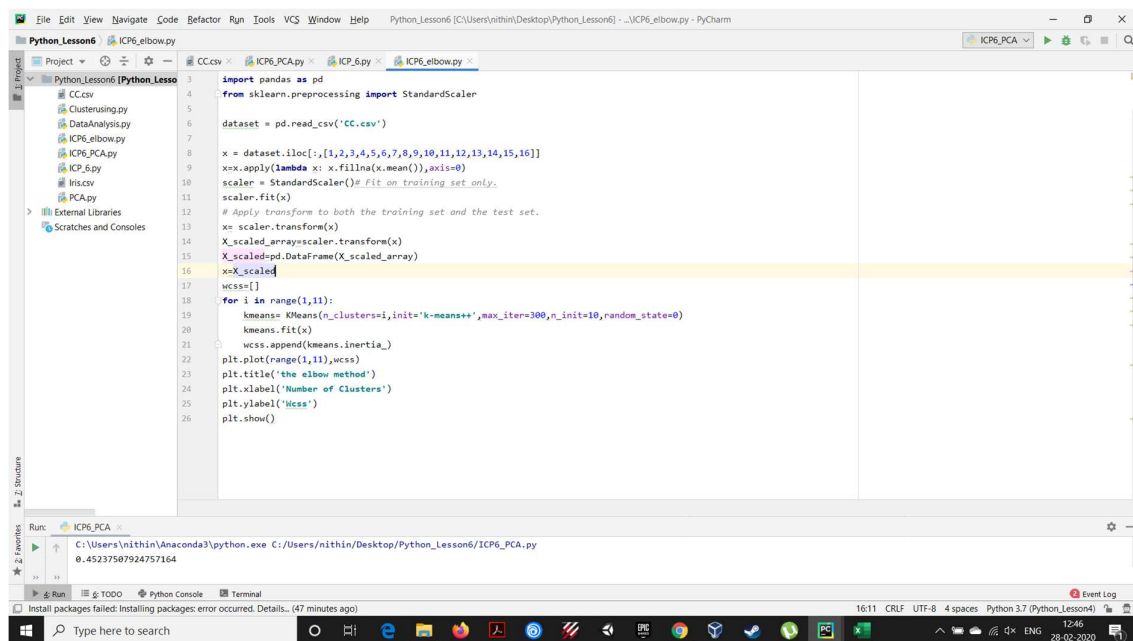**4.Apply PCA on the same dataset. Data Description can be found in**

**https://umkc.box.com/s/cjeenva7pyj6s0vz8s8zlpbwxtmg6d17*****

**2. Bonus points Apply kmeans algorithm on the PCA result and report your observation if the score improved or not?**

**BY**

**DUKKIPATI SRI SAI NITHIN CHOWDARY**

**CLASS ID: 4**