Creating the dataset import pandas as pd ratings =[['John',5,5,2,1],['Mary',4,5,3,2],['Bob',4,4,4,3],['Lisa',2,2,4,5],['Lee',1,2,3,4],['Harry',2,1,5, ratings Out[3]: [['John', 5, 5, 2, 1], ['Mary', 4, 5, 3, 2], ['Bob', 4, 4, 4, 3], ['Lisa', 2, 2, 4, 5], ['Lee', 1, 2, 3, 4], ['Harry', 2, 1, 5, 5]] In [4]: titles = ['User', 'Jaws', 'Star Wars', 'Exorcist', 'Omen'] titles Out[5]: ['User', 'Jaws', 'Star Wars', 'Exorcist', 'Omen'] movies=pd.DataFrame(ratings, columns=titles) movies Out[7]: User Jaws Star Wars Exorcist Omen **0** John 5 5 2 1 Mary 4 3 Bob Lisa 2 4 Lee 1 Harry # comments first users- action movies - Jaws, Star Wars # last 3 users - Horror movies # Our goal is to use k-means clusting on the users to identify groups of users with similar # movie preferences k=2from sklearn import cluster data= movies.drop('User', axis=1) Jaws Star Wars Exorcist Omen 5 5 2 1 2 1 2 4 4 4 3 3 2 3 4 1 4 k_means=cluster.KMeans(n_clusters=2, max_iter=50, random_state=1) In [14]: k means Out[14]: KMeans(max_iter=50, n_clusters=2, random_state=1) k_means.fit(data) Out[15]: KMeans(max_iter=50, n_clusters=2, random_state=1) labels=k_means.labels_ labels Out[17]: array([0, 0, 0, 1, 1, 1]) pd.DataFrame(labels,index=movies.User, columns=['Cluster ID']) **Cluster ID** User 0 John Mary 0 Bob Lisa 1 Lee Harry # K-means has assigned first users to one cluster and the last 3 users to the second cluster. # This is consistent with our expectation. #We can display the centroid for each of the two cluster centroid=k means.cluster centers In [24]: centroid Out[24]: array([[4.33333333, 4.66666667, 3. [1.66666667, 1.66666667, 4. , 4.66666667]]) pd.DataFrame(centroid, columns=data.columns) Jaws Star Wars Exorcist **0** 4.333333 4.666667 3.0 2.000000 **1** 1.666667 1.666667 4.0 4.666667 import numpy as np testData=np.array([[4,5,1,2],[3,2,4,4],[2,3,4,1],[3,2,3,3],[5,4,1,4]]) testData Out[28]: array([[4, 5, 1, 2], [3, 2, 4, 4], [2, 3, 4, 1], [3, 2, 3, 3], [5, 4, 1, 4]]) labels=k means.predict(testData) labels Out[30]: array([0, 1, 0, 1, 0]) labels=labels.reshape(-1,1) labels Out[32]: array([[0], [0], [1], [0]]) usernames=np.array(['Paul','Kim','Liz','Tom','Bill']).reshape(-1,1) In [34]: usernames Out[34]: array([['Paul'], ['Kim'], ['Liz'], ['Tom'], ['Bill']], dtype='<U4') cols=movies.columns.to list()

Introduction

K-means Clustering

• Recompute the centroid of each cluster.

similar to each other than to instances that belong to other clusters.

into k disjoint clusters by iteratively applying the following steps

Form K clusters by assigning each instance to its nearest centroid

Cluster analysis aims at partitioning the input data into groups of closely related instances that belong to the same cluster and are more

The K-means clustering algorithm represents each cluster by its correspoinding cluster centroid. The algorithm partitions the input data

Out[36]: ['User', 'Jaws', 'Star Wars', 'Exorcist', 'Omen'] Out[38]: ['User', 'Jaws', 'Star Wars', 'Exorcist', 'Omen', 'Cluster ID'] In [40]: Out[40]:

cols

cols

newuser

Paul

Kim

Liz

Tom

Bill

SSE=[]

SSE

2

In [41]:

In [42]:

In [43]:

In [44]:

In [45]:

In [46]:

In [47]:

In [49]:

2

5

%matplotlib inline

for k in numCluster:

k_means.fit(data)

plt.plot(numCluster, SSE)

plt.ylabel('SSE')

Out[47]: Text(0, 0.5, 'SSE')

40

30

20

10

digits

plt.xlabel('Number of Clusters')

cols.append('Cluster ID')

User Jaws Star Wars Exorcist Omen Cluster ID

4

1

2

1

4

5

3

4

import matplotlib.pyplot as plt

numCluster = [1, 2, 3, 4, 5, 6]

k_means=cluster.KMeans(n_clusters=k)

Out[46]: [45.6666666666667, 9.33333333333334, 5.5, 2.5, 1.0, 0.0]

Number of Clusters

from sklearn.datasets import load_digits

Out[51]: {'data': array([[0., 0., 5., ..., 0., 0., 0.],

'target': array([0, 1, 2, ..., 8, 9, 8]),

'feature_names': ['pixel_0_0',

[0., 0., 0., ..., 10., 0., 0.], [0., 0., 0., ..., 16., 9., 0.],

[0., 0., 1., ..., 6., 0., 0.], [0., 0., 2., ..., 12., 0., 0.], [0., 0., 10., ..., 12., 1., 0.]]),

digits=load_digits()

'frame': None,

'pixel_0_1', 'pixel_0_2', 'pixel 0 3', 'pixel 0 4', 'pixel 0 5', 'pixel 0 6', 'pixel_0_7', 'pixel_1_0', 'pixel_1_1', 'pixel 1 2', 'pixel 1 3', 'pixel 1 4', 'pixel 1 5', 'pixel 1 6',

Loading [MathJax]/extensions/Safe.js

SSE.append(k_means.inertia_)

newuser=pd.DataFrame(np.concatenate((usernames,testData,labels), axis=1),columns=cols)

0

1

0

1 0

Elbew in the plow of SSE vs the number of clusters can be used to estimate the number of clusters.

To determine the number of clusters, we can vary the clusters 1-6 and then computer the corresponding sum of squared error. The

'pixel_1_7', 'pixel_2_0', 'pixel_2_1', 'pixel_2_2', 'pixel 2 3', 'pixel 2 4', 'pixel 2 5', 'pixel 2 6', 'pixel_2_7', 'pixel_3_0', 'pixel_3_1', 'pixel 3 2', 'pixel 3 3', 'pixel 3 4', 'pixel 3 5', 'pixel 3 6', 'pixel_3_7', 'pixel_4_0', 'pixel_4_1', 'pixel 4 2', 'pixel 4 3', 'pixel 4 4', 'pixel 4 5', 'pixel 4 6', 'pixel 4 7', 'pixel_5_0', 'pixel_5_1', 'pixel 5 2', 'pixel 5 3', 'pixel 5 4', 'pixel 5 5', 'pixel 5 6', 'pixel_5_7', 'pixel_6_0', 'pixel_6_1', 'pixel_6_2', 'pixel 6 3', 'pixel 6 4', 'pixel 6 5', 'pixel 6 6', 'pixel 6 7', 'pixel_7_0', 'pixel_7_1', 'pixel_7_2', 'pixel 'pixel_7_4', 'pixel 7 5', 'pixel 7 6', 'pixel_7_7'], 'target_names': array([0, 1, 2, 3, 4, 5, 6, 7, 8, 9]), 'images': array([[[0., 0., 5., ..., 1., 0., 0.], [0., 0., 13., ..., 15., 5., 0.], [0., 3., 15., ..., 11., 8., 0.], [0., 4., 11., ..., 12., 7., 0.], 2., 14., ..., 12., 0., 0.], [0., 0., 6., ..., 0., 0., 0.]], [[0., 0., 0., ..., 5., 0., [0., 0., 0., ..., 9., 0., 0.], 0., 3., ..., 6., [0., 0., 0.], [0., 0., 1., ..., 6., 0., 0.], [0., 0., 1., ..., 6., 0., 0.], [0., 0., 0., ..., 10., 0., 0.]], [[0., 0., 0., ..., 12., 0., 0.], [0., 0., 3., ..., 14., 0., 0.], [0., 0., 8., ..., 16., 0., 0.], $[0., 9., 16., \ldots, 0., 0., 0.],$ [0., 3., 13., ..., 11., 5., 0.], $[0., 0., 0., \ldots, 16., 9., 0.]],$. . . , [[0., 0., 1., ..., 1., 0., 0.], [0., 0., 13., ..., 2., 1., 0.], [0., 0., 16., ..., 16., 5., 0.],0., 16., ..., 15., 0., 0.], [0., [0., 0., 15., ..., 16., 0., 0.], [0., 0., 2., ..., 6., 0., 0.]], [[0., 0., 2., ..., 0., 0., [0., 0., 14., ..., 15., 1., [0., 4., 16., ..., 16., 7., [0., 0., 0., ..., 16., 2., 0.], 0.], [0., 0., 4., ..., 16., 2., [0., 0., 5., ..., 12., 0., 0.]], [[0., 0., 10., ..., 1., 0., 0.], [0., 2., 16., ..., 1., 0., 0.], [0., 0., 15., ..., 15., 0., 0.], [0., 4., 16., ..., 16., 6., 0.], [0., 8., 16., ..., 16., 8., 0.], [0., 1., 8., ..., 12., 1., 0.]]]), 'DESCR': ".. _digits_dataset:\n\nOptical recognition of handwritten digits dataset\n----------\n\n**Data Set Characteristics:**\n\n :Number of Instances: 5620\n :Number of Attributes: 64\n :Attribute Information: 8x8 image of integer pixels in the range 0..16.\n Attribute Values: None\n :Creator: E. Alpaydin (alpaydin '@' boun.edu.tr)\n :Date: July; 1998\n\nThis is a copy of the test set of the UCI ML hand-written digits datasets\nhttps://archive.ics.uci.edu/ml/dataset s/Optical+Recognition+of+Handwritten+Digits\n\nThe data set contains images of hand-written digits: 10 class es where\neach class refers to a digit.\n\nPreprocessing programs made available by NIST were used to extrac t\nnormalized bitmaps of handwritten digits from a preprinted form. From a\ntotal of 43 people, 30 contribut ed to the training set and different 13\nto the test set. 32x32 bitmaps are divided into nonoverlapping bloc ks of\n4x4 and the number of on pixels are counted in each block. This generates\nan input matrix of 8x8 whe re each element is an integer in the range n0..16. This reduces dimensionality and gives invariance to small \ndistortions.\n\nFor info on NIST preprocessing routines, see M. D. Garris, J. L. Blue, G.\nT. Candela, D. L. Dimmick, J. Geist, P. J. Grother, S. A. Janet, and C.\nL. Wilson, NIST Form-Based Handprint Recognition S ystem, NISTIR 5469,\n1994.\n\n.. topic:: References\n\n - C. Kaynak (1995) Methods of Combining Multiple Cl assifiers and Their\n Applications to Handwritten Digit Recognition, MSc Thesis, Institute of\n te Studies in Science and Engineering, Bogazici University.\n - E. Alpaydin, C. Kaynak (1998) Cascading Cla ssifiers, Kybernetika.\n - Ken Tang and Ponnuthurai N. Suganthan and Xi Yao and A. Kai Qin.\n ensionalityreduction using relevance weighted LDA. School of\n Electrical and Electronic Engineering Nany ang Technological University.\n 2005.\n - Claudio Gentile. A New Approximate Maximal Margin Classificati Algorithm. NIPS. 2000."}

