

Module 4

Lesson 1 - What Is k-Means Cluster Analysis?

Cluster Analysis

Goal: group similar observations together

The goal of cluster analysis is to group or cluster observations into subsets based

A photograph of a bustling indoor market. Shoppers are seen from behind, walking through narrow aisles filled with various goods. On the left, there are stacks of colorful fabrics and textiles. In the center, shelves are packed with small items like figurines or trinkets. To the right, there's a display of larger objects, possibly ceramics or traditional artifacts. The lighting is warm and somewhat dim, creating a cozy atmosphere.

Market Segmentation

**Less variance
within clusters**

**More variance
between clusters**

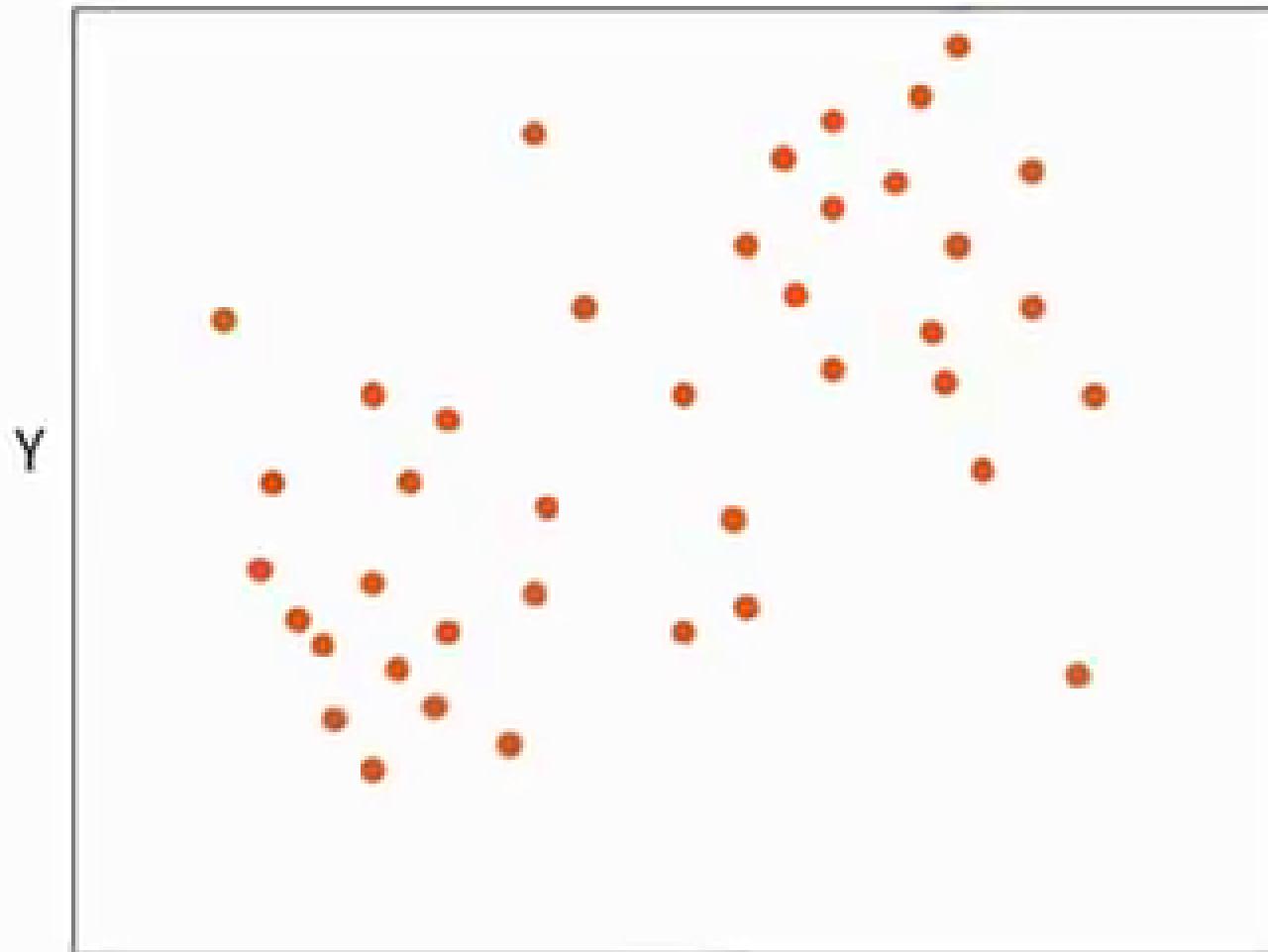
Data Reduction

- Reduce # of variables

X=2+ This categorical variable can then be used in other analysis

K-Means Cluster Analysis

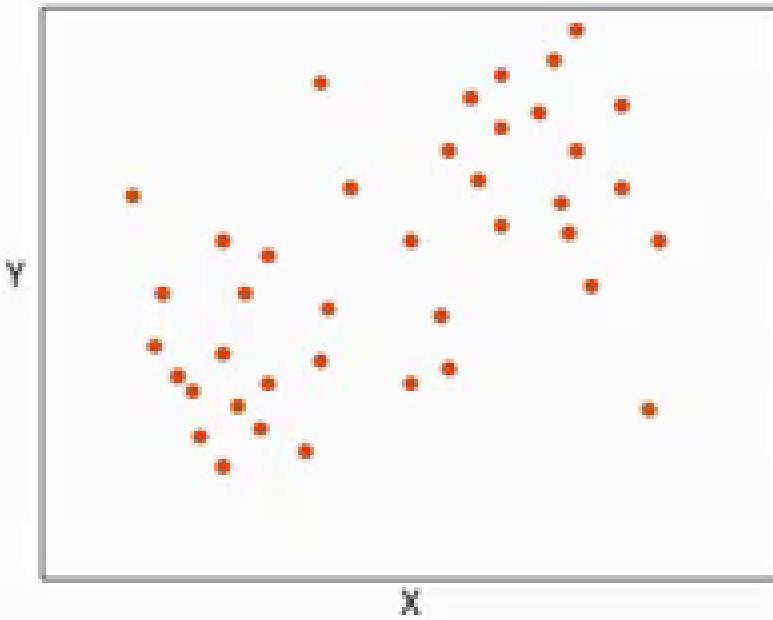
- Creates a p -dimensional space where p is the number of input variables
- Euclidean distance



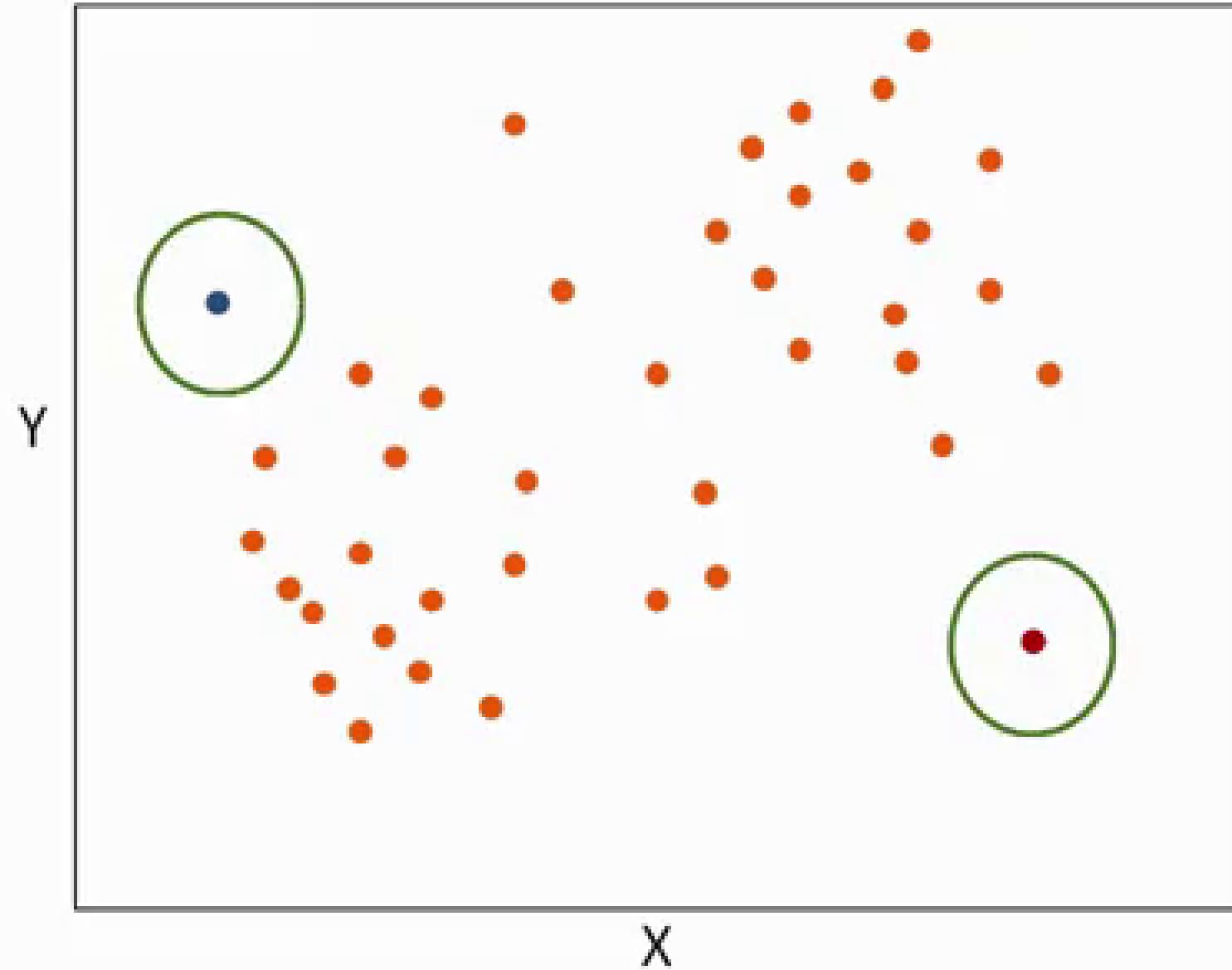
Scatterplot
 $p=2$ variables

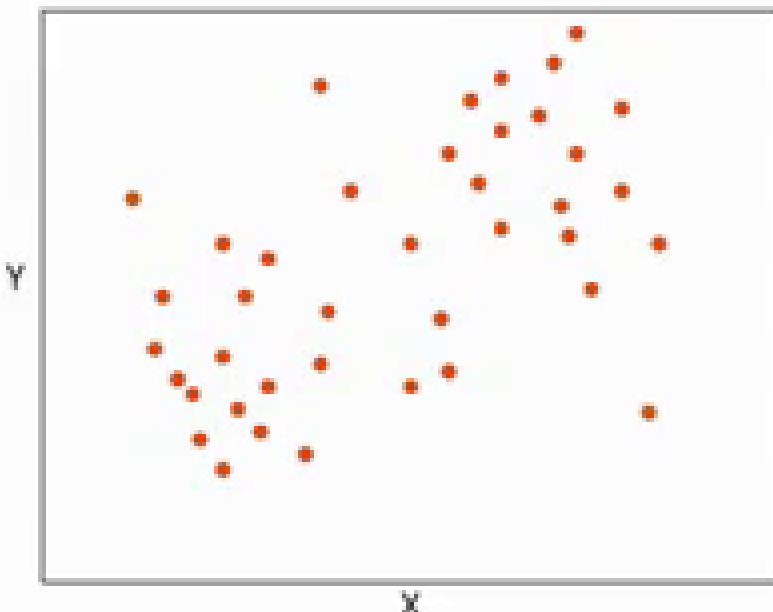
Each point is
an observation

where each of these points represents
an observation in the dataset

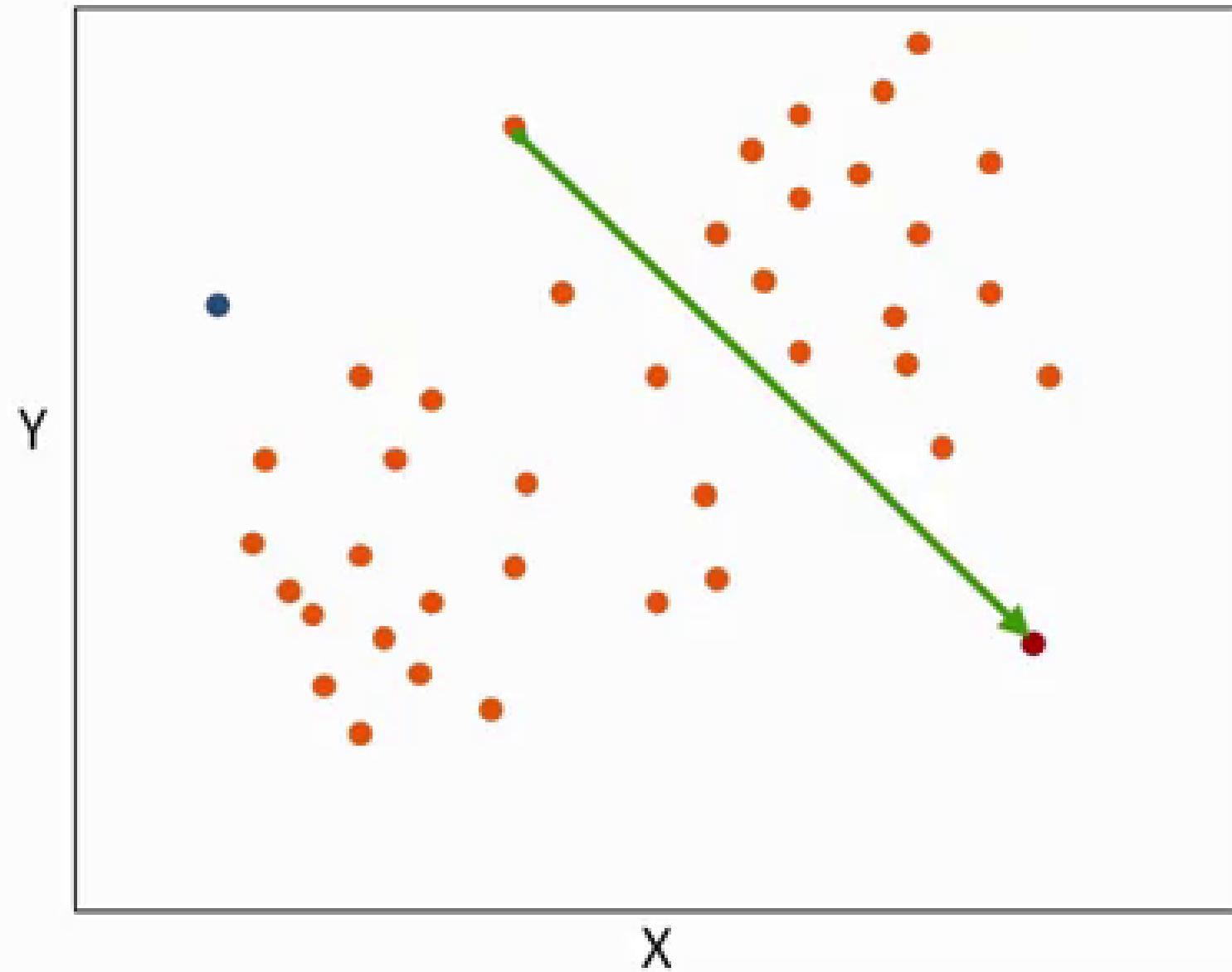


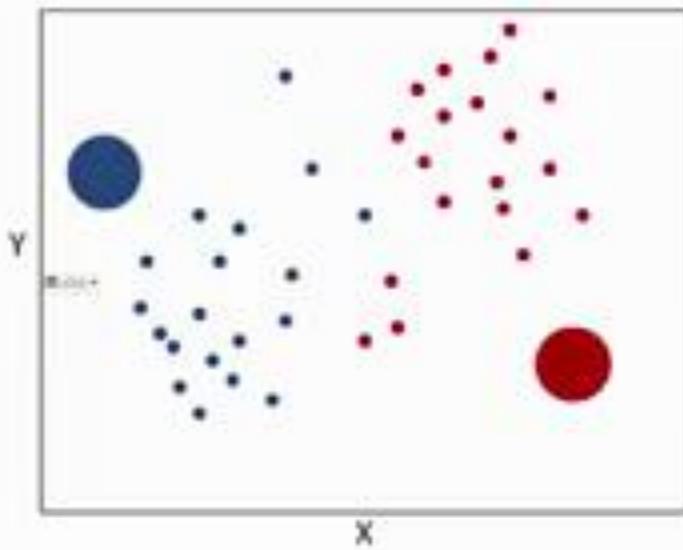
Step 1: randomly choose 2 points as the initial centroids (aka seeds) and calculate distance between points and the centroids



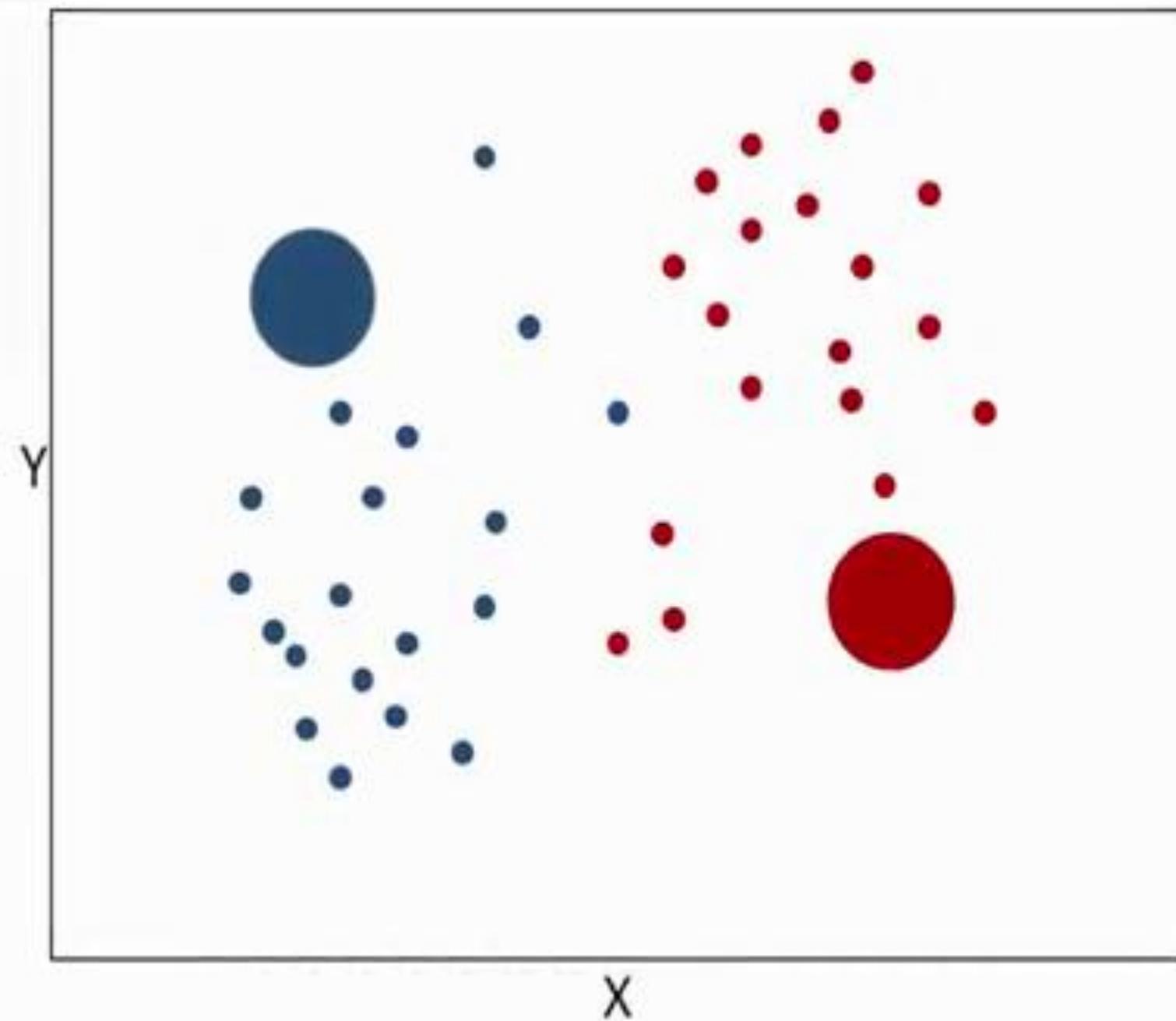


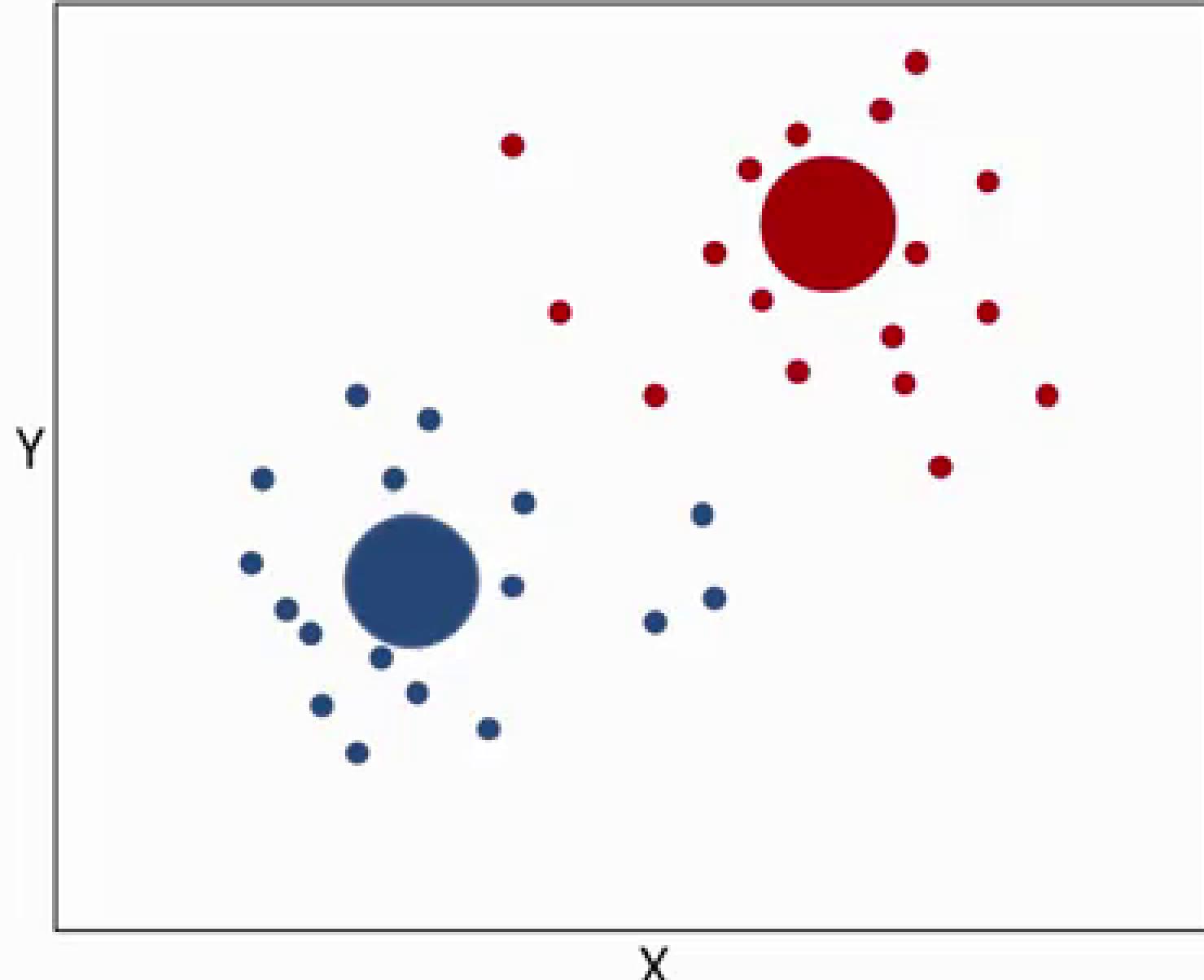
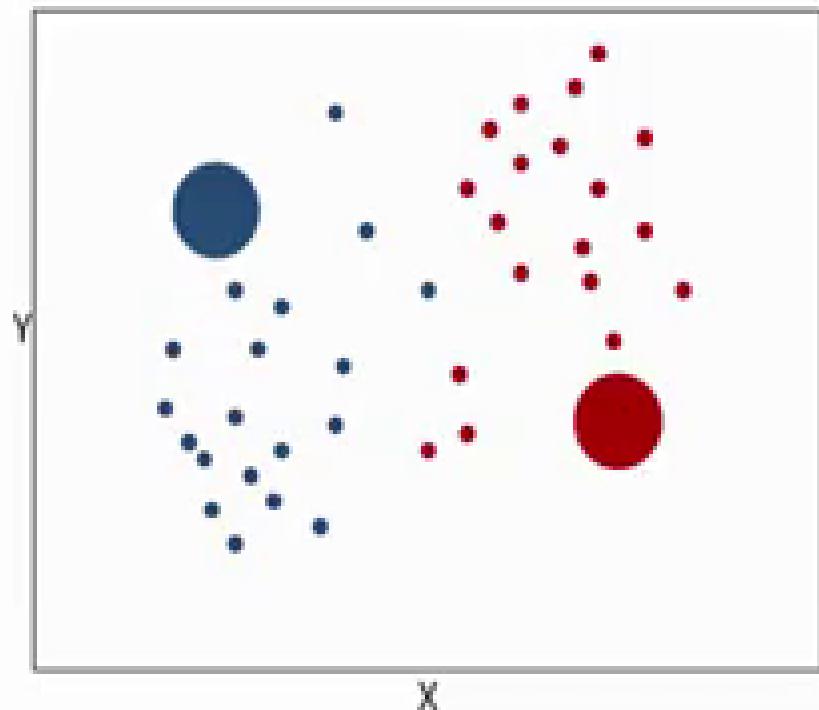
Step 1: randomly choose 2 points as the initial centroids (aka seeds) and calculate distance between points and the centroids



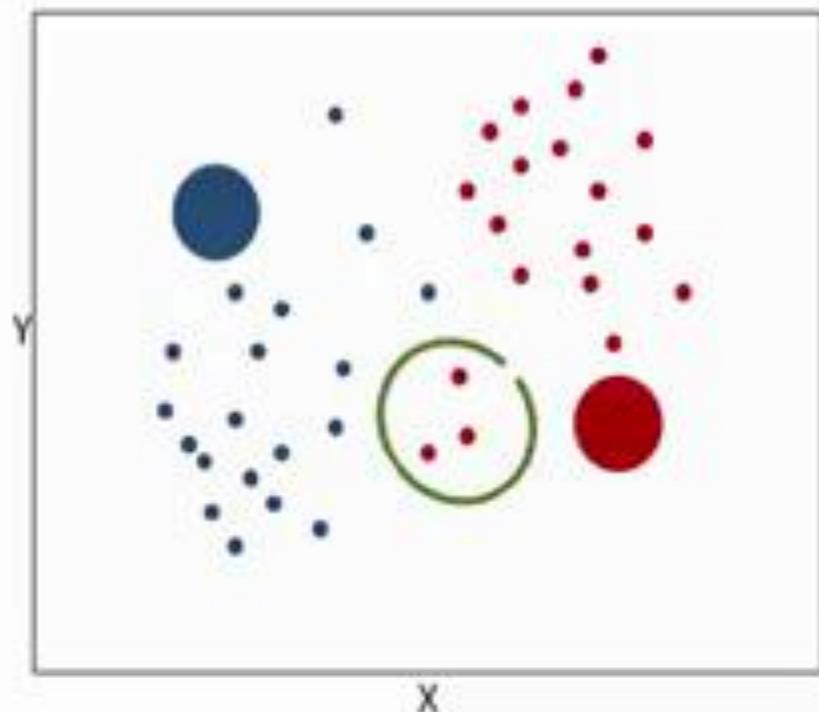


*Step 2: each centroid
is recalculated based
on the location of the
points that were
assigned to it*

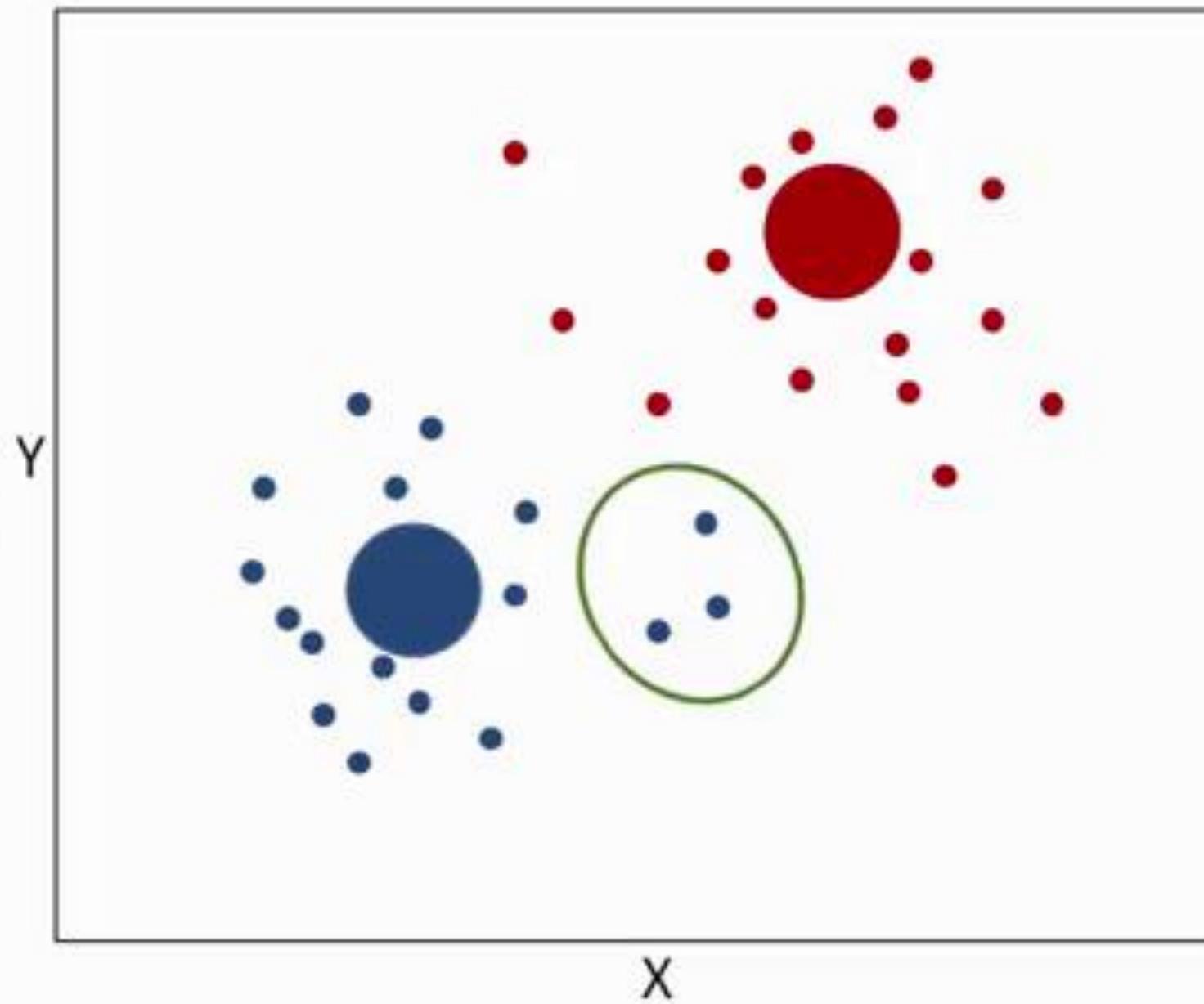


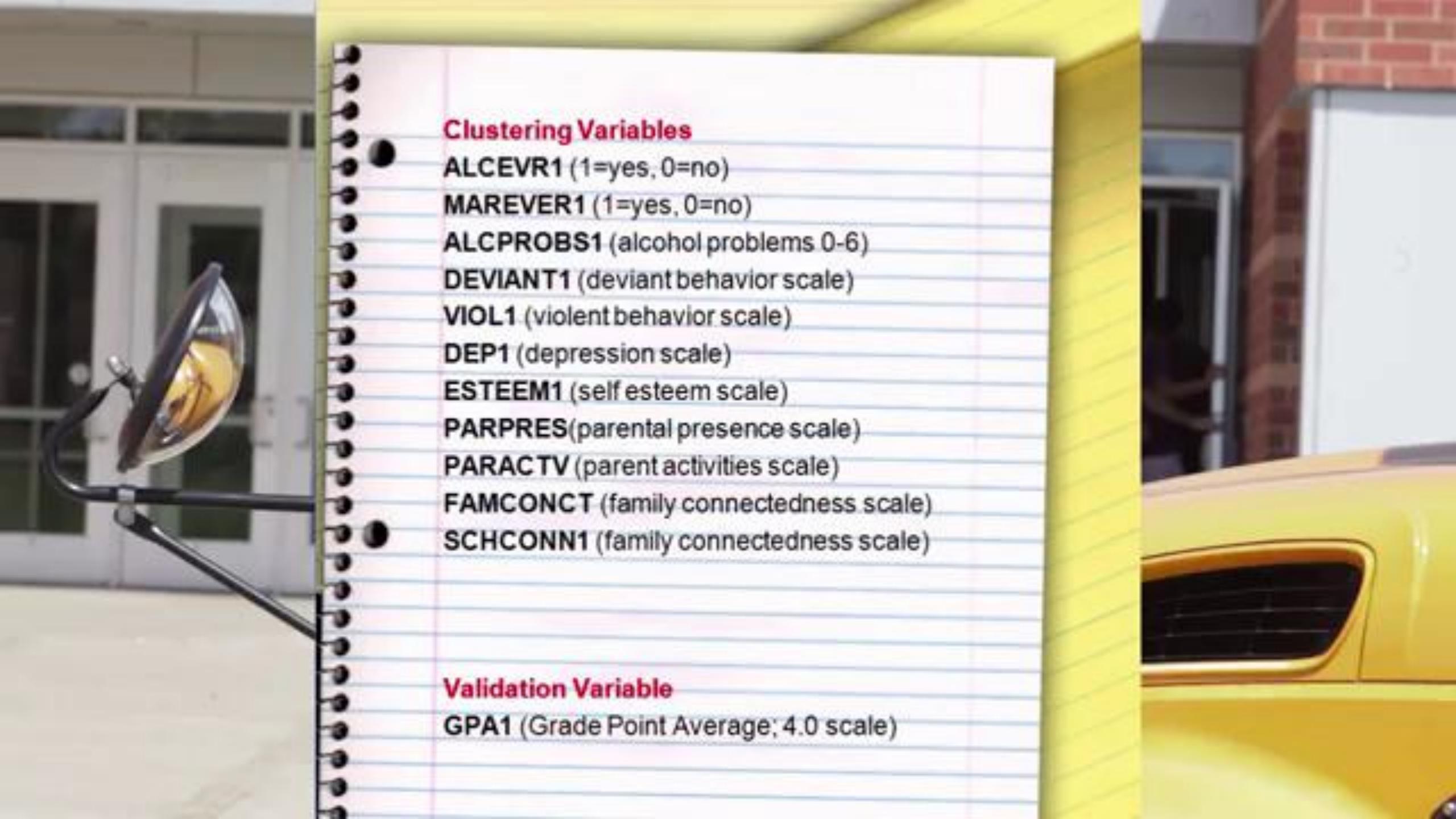


*Step 3: process is
repeated until the
location of the
centroids doesn't
change very much*



Step 3: process is repeated until the location of the centroids doesn't change very much





Clustering Variables

ALCEVR1 (1=yes, 0=no)
MAREVER1 (1=yes, 0=no)
ALCPROBS1 (alcohol problems 0-6)
DEVIANT1 (deviant behavior scale)
VIOL1 (violent behavior scale)
DEP1 (depression scale)
ESTEEM1 (self esteem scale)
PARPRES(parental presence scale)
PARACTV (parent activities scale)
FAMCONCT (family connectedness scale)
SCHCONN1 (family connectedness scale)

Validation Variable

GPA1 (Grade Point Average; 4.0 scale)



If we do find differences, then we have some evidence that the clusters are valid

Module 4

Lesson 2 - Running a k-Means Cluster Analysis in Python, pt. 1

```
1 # -*- coding: utf-8 -*-
2 """
3 Created on Mon Jan 18 19:51:29 2016
4
5 @author: jrose01
6 """
7
8 from pandas import Series, DataFrame
9 import pandas as pd
10 import numpy as np
11 import matplotlib.pyplot as plt
12 from sklearn.cross_validation import train_test_split
13 from sklearn import preprocessing
14 from sklearn.cluster import KMeans
15
16 """
17 Data Management
18 """
19 data = pd.read_csv("treeaddhealth.csv")
20
21 #data.reset_index(Level=0, inplace=True)
22
23 #upper-case all DataFrame column names
24 data.columns = map(str.upper, data.columns)
25
26 """
27 # Data Management
28
29 data_clean = data.dropna()
30 #recode1 = {1:1, 2:0}
```

First, we will call in
the libraries that we will need.

Source Console Object

Usage

Here you
either or

Help can
to an ob
Inspecto

Object inspector Variable explorer File

IPython console

Console I/A

64 bit (AMD64)

Type "copyright", "credit"

IPython 3.2.0 -- An enhan
Anaconda is brought to yo
Please check out: http://

```
18 """
19 data = pd.read_csv("treeaddhealth.csv")
20
21 #data.reset_index(Level=0, inplace=True)
22
23 #upper-case all DataFrame column names
24 data.columns = map(str.upper, data.columns)
25
26 # Data Management
27
28 data_clean = data.dropna()
29 #recode1 = {1:1, 2:0}
30 #data_clean['MALE']= data_clean['BIO SEX'].map(recode1)
31
32 # subset clustering variables
33 cluster=data_clean[['ALCEVR1','MAREVER1','ALCPROBS1','DEVIANT1','VIOL1',
34 'DEP1','ESTEEM1','SCHCONN1','PARACTV', 'PARPRES','FAMCONCT']]
35 cluster.describe()
36
37
38 # standardize clustering variables to have mean=0 and sd=1
39 clustervar=cluster.copy()
40 clustervar['ALCEVR1']=preprocessing.scale(clustervar['ALCEVR1'].astype('float64'))
41 clustervar['ALCPROBS1']=preprocessing.scale(clustervar['ALCPROBS1'].astype('float64'))
42 clustervar['MAREVER1']=preprocessing.scale(clustervar['MAREVER1'].astype('float64'))
```

Object inspect
IPython console
Cell Co
64 bit
Type "c"
IPython
Anaconda
Please
?
%quickref
help
object?
%guiref
These c
=>>> fro
>>> fro
>>> x,
>>> k,
>>> f,
>>> In [1]:

C4W4 python code for video examples.py

```
32 # subset clustering variables
33 cluster=data_clean[['ALCEVR1','MAREVER1','ALCPROBS1','DEVIANT1','VIOL1',
34 'DEP1','ESTEEM1','SCHCONN1','PARACTV', 'PARPRES','FAMCONCT']]
35 cluster.describe()
36
37
38 # standardize clustering variables to have mean=0 and sd=1
39 clustervar=cluster.copy()
40 clustervar['ALCEVR1']=preprocessing.scale(clustervar['ALCEVR1'].astype('float64'))
41 clustervar['ALCPROBS1']=preprocessing.scale(clustervar['ALCPROBS1'].astype('float64'))
42 clustervar['MAREVER1']=preprocessing.scale(clustervar['MAREVER1'].astype('float64'))
43 clustervar['DEP1']=preprocessing.scale(clustervar['DEP1'].astype('float64'))
44 clustervar['ESTEEM1']=preprocessing.scale(clustervar['ESTEEM1'].astype('float64'))
45 clustervar['VIOL1']=preprocessing.scale(clustervar['VIOL1'].astype('float64'))
46 clustervar['DEVIANT1']=preprocessing.scale(clustervar['DEVIANT1'].astype('float64'))
47 clustervar['FAMCONCT']=preprocessing.scale(clustervar['FAMCONCT'].astype('float64'))
48 clustervar['SCHCONN1']=preprocessing.scale(clustervar['SCHCONN1'].astype('float64'))
49 clustervar['PARACTV']=preprocessing.scale(clustervar['PARACTV'].astype('float64'))
50 clustervar['PARPRES']=preprocessing.scale(clustervar['PARPRES'].astype('float64'))
51
52 # split data into train and test sets
53 clus_train, clus_test = train_test_split(clustervar, test_size=.3, random_state=123)
54
55 # k-means clustering
56 from scipy.spatial.distance import cdist
57 clusters=range(1,10)
58 meandist=[]
59
```

We use the following code to standardize the clustering variables to have a mean of

Object inspector Variable explorer

Python console

Console 1/A

C4W4 python code for video examples.py

```
30 #data_clean['MALE']= data_clean['BIO_SEX'].map(recode1)
31
32 # subset clustering variables
33 cluster=data_clean[['ALCEVR1','MAREVER1','ALCPROBS1','DEVIANT1','VIOL1',
34 'DEP1','ESTEEM1','SCHCONN1','PARACTV', 'PARPRES','FAMCONCT']]
35 cluster.describe()
36
37
38 # standardize clustering variables to have mean=0 and sd=1
39 clustervar=cluster.copy()
40 clustervar['ALCEVR1']=preprocessing.scale(clustervar['ALCEVR1'].astype('float64'))
41 clustervar['ALCPROBS1']=preprocessing.scale(clustervar['ALCPROBS1'].astype('float64'))
42 clustervar['MAREVER1']=preprocessing.scale(clustervar['MAREVER1'].astype('float64'))
43 clustervar['DEP1']=preprocessing.scale(clustervar['DEP1'].astype('float64'))
44 clustervar['ESTEEM1']=preprocessing.scale(clustervar['ESTEEM1'].astype('float64'))
45 clustervar['VIOL1']=preprocessing.scale(clustervar['VIOL1'].astype('float64'))
46 clustervar['DEVIANT1']=preprocessing.scale(clustervar['DEVIANT1'].astype('float64'))
47 clustervar['FAMCONCT']=preprocessing.scale(clustervar['FAMCONCT'].astype('float64'))
48 clustervar['SCHCONN1']=preprocessing.scale(clustervar['SCHCONN1'].astype('float64'))
49 clustervar['PARACTV']=preprocessing.scale(clustervar['PARACTV'].astype('float64'))
50 clustervar['PARPRES']=preprocessing.scale(clustervar['PARPRES'].astype('float64'))
51
52 # split data into train and test sets
53 clus_train, clus_test = train_test_split(clustervar, test_size=0.25, random_state=123)
54
55 # k-means cluster analysis
56 from scipy.spatial.distance import cdist
57 clusters=range(1,10)
```

a test data set consisting of
the other 30% of the observations.

```
42 clustervar['MAREVER1']=preprocessing.scale(clustervar['MAREVER1'].astype('float64'))
43 clustervar['DEP1']=preprocessing.scale(clustervar['DEP1'].astype('float64'))
44 clustervar['ESTEEM1']=preprocessing.scale(clustervar['ESTEEM1'].astype('float64'))
45 clustervar['VIOL1']=preprocessing.scale(clustervar['VIOL1'].astype('float64'))
46 clustervar['DEVIANT1']=preprocessing.scale(clustervar['DEVIANT1'].astype('float64'))
47 clustervar['FAMCONCT']=preprocessing.scale(clustervar['FAMCONCT'].astype('float64'))
48 clustervar['SCHCONN1']=preprocessing.scale(clustervar['SCHCONN1'].astype('float64'))
49 clustervar['PARACTV']=preprocessing.scale(clustervar['PARACTV'].astype('float64'))
50 clustervar['PARPRES']=preprocessing.scale(clustervar['PARPRES'].astype('float64'))
51
52 # split data into train and test sets
53 clus_train, clus_test = train_test_split(clustervar, test_size=.3, random_state=123)
54
55 # k-means cluster analysis for 1-9 clusters
56 from scipy.spatial.distance import cdist
57 clusters=range(1,10)
58 meandist=[]
59
60 for k in clusters:
61     model=KMeans(n_clusters=k)
62     model.fit(clus_train)
63     clusassign=model.predict(clus_train)
64     meandist.append(sum([euclidean(clusassign[i], model.cluster_centers_[mean]), axis=1]) / clus_train.shape[0])
65
66 """
67 Plot average distance from observations from the cluster centroid
68
```

we don't know how many clusters
actually exist in the population for

Object Inspector
IPython console
Cell
...
cluster
4'))

```
42 clustervar['MAREVER1']=preprocessing.scale(clustervar['MAREVER1'].astype('float64'))
43 clustervar['DEP1']=preprocessing.scale(clustervar['DEP1'].astype('float64'))
44 clustervar['ESTEEM1']=preprocessing.scale(clustervar['ESTEEM1'].astype('float64'))
45 clustervar['VIOL1']=preprocessing.scale(clustervar['VIOL1'].astype('float64'))
46 clustervar['DEVIANT1']=preprocessing.scale(clustervar['DEVIANT1'].astype('float64'))
47 clustervar['FAMCONCT']=preprocessing.scale(clustervar['FAMCONCT'].astype('float64'))
48 clustervar['SCHCONN1']=preprocessing.scale(clustervar['SCHCONN1'].astype('float64'))
49 clustervar['PARACTV']=preprocessing.scale(clustervar['PARACTV'].astype('float64'))
50 clustervar['PARPRES']=preprocessing.scale(clustervar['PARPRES'].astype('float64'))
51
52 # split data into train and test sets
53 clus_train, clus_test = train_test_split(clustervar, test_size=.3, random_state=123)
54
55 # k-means cluster analysis for 1-9 clusters
56 from scipy.spatial.distance import cdist
57 clusters=range(1,10)
58 meandist=[]
59
60 for k in clusters:
61     model=KMeans(n_clusters=k)
62     model.fit(clus_train)
63     clusassign=model.predict(clus_train)
64     meandist.append(np.mean(cdist(clusassign, model.cluster_centers_), axis=1))
65     / clus_train.shape[0])
66
67 """
68 Plot average distance from observations from the cluster centroid
    a range of values on the number of
    clusters before we begin we'll import
```

```
42 clustervar['MAREVER1']=preprocessing.scale(clustervar['MAREVER1'].astype('float64'))
43 clustervar['DEP1']=preprocessing.scale(clustervar['DEP1'].astype('float64'))
44 clustervar['ESTEEM1']=preprocessing.scale(clustervar['ESTEEM1'].astype('float64'))
45 clustervar['VIOL1']=preprocessing.scale(clustervar['VIOL1'].astype('float64'))
46 clustervar['DEVIANT1']=preprocessing.scale(clustervar['DEVIANT1'].astype('float64'))
47 clustervar['FAMCONCT']=preprocessing.scale(clustervar['FAMCONCT'].astype('float64'))
48 clustervar['SCHCONN1']=preprocessing.scale(clustervar['SCHCONN1'].astype('float64'))
49 clustervar['PARACTV']=preprocessing.scale(clustervar['PARACTV'].astype('float64'))
50 clustervar['PARPRES']=preprocessing.scale(clustervar['PARPRES'].astype('float64'))
51
52 # split data into train and test sets
53 clus_train, clus_test = train_test_split(clustervar, test_size=.3, random_state=123)
54
55 # k-means cluster analysis for 1-9 clusters
56 from scipy.spatial.distance import cdist
57 clusters=range(1,10)
58 meandist=[]
59
60 for k in clusters:
61     model=KMeans(n_clusters=k)
62     model.fit(clus_train)
63     clusassign=model.predict(clus_train)
64     meandist.append(sum(np.min(cdist(clus_train, model.cluster_centers_, 'euclidean'), axis=1)) / clus_train.shape[0])
65
66 """
67 Plot average distance from observations from the cluster centroid
68
```

```
42 clustervar['MAREVER1']=preprocessing.scale(clustervar['MAREVER1'].astype('float64'))
43 clustervar['DEP1']=preprocessing.scale(clustervar['DEP1'].astype('float64'))
44 clustervar['ESTEEM1']=preprocessing.scale(clustervar['ESTEEM1'].astype('float64'))
45 clustervar['VIOL1']=preprocessing.scale(clustervar['VIOL1'].astype('float64'))
46 clustervar['DEVIANT1']=preprocessing.scale(clustervar['DEVIANT1'].astype('float64'))
47 clustervar['FAMCONCT']=preprocessing.scale(clustervar['FAMCONCT'].astype('float64'))
48 clustervar['SCHCONN1']=preprocessing.scale(clustervar['SCHCONN1'].astype('float64'))
49 clustervar['PARACTV']=preprocessing.scale(clustervar['PARACTV'].astype('float64'))
50 clustervar['PARPRES']=preprocessing.scale(clustervar['PARPRES'].astype('float64'))
51
52 # split data into train and test sets
53 clus_train, clus_test = train_test_split(clustervar, test_size=.3, random_state=123)
54
55 # k-means cluster analysis for 1-9 clusters
56 from scipy.spatial.distance import cdist
57 clusters=range(1,10)
58 meandist=[]
59
60 for k in clusters:
61     model=KMeans(n_clusters=k)
62     model.fit(clus_train)
63     clusassign=model.predict(clus_train)
64     meandist.append(np.mean(cdist(clus_train, model.cluster_centers_, axis=1)) / clus_train.shape[0])
65
66 """
67 Plot average distance from observations from the cluster centroid
68
```

which will give us the cluster solutions, for k equals 1 to k equals 9 clusters.

Object Inspector
IPython console
Cell
...
cluster
4'))

C4W4 python code for video examples.py

```
42 clustervar['MAREVER1']=preprocessing.scale(clustervar['MAREVER1'].astype('float64'))
43 clustervar['DEP1']=preprocessing.scale(clustervar['DEP1'].astype('float64'))
44 clustervar['ESTEEM1']=preprocessing.scale(clustervar['ESTEEM1'].astype('float64'))
45 clustervar['VIOL1']=preprocessing.scale(clustervar['VIOL1'].astype('float64'))
46 clustervar['DEVIANT1']=preprocessing.scale(clustervar['DEVIANT1'].astype('float64'))
47 clustervar['FAMCONCT']=preprocessing.scale(clustervar['FAMCONCT'].astype('float64'))
48 clustervar['SCHCONN1']=preprocessing.scale(clustervar['SCHCONN1'].astype('float64'))
49 clustervar['PARACTV']=preprocessing.scale(clustervar['PARACTV'].astype('float64'))
50 clustervar['PARPRES']=preprocessing.scale(clustervar['PARPRES'].astype('float64'))
51
52 # split data into train and test sets
53 clus_train, clus_test = train_test_split(clustervar, test_size=.3, random_state=123)
54
55 # k-means cluster analysis for 1-9 clusters
56 from scipy.spatial.distance import cdist
57 clusters=range(1,10)
58 meandist=[]
59
60 for k in clusters:
61     model=KMeans(n_clusters=k)
62     model.fit(clus_train)
63     clusassign=model.predict(clus_train)
64     meandist.append(np.mean(cdist(clusassign.reshape(-1,1), clus_train, axis=1)))
65 / clus_train.shape[0]
66
67 ....
68 Plot average distance from observations from the cluster centroid
```

The formula first calculates the distance between each observation and the cluster

C4W4 python code for video examples.py

```
42 clustervar['MAREVER1']=preprocessing.scale(clustervar['MAREVER1'].astype('float64'))
43 clustervar['DEP1']=preprocessing.scale(clustervar['DEP1'].astype('float64'))
44 clustervar['ESTEEM1']=preprocessing.scale(clustervar['ESTEEM1'].astype('float64'))
45 clustervar['VIOL1']=preprocessing.scale(clustervar['VIOL1'].astype('float64'))
46 clustervar['DEVIANT1']=preprocessing.scale(clustervar['DEVIANT1'].astype('float64'))
47 clustervar['FAMCONCT']=preprocessing.scale(clustervar['FAMCONCT'].astype('float64'))
48 clustervar['SCHCONN1']=preprocessing.scale(clustervar['SCHCONN1'].astype('float64'))
49 clustervar['PARACTV']=preprocessing.scale(clustervar['PARACTV'].astype('float64'))
50 clustervar['PARPRES']=preprocessing.scale(clustervar['PARPRES'].astype('float64'))
51
52 # split data into train and test sets
53 clus_train, clus_test = train_test_split(clustervar, test_size=.3, random_state=123)
54
55 # k-means cluster analysis for 1-9 clusters
56 from scipy.spatial.distance import cdist
57 clusters=range(1,10)
58 meandist=[]
59
60 for k in clusters:
61     model=KMeans(n_clusters=k)
62     model.fit(clus_train)
63     clusassign=model.predict(clus_train)
64     meandist.append(np.mean(cdist(clus_train, model.cluster_centers_, 'euclidean'), axis=1))
65 / clus_train.shape[0])
66
67 ....
68 Plot average distance from observations from the cluster centroid
```

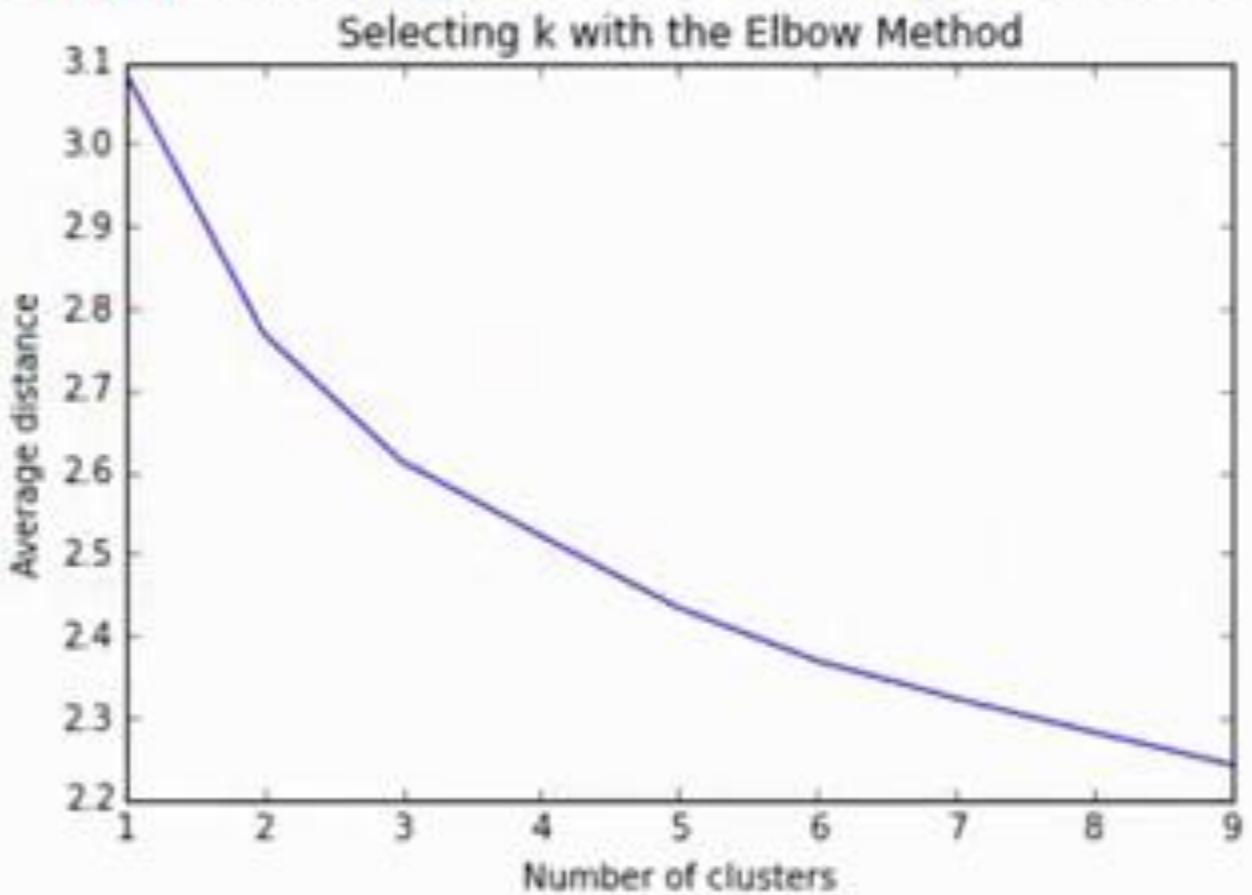
each centroid taking the smallest distance as the value of the minimum,

```
42 clustervar['MAREVER1']=preprocessing.scale(clustervar['MAREVER1'].astype('float64'))
43 clustervar['DEP1']=preprocessing.scale(clustervar['DEP1'].astype('float64'))
44 clustervar['ESTEEM1']=preprocessing.scale(clustervar['ESTEEM1'].astype('float64'))
45 clustervar['VIOL1']=preprocessing.scale(clustervar['VIOL1'].astype('float64'))
46 clustervar['DEVIANT1']=preprocessing.scale(clustervar['DEVIANT1'].astype('float64'))
47 clustervar['FAMCONCT']=preprocessing.scale(clustervar['FAMCONCT'].astype('float64'))
48 clustervar['SCHCONN1']=preprocessing.scale(clustervar['SCHCONN1'].astype('float64'))
49 clustervar['PARACTV']=preprocessing.scale(clustervar['PARACTV'].astype('float64'))
50 clustervar['PARPRES']=preprocessing.scale(clustervar['PARPRES'].astype('float64'))
51
52 # split data into train and test sets
53 clus_train, clus_test = train_test_split(clustervar, test_size=.3, random_state=123)
54
55 # k-means cluster analysis for 1-9 clusters
56 from scipy.spatial.distance import cdist
57 clusters=range(1,10)
58 meandist=[]
59
60 for k in clusters:
61     model=KMeans(n_clusters=k)
62     model.fit(clus_train)
63     clusassign=model.predict(clus_train)
64     meandist.append((model.cluster_centers_, axis=1))
65     / clus_train.shape[0])
66 """
67
68 Plot average distance from observations from the cluster centroid
```

by the number of observations in the
clus_train data set where the .shape with

```
.... plt.ylabel('Average distance')
.... plt.title('Selecting k with the Elbow Method')
....
```

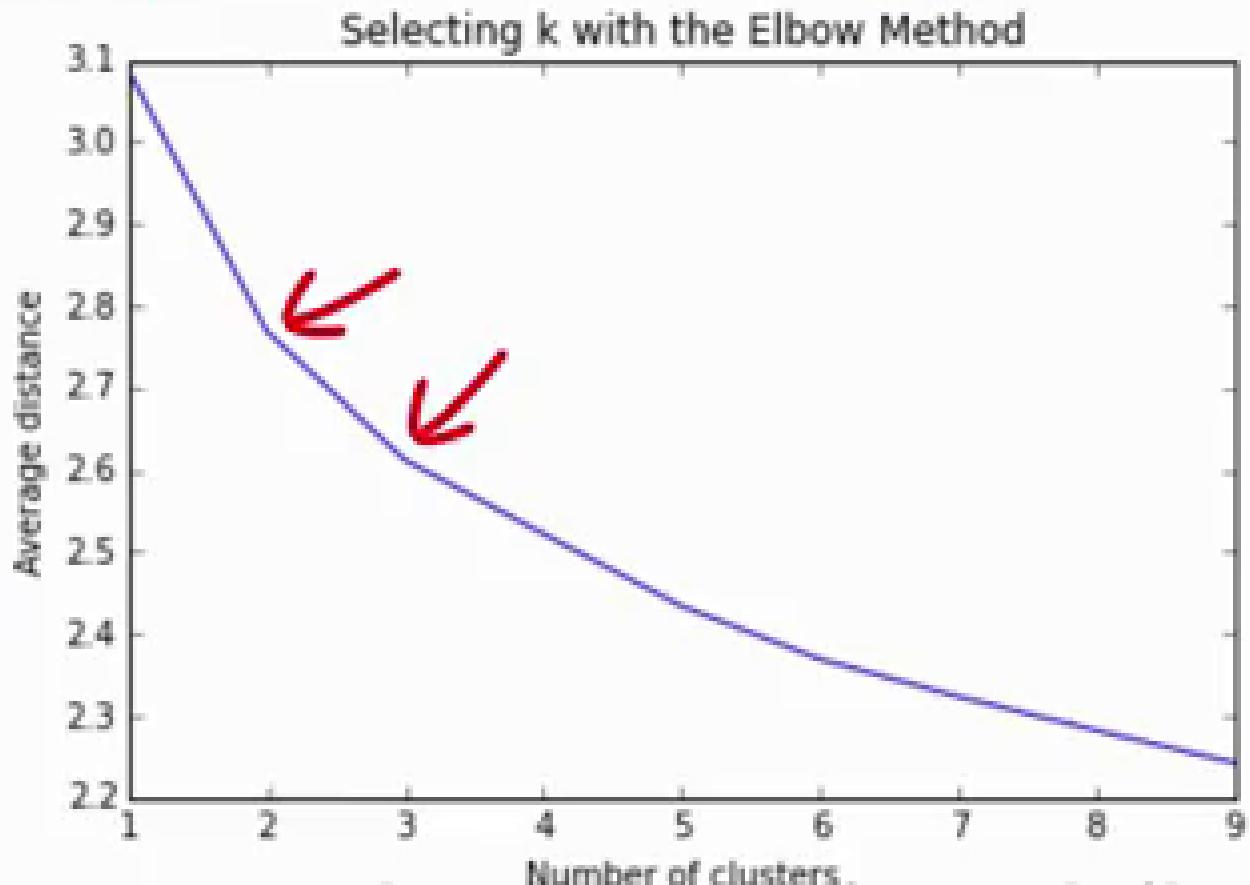
Out[5]: <matplotlib.text.Text at 0xd761588>



In [6]:

```
...: plt.ylabel('Average distance')
...: plt.title('Selecting k with the Elbow Method')
...:
```

Out[5]: <matplotlib.text.Text at 0xd761588>

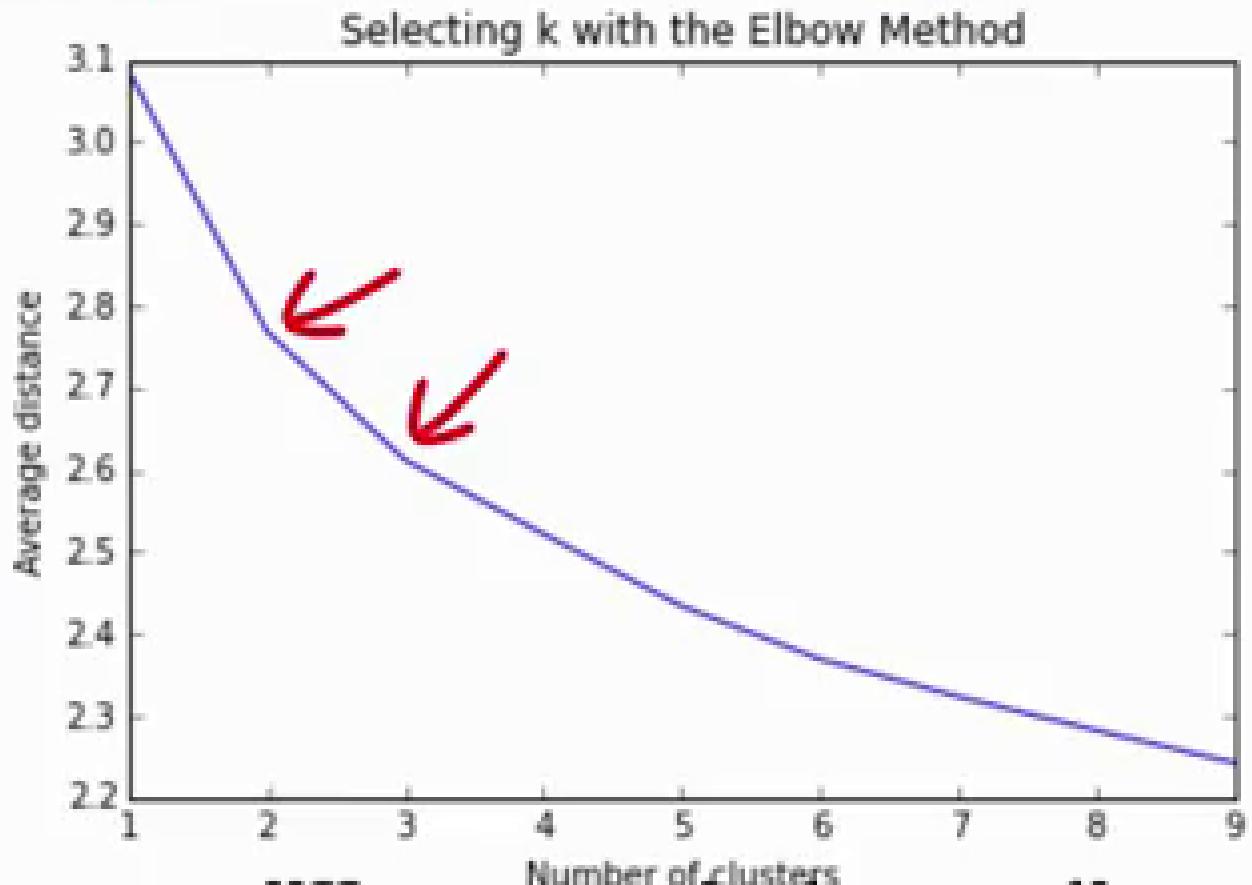


examine the cluster solutions for at least

In [6] the two and three cluster solutions to see

```
...: plt.ylabel('Average distance')
...: plt.title('Selecting k with the Elbow Method')
...:
```

Out[5]: <matplotlib.text.Text at 0xd761588>

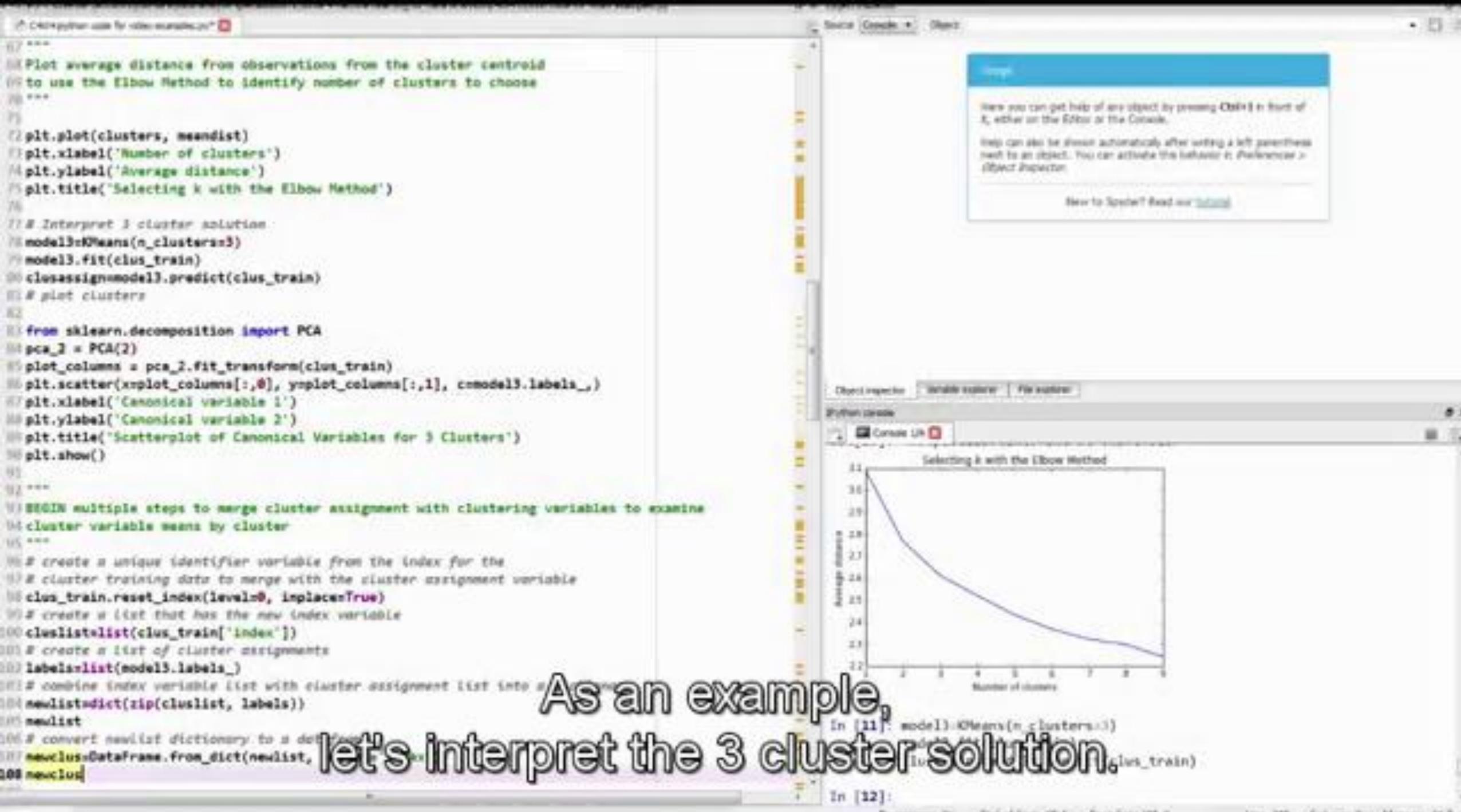


differences between the clusters on

In [6]: our external validation variable GPA.

Module 4

Lesson 3 - Running a k-Means Cluster Analysis in Python, pt. 2



```
67 """
68 Plot average distance from observations from the cluster centroid
69 to use the Elbow Method to identify number of clusters to choose
70 """
71
72 plt.plot(clusters, meandist)
73 plt.xlabel('Number of clusters')
74 plt.ylabel('Average distance')
75 plt.title('Selecting k with the Elbow Method')
76
77 # Interpret 3 cluster solution
78 model3=KMeans(n_clusters=3)
79 model3.fit(clus_train)
80 clusassign=model3.predict(clus_train)
81
82
83 from sklearn.decomposition import PCA
84 pca_2 = PCA(2)
85 plot_columns = pca_2.fit_transform(clus_train)
86 plt.scatter(x=plot_columns[:,0], y=plot_columns[:,1], c=model3.labels_)
87 plt.xlabel('Canonical variable 1')
88 plt.ylabel('Canonical variable 2')
89 plt.title('Scatter plot of Canonical Variables for 3 Clusters')
90 plt.show()
91
92 """
93 BEGIN multiple steps to merge cluster assignment with clustering variables to examine
```

So we create an object, 'model 3', which will contain the results from the cluster

```
67 """
68 Plot average distance from observations from the cluster centroid
69 to use the Elbow Method to identify number of clusters to choose
70 """
71
72 plt.plot(clusters, meandist)
73 plt.xlabel('Number of clusters')
74 plt.ylabel('Average distance')
75 plt.title('Selecting k with the Elbow Method')
76
77 # Interpret 3 cluster solution
78 model3=KMeans(n_clusters=3)
79 model3.fit(clus_train)
80 clusassign=model3.predict(clus_train)
81 # plot clusters
82
83 from sklearn.decomposition import PCA
84 pca_2 = PCA(2)
85 plot_columns = pca_2.fit_transform(clus_train)
86 plt.scatter(x=plot_columns[:,0], y=plot_columns[:,1], c=model3.labels_)
87 plt.xlabel('Canonical variable 1')
88 plt.ylabel('Canonical variable 2')
89 plt.title('Scatterplot of the first two canonical variables')
90 plt.show()
91 """
92 """
93 BEGIN multiple steps to merge cluster assignment with clustering variables to examine
```

analysis with 3 clusters =KMeans,
and in parenthesis, n_clusters=3.

```
67 """
68 Plot average distance from observations from the cluster centroid
69 to use the Elbow Method to identify number of clusters to choose
70 """
71
72 plt.plot(clusters, meandist)
73 plt.xlabel('Number of clusters')
74 plt.ylabel('Average distance')
75 plt.title('Selecting k with the Elbow Method')
76
77 # Interpret 3 cluster solution
78 model3=KMeans(n_clusters=3)
79 model3.fit(clus_train)
80 clusassign=model3.predict(clus_train)
81 # plot clusters
82
83 from sklearn.decomposition import PCA
84 pca_2 = PCA(2)
85 plot_columns = pca_2.fit_transform(clus_train)
86 plt.scatter(x=plot_columns[:,0], y=plot_columns[:,1], c=model3.labels_)
87 plt.xlabel('Canonical variable 1')
88 plt.ylabel('Canonical variable 2')
89 plt.title('Scatterplot of the first two canonical variables')
90 plt.show()
91
92 """
93 BEGIN multiple steps to merge cluster assignment with clustering variables to examine
```

The first thing we want to try is to graph the clusters in a scatter plot to

```
67 """
68 Plot average distance from observations from the cluster centroid
69 to use the Elbow Method to identify number of clusters to choose
70 """
71
72 plt.plot(clusters, meandist)
73 plt.xlabel('Number of clusters')
74 plt.ylabel('Average distance')
75 plt.title('Selecting k with the Elbow Method')
76
77 # Interpret 3 cluster solution
78 model3=KMeans(n_clusters=3)
79 model3.fit(clus_train)
80 clusassign=model3.predict(clus_train)
81 # plot clusters
82
83 from sklearn.decomposition import PCA
84 pca_2 = PCA(2)
85 plot_columns = pca_2.fit_transform(clus_train)
86 plt.scatter(x=plot_columns[:,0], y=plot_columns[:,1], c=model3.labels_,)
87 plt.xlabel('Canonical variable 1')
88 plt.ylabel('Canonical variable 2')
89 plt.title('Scatterplot of Canonical variables for 3 Clusters')
90 plt.show()
91 """
92 """
93 BEGIN multiple steps to merge cluster assignment with clustering variables to examine
```

see whether or not they overlap with each other in terms of their location in

```
67 """
68 Plot average distance from observations from the cluster centroid
69 to use the Elbow Method to identify number of clusters to choose
70 """
71
72 plt.plot(clusters, meandist)
73 plt.xlabel('Number of clusters')
74 plt.ylabel('Average distance')
75 plt.title('Selecting k with the Elbow Method')
76
77 # Interpret 3 cluster solution
78 model3=KMeans(n_clusters=3)
79 model3.fit(clus_train)
80 clusassign=model3.predict(clus_train)
81 # plot clusters
82
83 from sklearn.decomposition import PCA
84 pca_2 = PCA(2)
85 plot_columns = pca_2.fit_transform(clus_train)
86 plt.scatter(x=plot_columns[:,0], y=plot_columns[:,1], c=model3.labels_,)
87 plt.xlabel('Canonical variable 1')
88 plt.ylabel('Canonical variable 2')
89 plt.title('Scatterplot of first two canonical variables')
90 plt.show()
91 """
92 """
93 BEGIN multiple steps to merge cluster assignment with clustering variables to examine
```

However, with 11 clustering variables
that means we have 11 dimensions,

```
67 """
68 Plot average distance from observations from the cluster centroid
69 to use the Elbow Method to identify number of clusters to choose
70 """
71
72 plt.plot(clusters, meandist)
73 plt.xlabel('Number of clusters')
74 plt.ylabel('Average distance')
75 plt.title('Selecting k with the Elbow Method')
76
77 # Interpret 3 cluster solution
78 model3=KMeans(n_clusters=3)
79 model3.fit(clus_train)
80 clusassign=model3.predict(clus_train)
81 # plot clusters
82
83 from sklearn.decomposition import PCA
84 pca_2 = PCA(2)
85 plot_columns = pca_2.fit_transform(clus_train)
86 plt.scatter(x=plot_columns[:,0], y=plot_columns[:,1], c=model3.labels_)
87 plt.xlabel('Canonical variable 1')
88 plt.ylabel('Canonical variable 2')
89 plt.title('Scatterplot of Canonical Variables for 3 Clusters')
90 plt.show()
91 """
92 """
93 BEGIN multiple steps to merge cluster assignment with clustering variables to examine
```

which would be impossible to visualize.

```
67 """
68 Plot average distance from observations from the cluster centroid
69 to use the Elbow Method to identify number of clusters to choose
70 """
71
72 plt.plot(clusters, meandist)
73 plt.xlabel('Number of clusters')
74 plt.ylabel('Average distance')
75 plt.title('Selecting k with the Elbow Method')
76
77 # Interpret 3 cluster solution
78 model3=KMeans(n_clusters=3)
79 model3.fit(clus_train)
80 clusassign=model3.predict(clus_train)
81 # plot clusters
82
83 from sklearn.decomposition import PCA
84 pca_2 = PCA(2)
85 plot_columns = pca_2.fit_transform(clus_train)
86 plt.scatter(x=plot_columns[:,0], y=plot_columns[:,1], c=model3.labels_)
87 plt.xlabel('Canonical variable 1')
88 plt.ylabel('Canonical variable 2')
89 plt.title('Scatterplot')
90 plt.show()
91 """
92 """
93 BEGIN multiple steps to merge cluster assignment with clustering variables to examine
```

A scatter plot will work to visualize a few dimensions, but not 11 dimensions.

```
67 """
68 Plot average distance from observations from the cluster centroid
69 to use the Elbow Method to identify number of :
70 """
71
72 plt.plot(clusters, meandist)
73 plt.xlabel('Number of clusters')
74 plt.ylabel('Average distance')
75 plt.title('Selecting k with the Elbow Method')
76
77 # Interpret 3 cluster solution
78 model3=KMeans(n_clusters=3)
79 model3.fit(clus_train)
80 clusassign=model3.predict(clus_train)
81 # plot clusters
82
83 from sklearn.decomposition import PCA
84 pca_2 = PCA(2)
85 plot_columns = pca_2.fit_transform(clus_train)
86 plt.scatter(x=plot_columns[:,0], y=plot_columns[:,1])
87 plt.xlabel('Canonical variable 1')
88 plt.ylabel('Canonical variable 2')
89 plt.title('Scatter plot of canonical variables')
90 plt.show()
91 """
92 """
93 BEGIN multiple steps to merge cluster assignment with clustering variables to examine
```

Canonical Discriminant Analysis

- Creates a smaller number of variables
- Linear combinations of clustering variables
- Canonical variables are ordered by proportion of variance accounted for

called canonical variables, are ordered in terms of the proportion of variance and

```
67 """
68 Plot average distance from observations from the cluster centroid
69 to use the Elbow Method to identify number of :
70 """
71
72 plt.plot(clusters, meandist)
73 plt.xlabel('Number of clusters')
74 plt.ylabel('Average distance')
75 plt.title('Selecting k with the Elbow Method')
76
77 # Interpret 3 cluster solution
78 model3=KMeans(n_clusters=3)
79 model3.fit(clus_train)
80 clusassign=model3.predict(clus_train)
81 # plot clusters
82
83 from sklearn.decomposition import PCA
84 pca_2 = PCA(2)
85 plot_columns = pca_2.fit_transform(clus_train)
86 plt.scatter(x=plot_columns[:,0], y=plot_columns[:,1])
87 plt.xlabel('Canonical variable 1')
88 plt.ylabel('Canonical variable 2')
89 plt.title('Scatterplot of Canonical Variables')
90 plt.show()
91
92 """
93 BEGIN multiple steps to merge cluster assignment with clustering variables to examine
```

Canonical Discriminant Analysis

- Creates a smaller number of variables
- Linear combinations of clustering variables
- Canonical variables are ordered by proportion of variance accounted for

So the first canonical variable will count for

```
67 """
68 Plot average distance from observations from the cluster centroid
69 to use the Elbow Method to identify number of :
70 """
71
72 plt.plot(clusters, meandist)
73 plt.xlabel('Number of clusters')
74 plt.ylabel('Average distance')
75 plt.title('Selecting k with the Elbow Method')
76
77 # Interpret 3 cluster solution
78 model3=KMeans(n_clusters=3)
79 model3.fit(clus_train)
80 clusassign=model3.predict(clus_train)
81 # plot clusters
82
83 from sklearn.decomposition import PCA
84 pca_2 = PCA(2)
85 plot_columns = pca_2.fit_transform(clus_train)
86 plt.scatter(x=plot_columns[:,0], y=plot_columns[:,1])
87 plt.xlabel('Canonical variable 1')
88 plt.ylabel('Canonical variable 2')
89 plt.title('Scatterplot of Canonical Variables')
90 plt.show()
91
92 """
93 BEGIN multiple steps to merge cluster assignment with clustering variables to examine
```

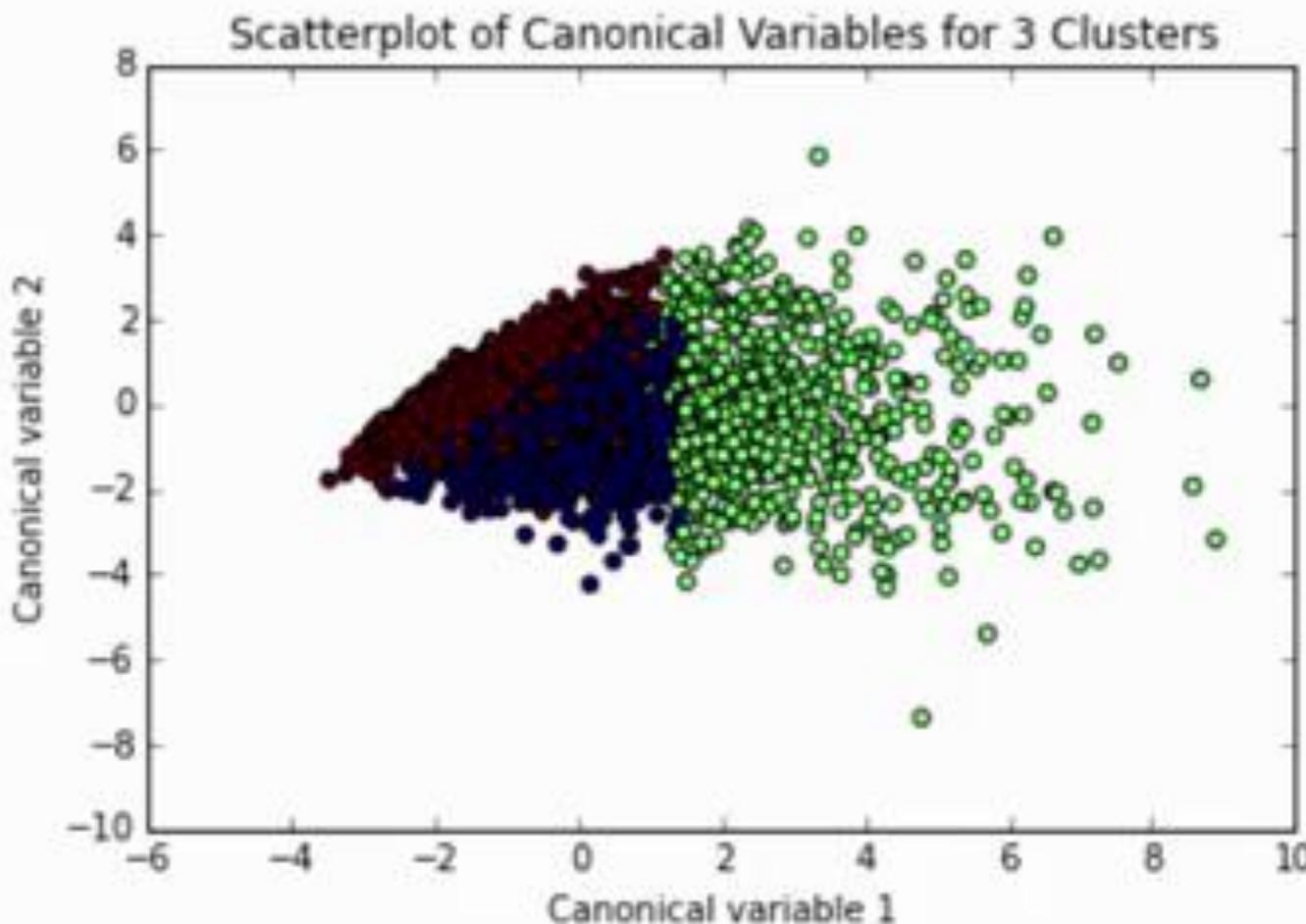
Canonical Discriminant Analysis

- Creates a smaller number of variables
- Linear combinations of clustering variables
- Canonical variables are ordered by proportion of variance accounted for
- Majority of variance is accounted for by first few canonical variables

```
67 """
68 Plot average distance from observations from the cluster centroid
69 to use the Elbow Method to identify number of clusters to choose
70 """
71
72 plt.plot(clusters, meandist)
73 plt.xlabel('Number of clusters')
74 plt.ylabel('Average distance')
75 plt.title('Selecting k with the Elbow Method')
76
77 # Interpret 3 cluster solution
78 model3=KMeans(n_clusters=3)
79 model3.fit(clus_train)
80 clusassign=model3.predict(clus_train)
81 # plot clusters
82
83 from sklearn.decomposition import PCA
84 pca_2 = PCA(2)
85 plot_columns = pca_2.fit_transform(clus_train)
86 plt.scatter(x=plot_columns[:,0], y=plot_columns[:,1], c=model3.labels_)
87 plt.xlabel('Canonical variable 1')
88 plt.ylabel('Canonical variable 2')
89 plt.title('Scatterplot of Canonical variables')
90 plt.show()
91 """
92 """
93 BEGIN multiple steps to merge cluster assignment with clustering variables to examine
```

PCA(2) asks Python to return the two first canonical variables.

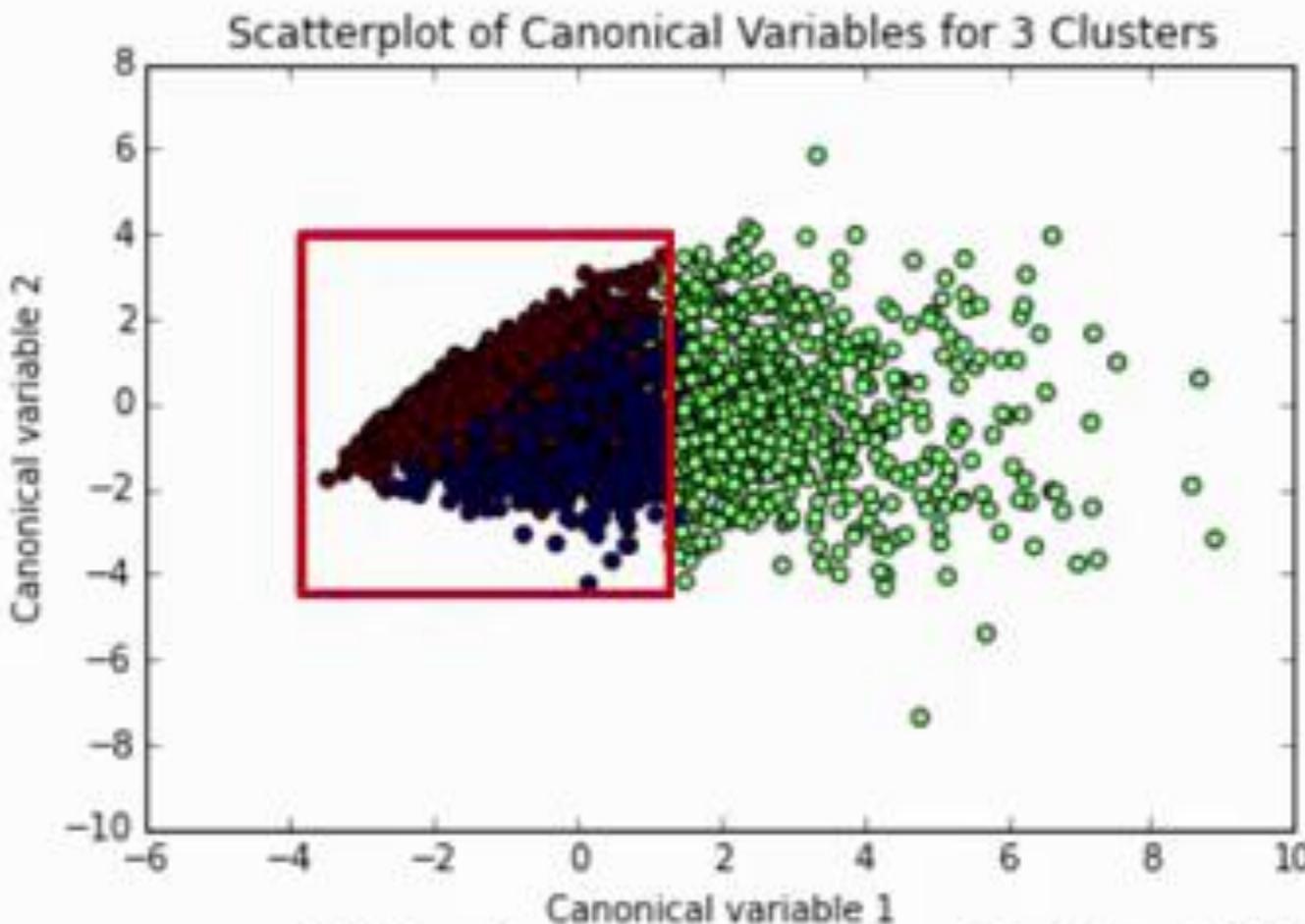
```
if self._edgecolors == str('face'):
```



In [13]:

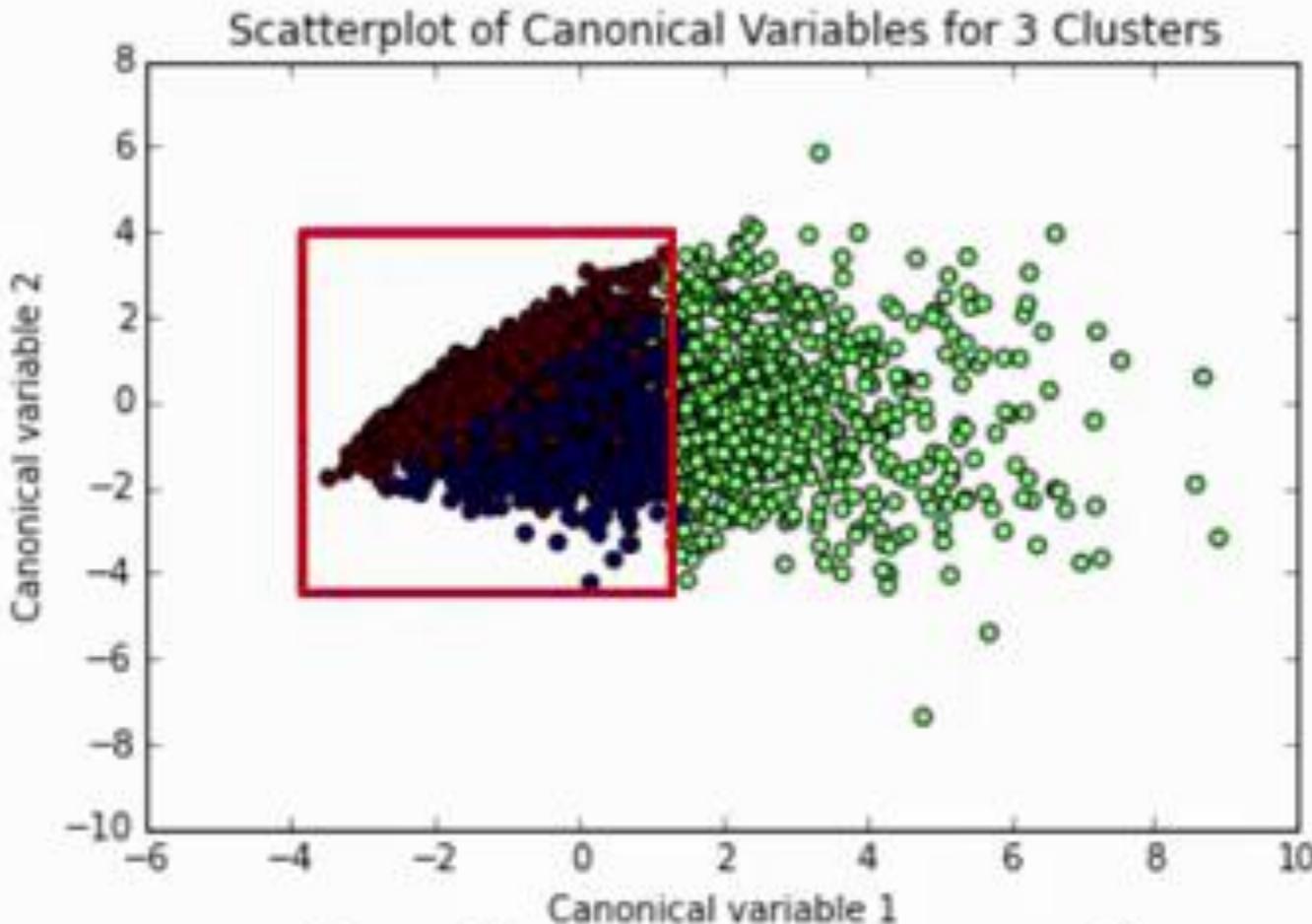
Here is the scatter plot.

```
if self._edgecolors == str('face'):
```



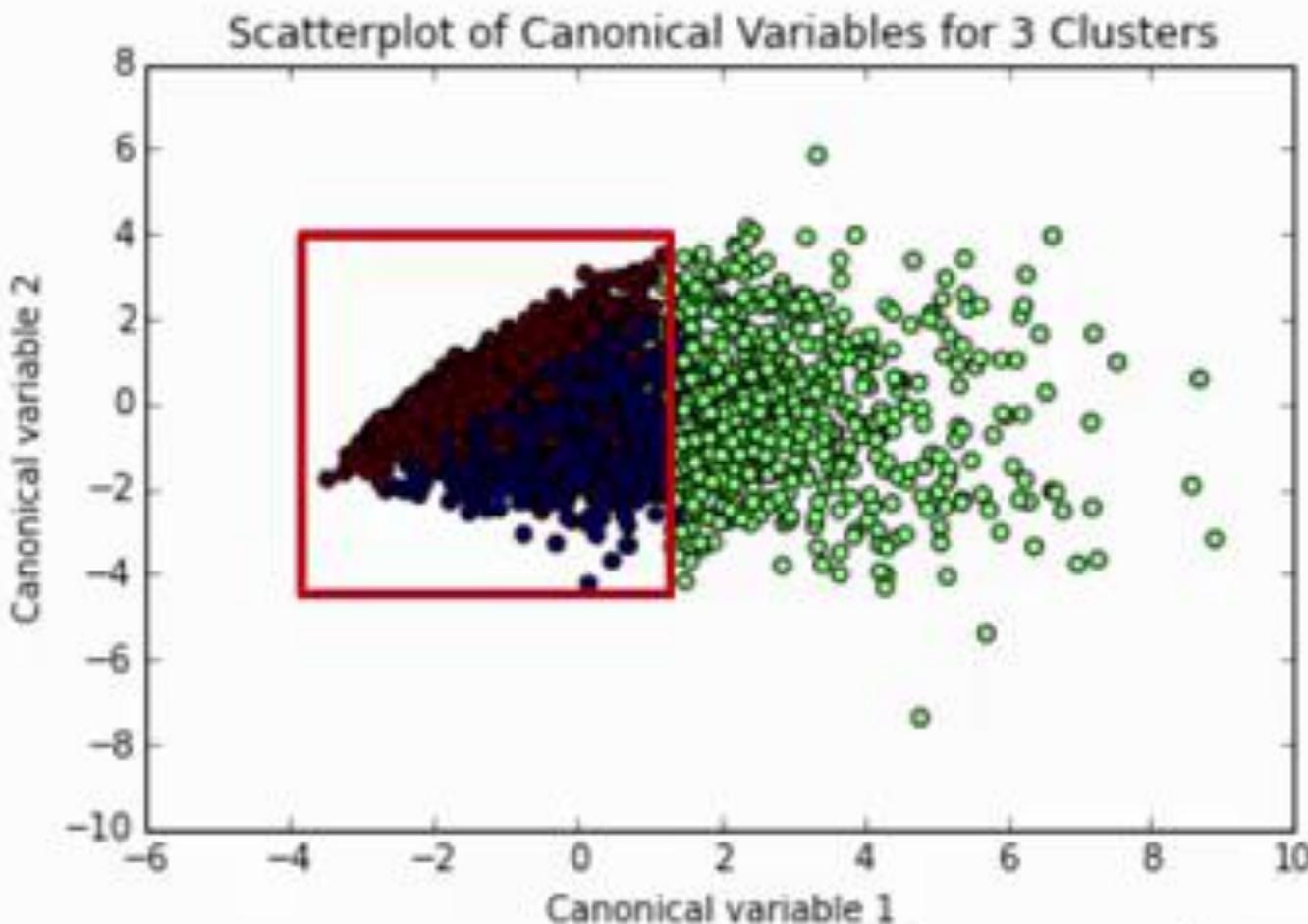
In [13]: What this shows is that these two clusters are densely packed,

```
if self._edgecolors == str('face'):
```



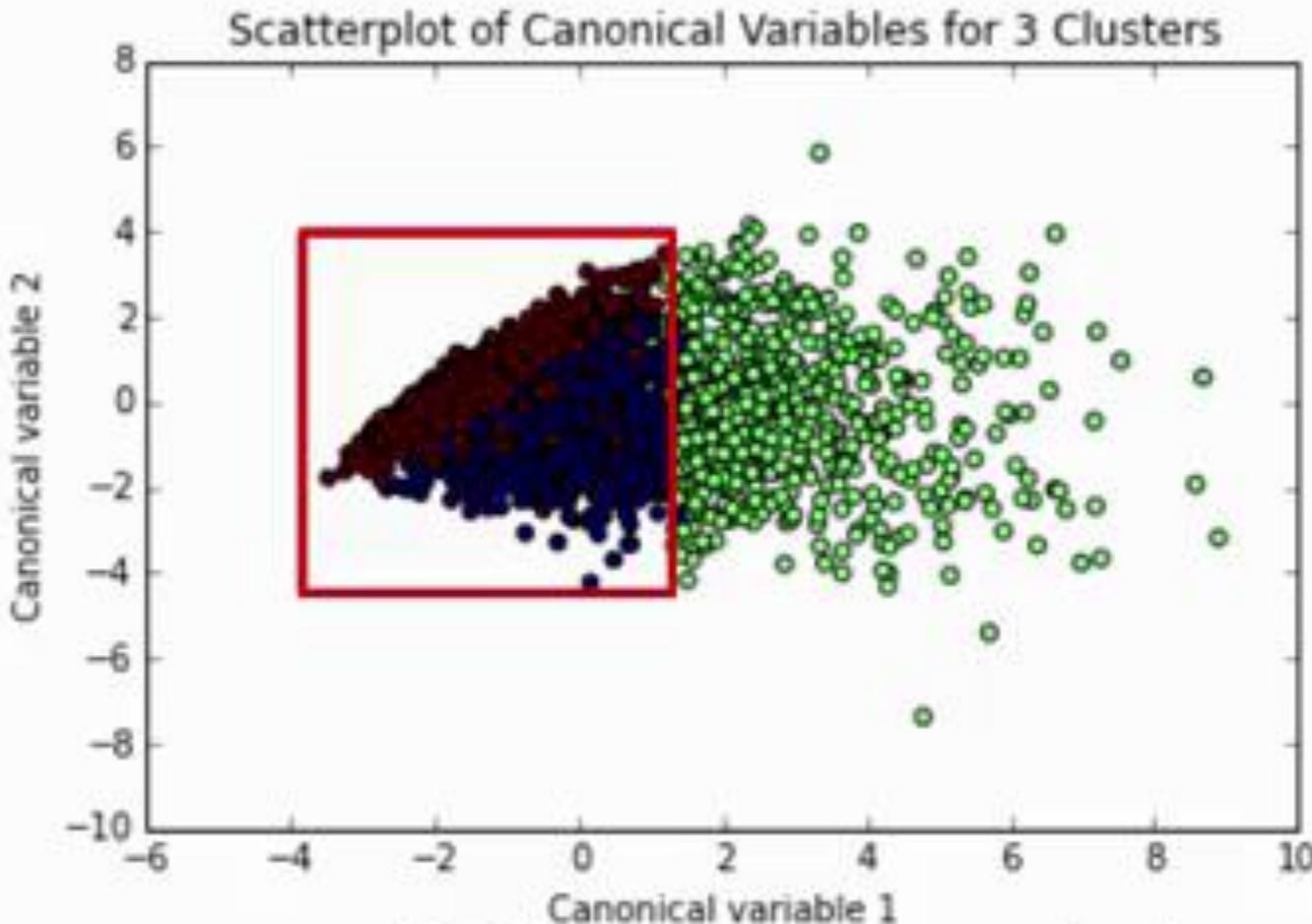
In [13]: meaning that the observations within the clusters are pretty highly correlated with

```
if self._edgecolors == str('face'):
```



each other, and
In [13]: within cluster variance is relatively low.

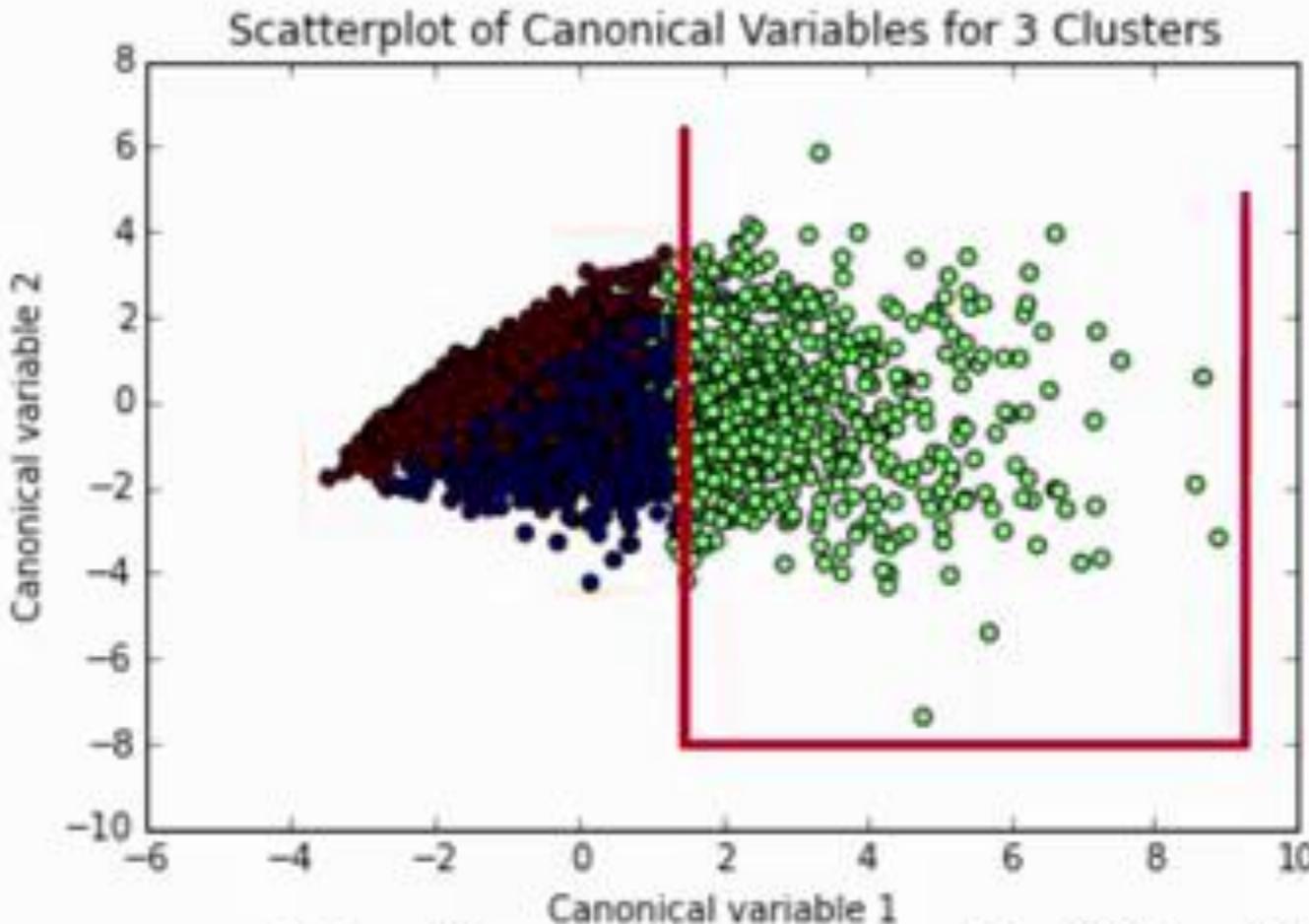
```
if self._edgecolors == str('face'):
```



In [13]:

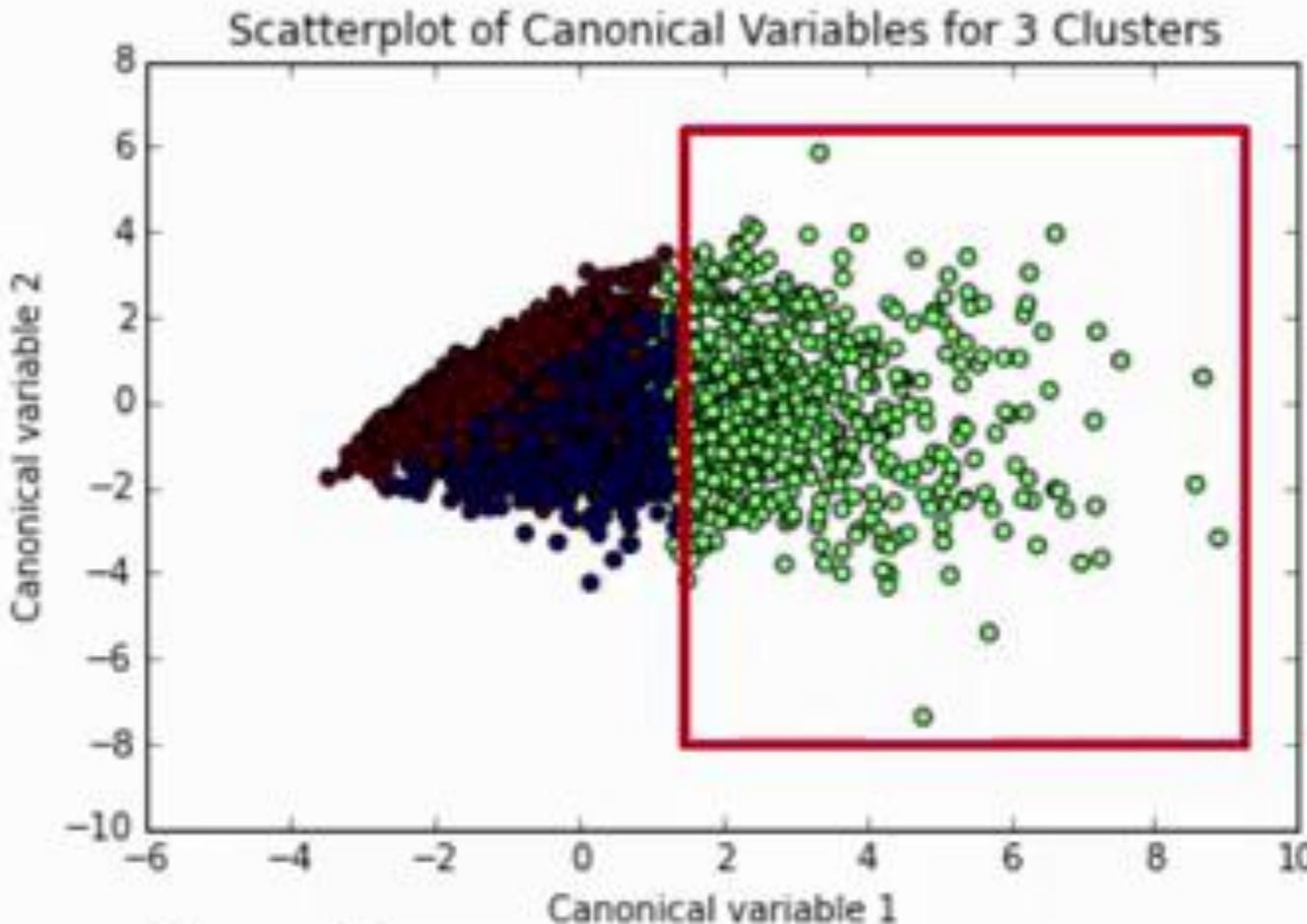
But they appear to have
a good deal of overlap,

```
if self._edgecolors == str('face'):
```



On the other hand, this cluster
In [13]: here shows better separation, but

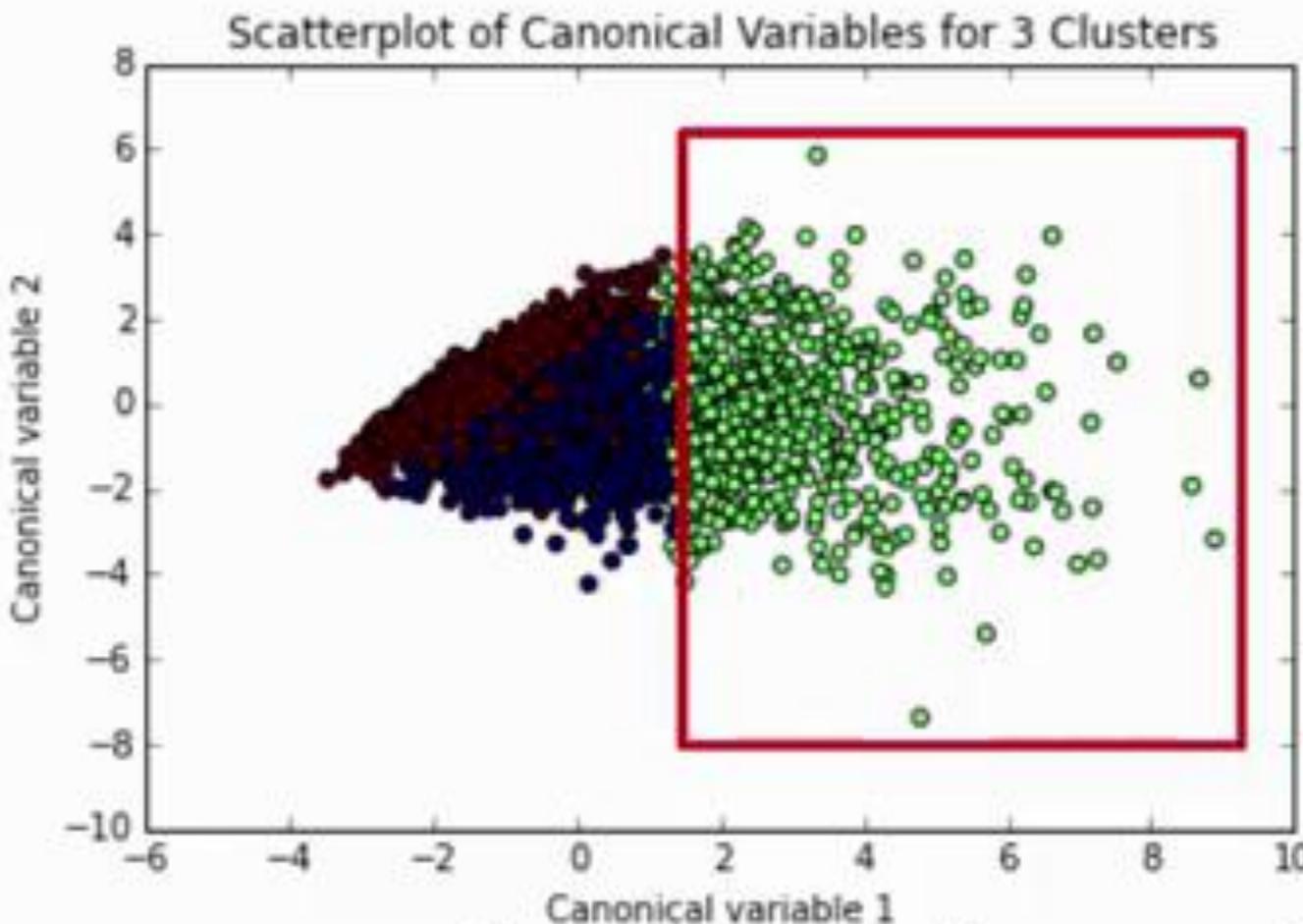
```
if self._edgecolors == str('face'):
```



the observations are more spread out

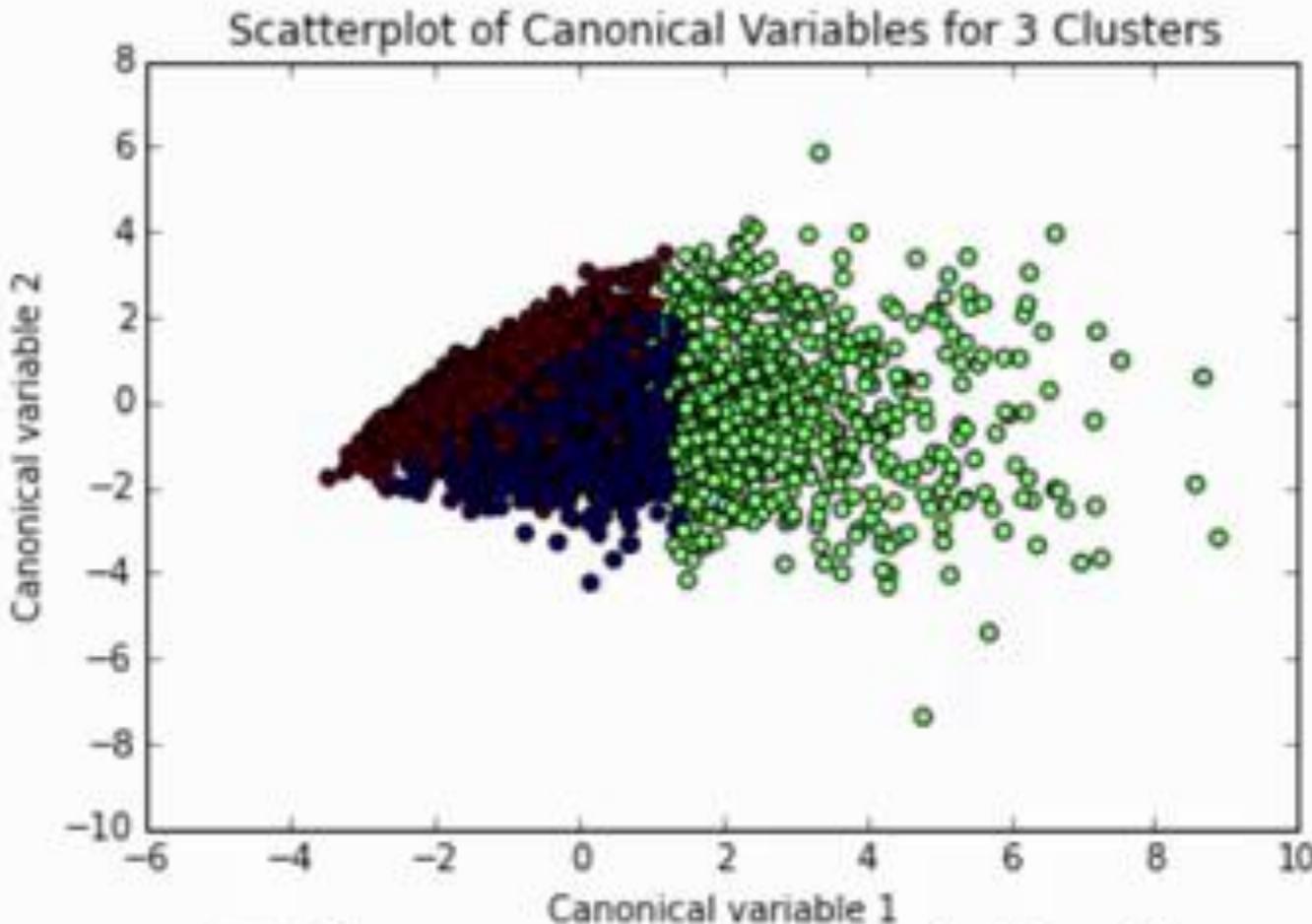
In [13]: indicating less correlation among

```
if self._edgecolors == str('face'):
```



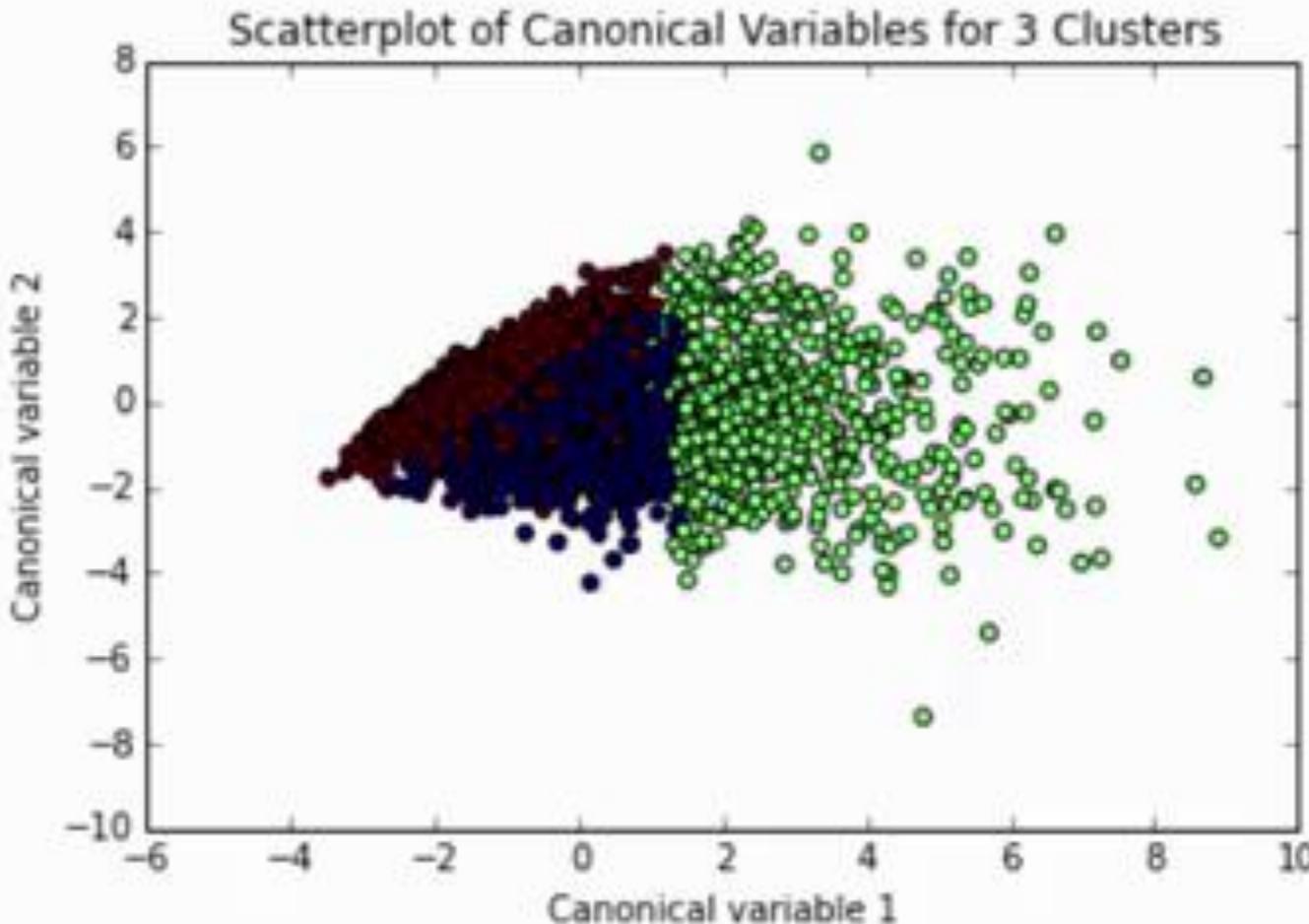
In [13]: the observations and higher within cluster variance.

```
if self._edgecolors == str('face'):
```



This suggests that the two cluster solution might be better, meaning that it

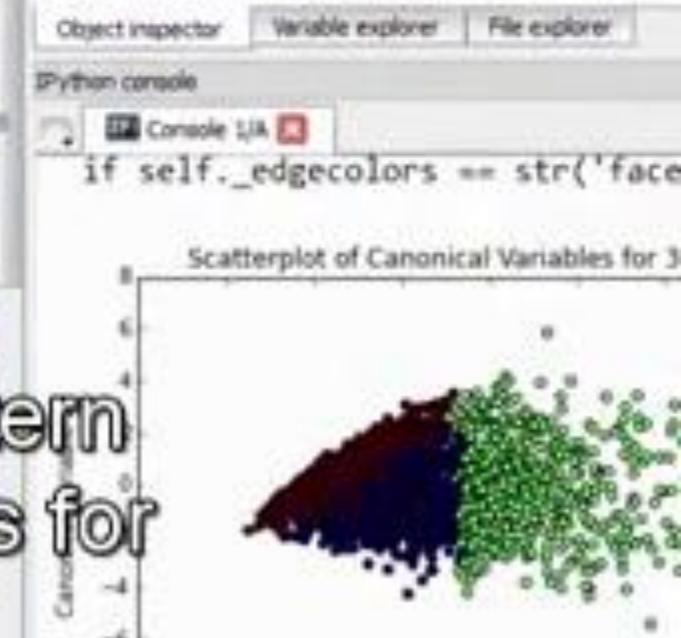
```
if self._edgecolors == str('face'):
```



In [13]: would be especially important to further evaluate the two cluster solution as well.

```
plt.show()
1
2 """
3 BEGIN multiple steps to merge cluster assignment with clustering variables to examine
4 cluster variable means by cluster
5 """
6 # create a unique identifier variable from the index for the
7 # cluster training data to merge with the cluster assignment variable
8 clus_train.reset_index(level=0, inplace=True)
9 # create a list that has the new index variable
10 cluslist=list(clus_train['index'])
11 # create a list of cluster assignments
12 labels=list(model3.labels_)
13 # combine index variable list with cluster assignment list into a dictionary
14 newlist=dict(zip(cluslist, labels))
15 newlist
16 # convert newlist dictionary to a dataframe
17 newclus=DataFrame.from_dict(newlist, orient='index')
18 newclus
19
20 # now do the same for the cluster assignment variable
21 # rename the cluster assignment column
22 newclus.columns = ['cluster']
23 # create a unique identifier variable from the index for the
24 # cluster assignment dataframe
25 # to merge with cluster training data
26 newclus.reset_index(level=0, inplace=True)
27 # merge the cluster assignment dataframe with the cluster training variable dataframe
28 # by the index variable
29 merged_train=pd.merge(clus_train, newclus, on='index')
30 merged_train.head(n=100)
31 # cluster frequencies
32 merged_train.cluster.value_counts()
33 """
```

Next we can take a look at the pattern of means on the clustering variables for



Usage

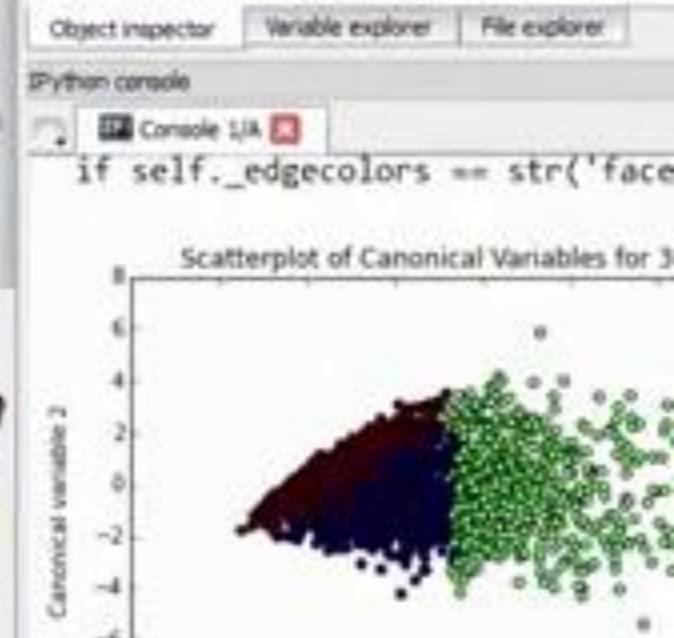
Here you can get help on it, either on the Editor or

Help can also be shown next to an object. You can use the Object Inspector.

New to

```
plt.show()
1
2 """
3 BEGIN multiple steps to merge cluster assignment with clustering variables to examine
4 cluster variable means by cluster
5 """
6 # create a unique identifier variable from the index for the
7 # cluster training data to merge with the cluster assignment variable
8 clus_train.reset_index(level=0, inplace=True)
9 # create a list that has the new index variable
10 cluslist=list(clus_train['index'])
11 # create a list of cluster assignments
12 labels=list(model3.labels_)
13 # combine index variable list with cluster assignment list into a dictionary
14 newlist=dict(zip(cluslist, labels))
15 newlist
16 # convert newlist dictionary to a dataframe
17 newclus=DataFrame.from_dict(newlist, orient='index')
18 newclus
19
20 # now do the same for the cluster assignment variable
21 # rename the cluster assignment column
22 newclus.columns = ['cluster']
23 # create a unique identifier variable from the index for the
24 # cluster assignment dataframe
25 # to merge with cluster training data
26 newclus.reset_index(level=0, inplace=True)
27 # merge the cluster assignment dataframe with the cluster training variable dataframe
28 # by the index variable
29 merged_train=pd.merge(clus_train, newclus, on='index')
30 merged_train.head(n=100)
31 # cluster frequencies
32 merged_train.cluster.value_counts()
33 """
```

each cluster to see whether they
are distinct and meaningful.



Usage

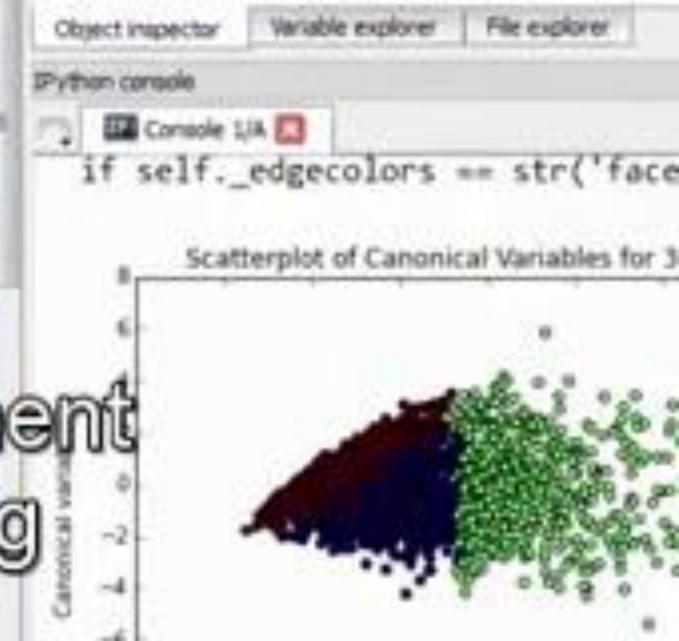
Here you can get help on it, either on the Editor or

Help can also be shown next to an object. You can use Object Inspector.

New to

```
plt.show()
1
2 """
3 BEGIN multiple steps to merge cluster assignment with clustering variables to examine
4 cluster variable means by cluster
5 """
6 # create a unique identifier variable from the index for the
7 # cluster training data to merge with the cluster assignment variable
8 clus_train.reset_index(level=0, inplace=True)
9 # create a list that has the new index variable
10 cluslist=list(clus_train['index'])
11 # create a list of cluster assignments
12 labels=list(model3.labels_)
13 # combine index variable list with cluster assignment list into a dictionary
14 newlist=dict(zip(cluslist, labels))
15 newlist
16 # convert newlist dictionary to a dataframe
17 newclus=DataFrame.from_dict(newlist, orient='index')
18 newclus
19
20 # now do the same for the cluster assignment variable
21 # rename the cluster assignment column
22 newclus.columns = ['cluster']
23 # create a unique identifier variable from the index for the
24 # cluster assignment dataframe
25 # to merge with cluster training data
26 newclus.reset_index(level=0, inplace=True)
27 # merge the cluster assignment dataframe with the cluster training variable dataframe
28 # by the index variable
29 merged_train=pd.merge(clus_train, newclus, on='index')
30 merged_train.head(n=100)
31 # cluster frequencies
32 merged_train.cluster.value_counts()
33 """
```

we have to link the cluster assignment
variable back to its corresponding



Usage

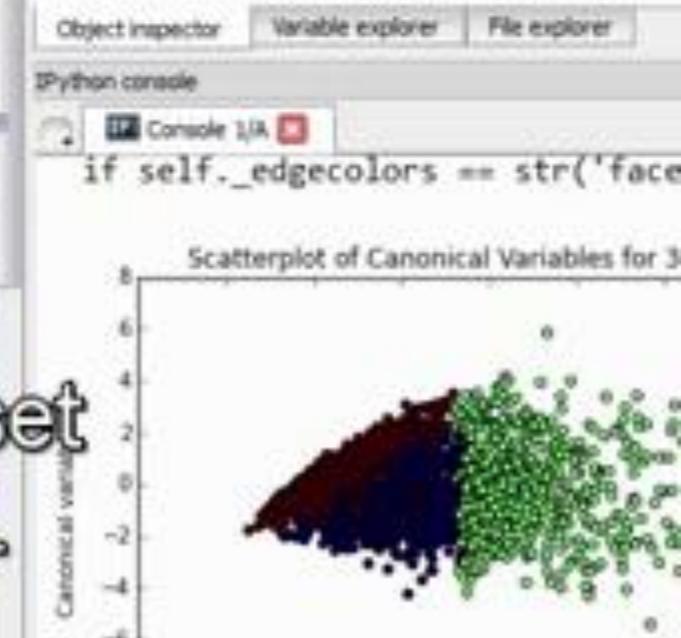
Here you can get help on it, either on the Editor or

Help can also be shown next to an object. You can use the Object Inspector.

New to

```
plt.show()
1
2 """
3 BEGIN multiple steps to merge cluster assignment with clustering variables to examine
4 cluster variable means by cluster
5 """
6 # create a unique identifier variable from the index for the
7 # cluster training data to merge with the cluster assignment variable
8 clus_train.reset_index(level=0, inplace=True)
9 # create a list that has the new index variable
10 cluslist=list(clus_train['index'])
11 # create a list of cluster assignments
12 labels=list(model3.labels_)
13 # combine index variable list with cluster assignment list into a dictionary
14 newlist=dict(zip(cluslist, labels))
15 newlist
16 # convert newlist dictionary to a dataframe
17 newclus=DataFrame.from_dict(newlist, orient='index')
18 newclus
19
20 # now do the same for the cluster assignment variable
21 # rename the cluster assignment column
22 newclus.columns = ['cluster']
23 # create a unique identifier variable from the index for the
24 # cluster assignment dataframe
25 # to merge with cluster training data
26 newclus.reset_index(level=0, inplace=True)
27 # merge the cluster assignment dataframe with the cluster training variable dataframe
28 # by the index variable
29 merged_train=pd.merge(clus@.head(n=100)
30 merged_train.head(n=100)
31 # cluster frequencies
32 merged_train.cluster.value_counts()
33 """
```

**observation in the clus_train dataset
that has the clustering variables.**



```
plt.show()
1
2 """
3 BEGIN multiple steps to merge cluster assignment with clustering variables to examine
4 cluster variable means by cluster
5 """
6 # create a unique identifier variable from the index for the
7 # cluster training data to merge with the cluster assignment variable
8 clus_train.reset_index(level=0, inplace=True)
9 # create a list that has the new index variable
10 cluslist=list(clus_train['index'])
11 # create a list of cluster assignments
12 labels=list(model3.labels_)
13 # combine index variable list with cluster assignment list into a dictionary
14 newlist=dict(zip(cluslist, labels))
15 newlist
16 # convert newlist dictionary to a dataframe
17 newclus=DataFrame.from_dict(newlist, orient='index')
18 newclus
19
20 # now do the same for the cluster assignment variable
21 # rename the cluster assignment column
22 newclus.columns = ['cluster']
23 # create a unique identifier variable from the index for the
24 # cluster assignment dataframe
25 # to merge with cluster training data
26 newclus.reset_index(level=0, inplace=True)
27 # merge the cluster assignment dataframe with the cluster training variable dataframe
28 # by the index variable
29 merged_train=pd.merge(clus_train, newclus, on='index')
30 merged_train.head(n=100)
31 # cluster frequencies
32 merged_train.cluster.value_counts()
33 """
```

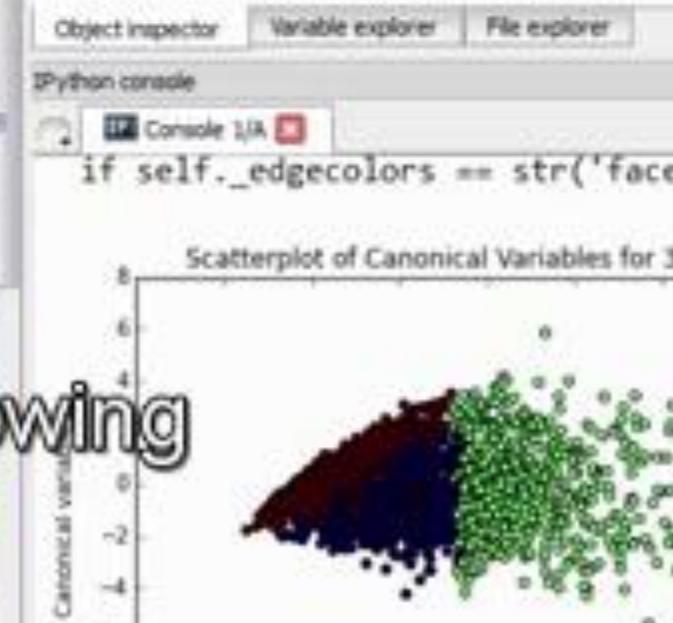
We will reset the index with the following code, `clus_train.reset_index`.

Usage

Here you can get help on it, either on the Editor or

Help can also be shown next to an object. You can use the Object Inspector.

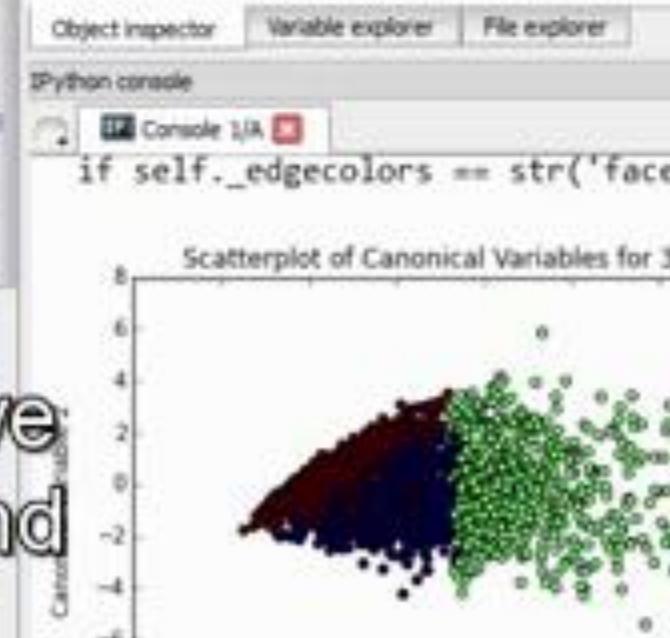
New to



```
plt.show()
1
2 """
3 BEGIN multiple steps to merge cluster assignment with clustering variables to examine
4 cluster variable means by cluster
5 """
6 # create a unique identifier variable from the index for the
7 # cluster training data to merge with the cluster assignment variable
8 clus_train.reset_index(level=0, inplace=True)
9 # create a list that has the new index variable
10 cluslist=list(clus_train['index'])
11 # create a list of cluster assignments
12 labels=list(model3.labels_)
13 # combine index variable list with cluster assignment list into a dictionary
14 newlist=dict(zip(cluslist, labels))
15 newlist
16 # convert newlist dictionary to a dataframe
17 newclus=DataFrame.from_dict(newlist, orient='index')
18 newclus
19
20 # now do the same for the cluster assignment variable
21 # rename the cluster assignment column
22 newclus.columns = ['cluster']
23 # create a unique identifier variable from the index for the
24 # cluster assignment dataframe
25 # to merge with cluster training data
26 newclus.reset_index(level=0, inplace=True)
27 # merge the cluster assignment dataframe with the cluster training variable dataframe
28 # by the index variable
29 merged_train=pd.merge(clus_t
30 merged_train.head(n=100)
31 # cluster frequencies
32 merged_train.cluster.value_counts()
33 """
```

level=0 tells Python to only remove
the given levels from the index, and

Usage
Here you can get help on an object. You can also search for objects in the Editor or Object Inspector.
Help can also be shown next to an object. You can click on the question mark icon to show it.
New to Python?



```
plt.show()
1
2 """
3 BEGIN multiple steps to merge cluster assignment with clustering variables to examine
4 cluster variable means by cluster
5 """
6 # create a unique identifier variable from the index for the
7 # cluster training data to merge with the cluster assignment variable
8 clus_train.reset_index(level=0, inplace=True)
9 # create a list that has the new index variable
10 cluslist=list(clus_train['index'])
11 # create a list of cluster assignments
12 labels=list(model3.labels_)
13 # combine index variable list with cluster assignment list into a dictionary
14 newlist=dict(zip(cluslist, labels))
15 newlist
16 # convert newlist dictionary to a dataframe
17 newclus=DataFrame.from_dict(newlist, orient='index')
18 newclus
19
20 # now do the same for the cluster assignment variable
21 # rename the cluster assignment column
22 newclus.columns = ['cluster']
23 # create a unique identifier variable from the index for the
24 # cluster assignment dataframe
25 # to merge with cluster training data
26 newclus.reset_index(level=0, inplace=True)
27 # merge the cluster assignment dataframe with the cluster training variable dataframe
28 # by the index variable
29 merged_train=pd.merge(clus_train, newclus, on='index')
30 merged_train.head(n=100)
31 # cluster frequencies
32 merged_train.cluster.value_counts()
33 """
```

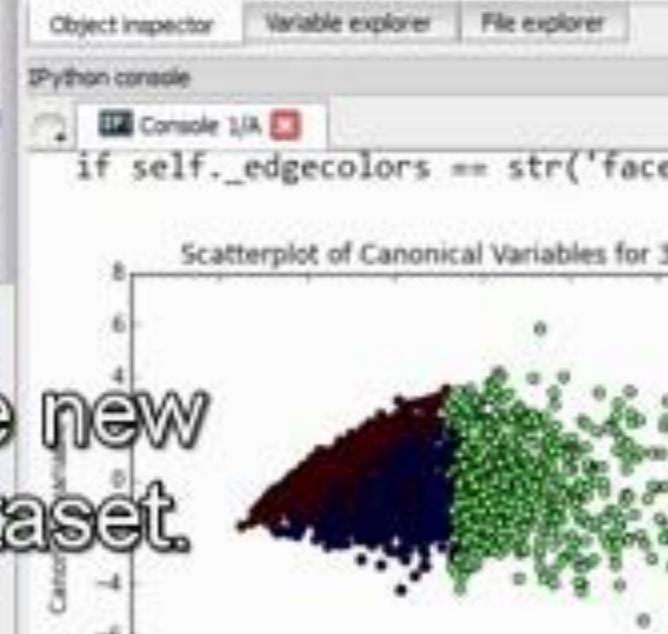
inplace=True, tells Python to add the new column to the existing `clus_train` dataset.

Usage

Here you can get help on it, either on the Editor or

Help can also be shown next to an object. You can use the Object Inspector.

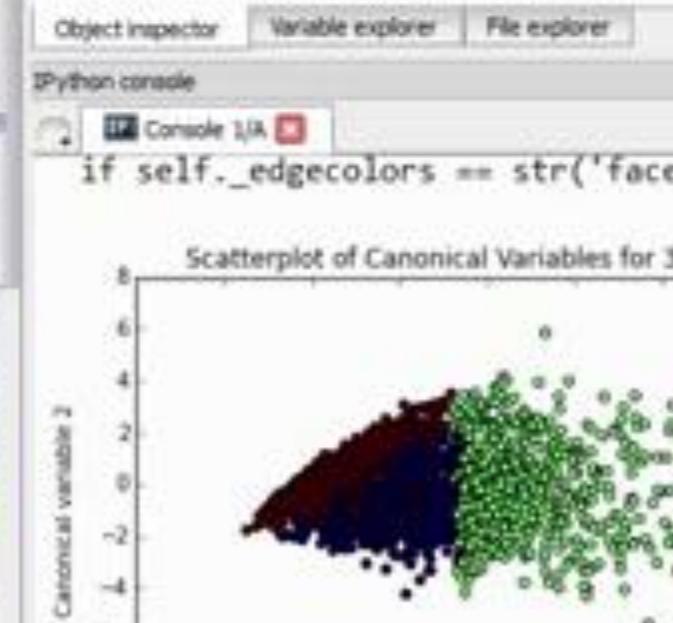
New to



```
plt.show()
1
2 """
3 BEGIN multiple steps to merge cluster assignment with clustering variables to examine
4 cluster variable means by cluster
5 """
6 # create a unique identifier variable from the index for the
7 # cluster training data to merge with the cluster assignment variable
8 clus_train.reset_index(level=0, inplace=True)
9 # create a list that has the new index variable
10 cluslist=list(clus_train['index'])
11 # create a list of cluster assignments
12 labels=list(model3.labels_)
13 # combine index variable list with cluster assignment list into a dictionary
14 newlist=dict(zip(cluslist, labels))
15 newlist
16 # convert newlist dictionary to a dataframe
17 newclus=DataFrame.from_dict(newlist, orient='index')
18 newclus
19
20 # now do the same for the cluster assignment variable
21 # rename the cluster assignment column
22 newclus.columns = ['cluster']
23 # create a unique identifier variable from the index for the
24 # cluster assignment dataframe
25 # to merge with cluster training data
26 newclus.reset_index(level=0, inplace=True)
27 # merge the cluster assignment dataframe with the cluster training variable dataframe
28 # by the index variable
29 merged_train=pd.merge(clus_train, newclus, on='index')
30 merged_train.head(n=100)
31 # cluster frequencies
32 merged_train.cluster.value_counts()
33 """
```

use the list function to
pull the new index variable

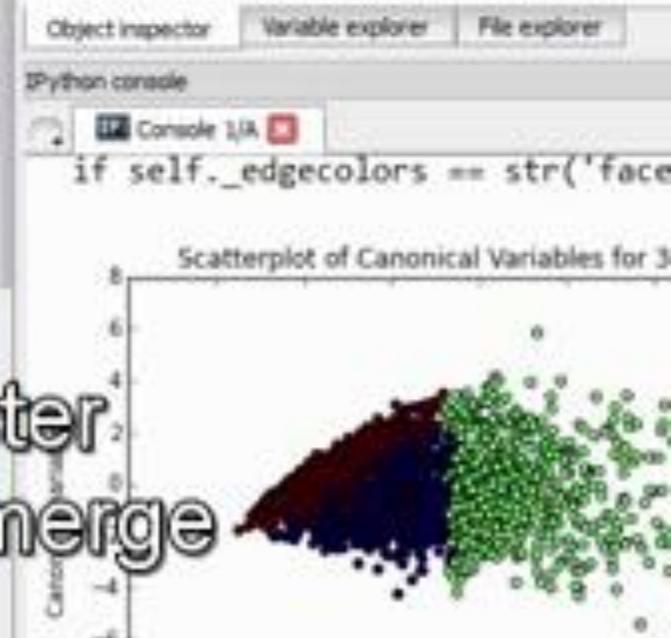
Usage
Here you can get help on it, either on the Editor or next to an object. You can also use the Object Inspector.



```
plt.show()
1
2 """
3 BEGIN multiple steps to merge cluster assignment with clustering variables to examine
4 cluster variable means by cluster
5 """
6 # create a unique identifier variable from the index for the
7 # cluster training data to merge with the cluster assignment variable
8 clus_train.reset_index(level=0, inplace=True)
9 # create a list that has the new index variable
10 cluslist=list(clus_train['index'])
11 # create a list of cluster assignments
12 labels=list(model3.labels_)
13 # combine index variable list with cluster assignment list into a dictionary
14 newlist=dict(zip(cluslist, labels))
15 newlist
16 # convert newlist dictionary to a dataframe
17 newclus=DataFrame.from_dict(newlist, orient='index')
18 newclus
19
20 # now do the same for the cluster assignment variable
21 # rename the cluster assignment column
22 newclus.columns = ['cluster']
23 # create a unique identifier variable from the index for the
24 # cluster assignment dataframe
25 # to merge with cluster training data
26 newclus.reset_index(level=0, inplace=True)
27 # merge the cluster assignment dataframe with the cluster training variable dataframe
28 # by the index variable
29 merged_train=pd.merge(clus_train, newclus, on='index')
30 merged_train.head(n=100)
31 # cluster frequencies
32 merged_train.cluster.value_counts()
33 """
```

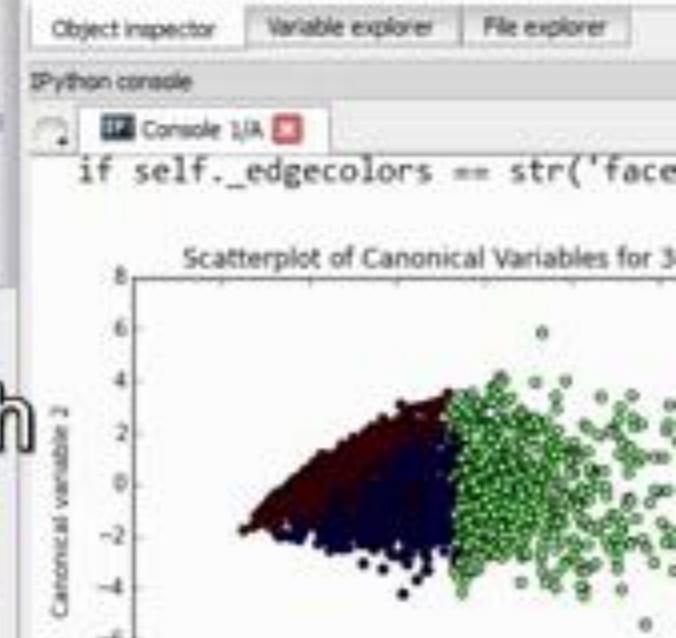
This will be combined with the cluster assignment variable, so that we can merge

Usage
Here you can get help on it, either on the Editor or next to an object. You can also use the Object Inspector.



```
plt.show()
1
2 """
3 BEGIN multiple steps to merge cluster assignment with clustering variables to examine
4 cluster variable means by cluster
5 """
6 # create a unique identifier variable from the index for the
7 # cluster training data to merge with the cluster assignment variable
8 clus_train.reset_index(level=0, inplace=True)
9 # create a list that has the new index variable
10 cluslist=list(clus_train['index'])
11 # create a list of cluster assignments
12 labels=list(model3.labels_)
13 # combine index variable list with cluster assignment list into a dictionary
14 newlist=dict(zip(cluslist, labels))
15 newlist
16 # convert newlist dictionary to a dataframe
17 newclus=DataFrame.from_dict(newlist, orient='index')
18 newclus
19
20 # now do the same for the cluster assignment variable
21 # rename the cluster assignment column
22 newclus.columns = ['cluster']
23 # create a unique identifier variable from the index for the
24 # cluster assignment dataframe
25 # to merge with cluster training data
26 newclus.reset_index(level=0, inplace=True)
27 # merge the cluster assignment dataframe with the cluster training variable dataframe
28 # by the index variable
29 merged_train=pd.merge(clus_train, newclus, on='index')
30 merged_train.head(n=100)
31 # cluster frequencies
32 merged_train.cluster.value_counts()
33 """
```

the two datasets together by each observation's unique identifier.



```
plt.show()
1
2 """
3 BEGIN multiple steps to merge cluster assignment with clustering variables to examine
4 cluster variable means by cluster
5 """
6 # create a unique identifier variable from the index for the
7 # cluster training data to merge with the cluster assignment variable
8 clus_train.reset_index(level=0, inplace=True)
9 # create a list that has the new index variable
10 cluslist=list(clus_train['index'])
11 # create a list of cluster assignments
12 labels=list(model3.labels_)
13 # combine index variable list with cluster assignment list into a dictionary
14 newlist=dict(zip(cluslist, labels))
15 newlist
16 # convert newlist dictionary to a dataframe
17 newclus=pd.DataFrame.from_dict(newlist, orient='index')
18 newclus
19
20 # now do the same for the cluster assignment variable
21 # rename the cluster assignment column
22 newclus.columns = ['cluster']
23 # create a unique identifier variable from the index for the
24 # cluster assignment dataframe
25 # to merge with cluster training data
26 newclus.reset_index(level=0, inplace=True)
27 # merge the cluster assignment dataframe with the cluster training variable dataframe
28 # by the index variable
29 merged_train=pd.merge(clus_train, newclus, on='index')
30 merged_train.head(n=100)
31 # cluster frequencies
32 merged_train.cluster.value_counts()
33 """
```

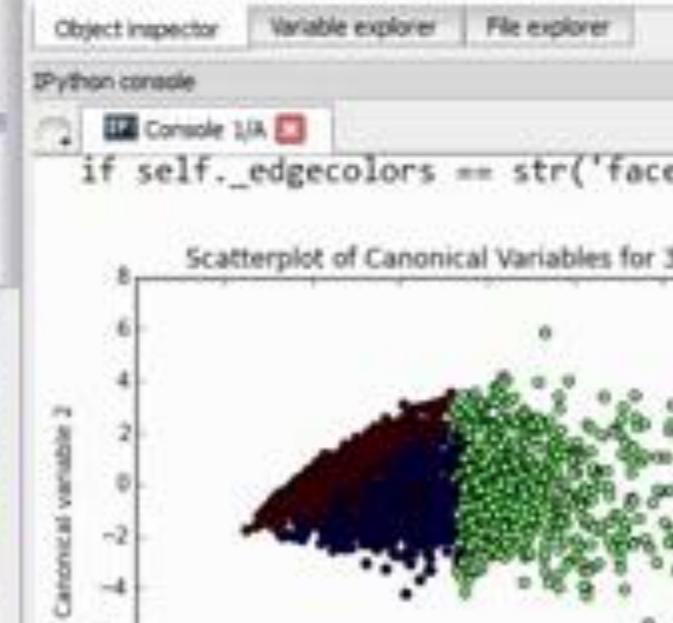
and cluster assignment
together into a single list.

Usage

Here you can get help on it, either on the Editor or

Help can also be shown next to an object. You can use the Object Inspector.

New to



```
plt.show()
1
2 """
3 BEGIN multiple steps to merge cluster assignment with clustering variables to examine
4 cluster variable means by cluster
5 """
6 # create a unique identifier variable from the index for the
7 # cluster training data to merge with the cluster assignment variable
8 clus_train.reset_index(level=0, inplace=True)
9 # create a list that has the new index variable
10 cluslist=list(clus_train['index'])
11 # create a list of cluster assignments
12 labels=list(model3.labels_)
13 # combine index variable list with cluster assignment list into a dictionary
14 newlist=dict(zip(cluslist, labels))
15 newlist
16 # convert newlist dictionary to a dataframe
17 newclus=DataFrame.from_dict(newlist, orient='index')
18 newclus
19
20 # now do the same for the cluster assignment variable
21 # rename the cluster assignment column
22 newclus.columns = ['cluster']
23 # create a unique identifier variable from the index for the
24 # cluster assignment dataframe
25 # to merge with cluster training data
26 newclus.reset_index(level=0, inplace=True)
27 # merge the cluster assignment dataframe with the cluster training variable dataframe
28 # by the index variable
29 merged_train=pd.merge(clus_train, newclus, on='index')
30 merged_train.head(n=100)
31 # cluster frequencies
32 merged_train.cluster.value_counts()
33 """
```

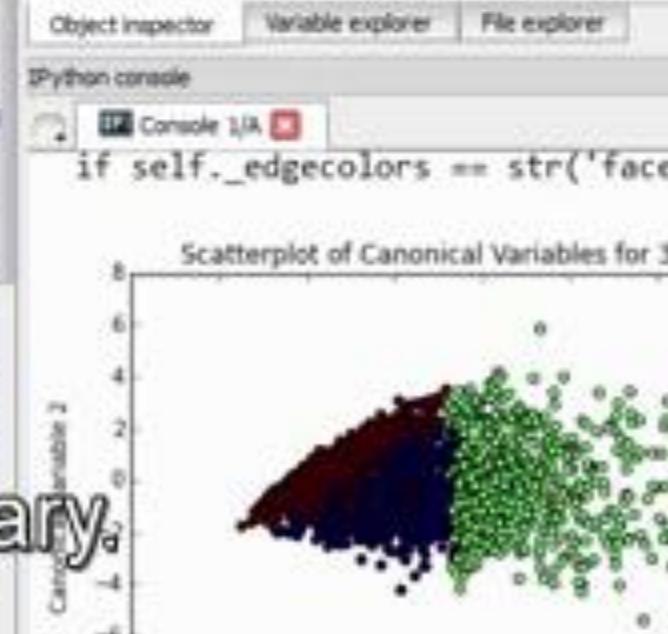
We'll call this list newlist using
a dict command to create a dictionary.

Usage

Here you can get help on it, either on the Editor or

Help can also be shown next to an object. You can use Object Inspector.

New to



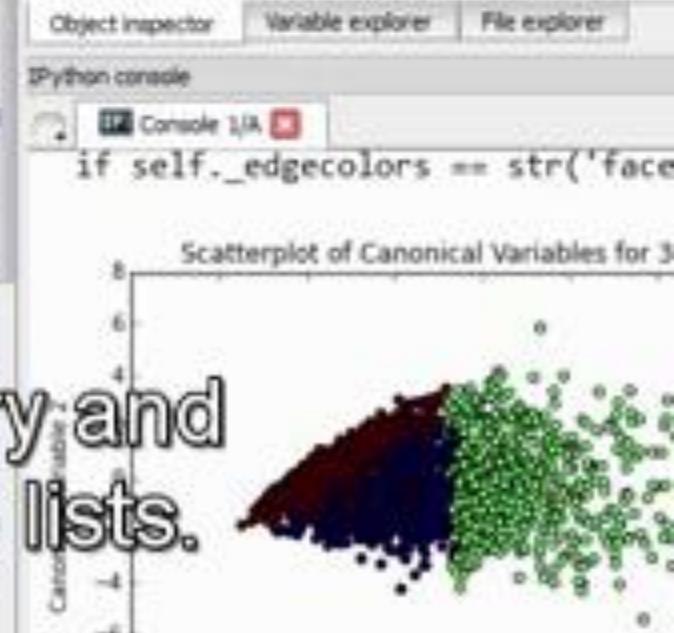
```
plt.show()
1
2 """
3 BEGIN multiple steps to merge cluster assignment with clustering variables to examine
4 cluster variable means by cluster
5 """
6 # create a unique identifier variable from the index for the
7 # cluster training data to merge with the cluster assignment variable
8 clus_train.reset_index(level=0, inplace=True)
9 # create a list that has the new index variable
10 cluslist=list(clus_train['index'])
11 # create a list of cluster assignments
12 labels=list(model3.labels_)
13 # combine index variable list with cluster assignment list into a dictionary
14 newlist=dict(zip(cluslist, labels))
15 newlist
16 # convert newlist dictionary to a dataframe
17 newclus=DataFrame.from_dict(newlist, orient='index')
18 newclus
19
20 # now do the same for the cluster assignment variable
21 # rename the cluster assignment column
22 newclus.columns = ['cluster']
23 # create a unique identifier variable from the index for the
24 # cluster assignment dataframe
25 # to merge with cluster training data
26 newclus.reset_index(level=0, inplace=True)
27 # merge the cluster assignment dataframe with the cluster training variable dataframe
28 # by the index vari
29 merged_train=pd.merge(clus_train, newclus)
30 merged_train.head(n=100)
31 # cluster frequencies
32 merged_train.cluster.value_counts()
33
34 """
35
```

Usage

Here you can get help on it, either on the Editor or

Help can also be shown next to an object. You can use the Object Inspector.

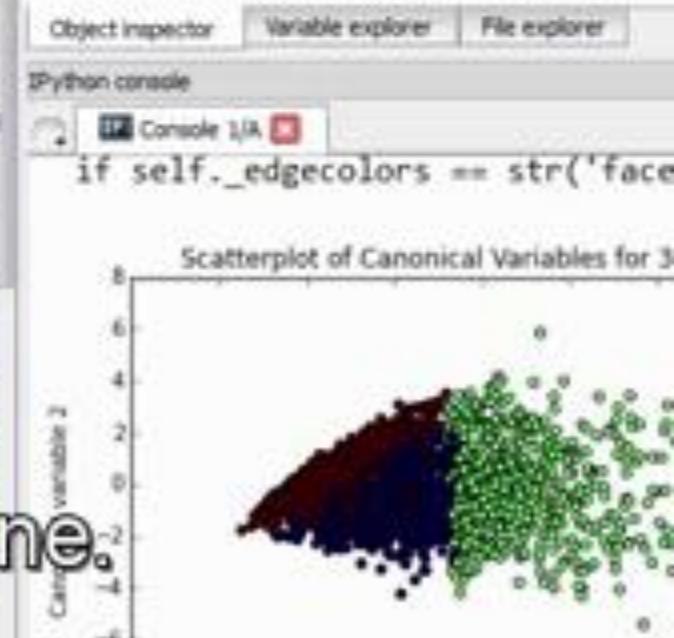
New to



The `dict` command creates a dictionary and the `zip` command is used to combine lists.

```
plt.show()
1
2 """
3 BEGIN multiple steps to merge cluster assignment with clustering variables to examine
4 cluster variable means by cluster
5 """
6 # create a unique identifier variable from the index for the
7 # cluster training data to merge with the cluster assignment variable
8 clus_train.reset_index(level=0, inplace=True)
9 # create a list that has the new index variable
10 cluslist=list(clus_train['index'])
11 # create a list of cluster assignments
12 labels=list(model3.labels_)
13 # combine index variable list with cluster assignment list into a dictionary
14 newlist=dict(zip(cluslist, labels))
15 newlist
16 # convert newlist dictionary to a dataframe
17 newclusDataframe.from_dict(newlist, orient='index')
18 newclus
19
20 # now do the same for the cluster assignment variable
21 # rename the cluster assignment column
22 newclus.columns = ['cluster']
23 # create a unique identifier variable from the index for the
24 # cluster assignment dataframe
25 # to merge with cluster training data
26 newclus.reset_index(level=0, inplace=True)
27 # merge the cluster assignment dataframe with the cluster training variable dataframe
28 # by the index variable
29 merged_train=pd.merge(clus_train, newclus, on='index')
30 merged_train.head(n=100)
31 # cluster frequencies
32 merged_train.cluster.value_counts()
33 """
```

So here we type `dict(zip)` and the two lists that we want to combine.



Usage

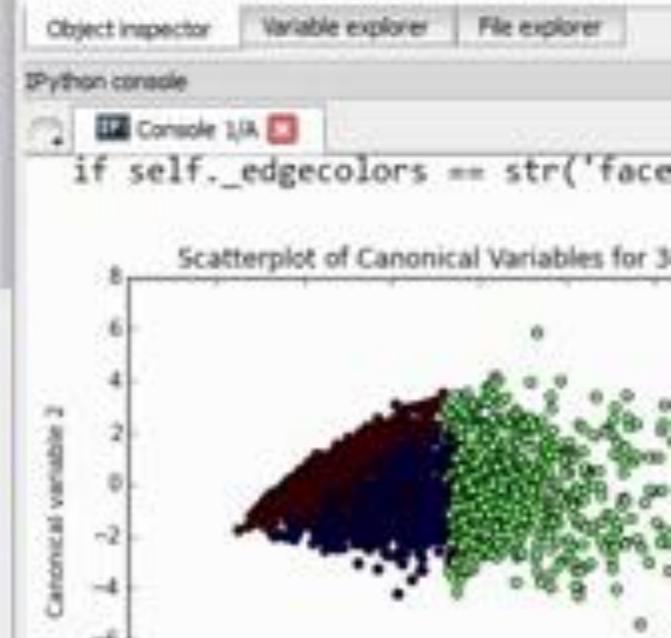
Here you can get help on it, either on the Editor or

Help can also be shown next to an object. You can use the Object Inspector.

New to

```
plt.show()
1
2 """
3 BEGIN multiple steps to merge cluster assignment with clustering variables to examine
4 cluster variable means by cluster
5 """
6 # create a unique identifier variable from the index for the
7 # cluster training data to merge with the cluster assignment variable
8 clus_train.reset_index(level=0, inplace=True)
9 # create a list that has the new index variable
10 cluslist=list(clus_train['index'])
11 # create a list of cluster assignments
12 labels=list(model3.labels_)
13 # combine index variable list with cluster assignment list into a dictionary
14 newlist=dict(zip(cluslist, labels))
15 newlist
16 # convert newlist dictionary to a dataframe
17 newclus=pd.DataFrame.from_dict(newlist, orient='index')
18 newclus
19
20 # now do the same for the cluster assignment variable
21 # rename the cluster assignment column
22 newclus.columns = ['cluster']
23 # create a unique identifier variable from the index for the
24 # cluster assignment dataframe
25 # to merge with cluster training data
26 newclus.reset_index(level=0, inplace=True)
27 # merge the cluster assignment dataframe with the cluster training variable dataframe
28 # by the index variable
29 merged_train=pd.merge(clus_train, newclus, on='index')
30 merged_train.head(n=100)
31 # cluster frequencies
32 merged_train.cluster.value_counts()
33 """
```

Finally, we convert the newlist dictionary to a data frame using



Usage

Here you can get help on it, either on the Editor or

Help can also be shown next to an object. You can use Object Inspector.

New to

```
plt.show()
1
2 """
3 BEGIN multiple steps to merge cluster assignment with clustering variables to examine
4 cluster variable means by cluster
5 """
6 # create a unique identifier variable from the index for the
7 # cluster training data to merge with the cluster assignment variable
8 clus_train.reset_index(level=0, inplace=True)
9 # create a list that has the new index variable
10 cluslist=list(clus_train['index'])
11 # create a list of cluster assignments
12 labels=list(model3.labels_)
13 # combine index variable list with cluster assignment list into a dictionary
14 newlist=dict(zip(cluslist, labels))
15 newlist
16 # convert newlist dictionary to a dataframe
17 newclus=pd.DataFrame.from_dict(newlist, orient='index')
18 newclus
19
20 # now do the same for the cluster assignment variable
21 # rename the cluster assignment column
22 newclus.columns = ['cluster']
23 # create a unique identifier variable from the index for the
24 # cluster assignment dataframe
25 # to merge with cluster training data
26 newclus.reset_index(level=0, inplace=True)
27 # merge the cluster assignment dataframe with the cluster training variable dataframe
28 # by the index variable
29 merged_train=pd.merge(clus_train, newclus, on='index')
30 merged_train.head(n=100)
31 # cluster frequencies
32 merged_train.cluster.value_counts()
33 """
```

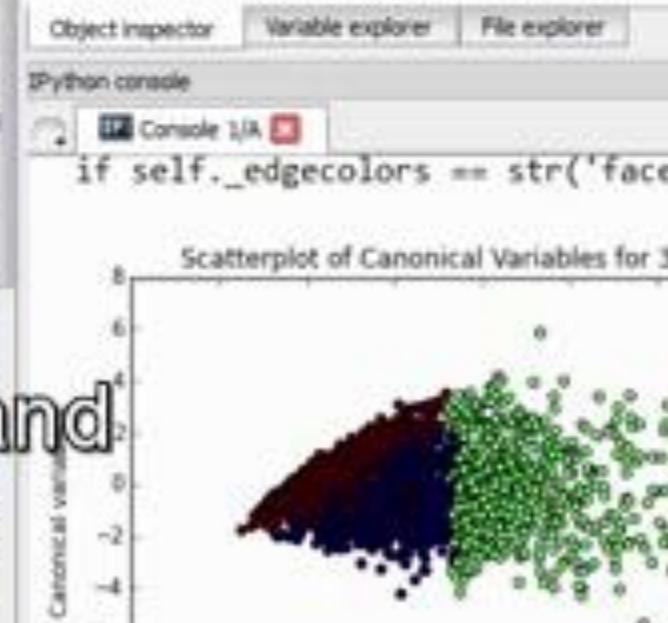
the `pd.DataFrame.from_dict` function and
rename the cluster assignment

Usage

Here you can get help on it, either on the Editor or

Help can also be shown next to an object. You can use the Object Inspector.

New to



```
plt.show()
1
2 """
3 BEGIN multiple steps to merge cluster assignment with clustering variables to examine
4 cluster variable means by cluster
5 """
6 # create a unique identifier variable from the index for the
7 # cluster training data to merge with the cluster assignment variable
8 clus_train.reset_index(level=0, inplace=True)
9 # create a list that has the new index variable
10 cluslist=list(clus_train['index'])
11 # create a list of cluster assignments
12 labels=list(model3.labels_)
13 # combine index variable list with cluster assignment list into a dictionary
14 newlist=dict(zip(cluslist, labels))
15 newlist
16 # convert newlist dictionary to a dataframe
17 newclus=DataFrame.from_dict(newlist, orient='index')
18 newclus
19
20 # now do the same for the cluster assignment variable
21 # rename the cluster assignment column
22 newclus.columns = ['cluster']
23 # create a unique identifier variable from the index for the
24 # cluster assignment dataframe
25 # to merge with cluster training data
26 newclus.reset_index(level=0, inplace=True)
27 # merge the cluster assignment dataframe with the cluster training variable dataframe
28 # by the index variable
29 merged_train=pd.merge(clus_train, newclus, on='index')
30 merged_train.head(n=100)
31 # cluster frequencies
32 merged_train.cluster.value_counts()
33 """
```

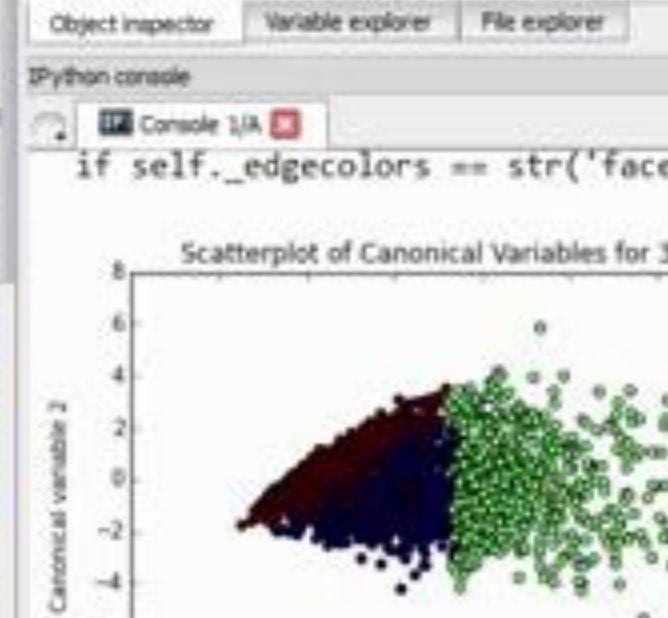
variable to cluster by typing
newclus.columns=cluster.

Usage

Here you can get help on it, either on the Editor or

Help can also be shown next to an object. You can use Object Inspector.

New to



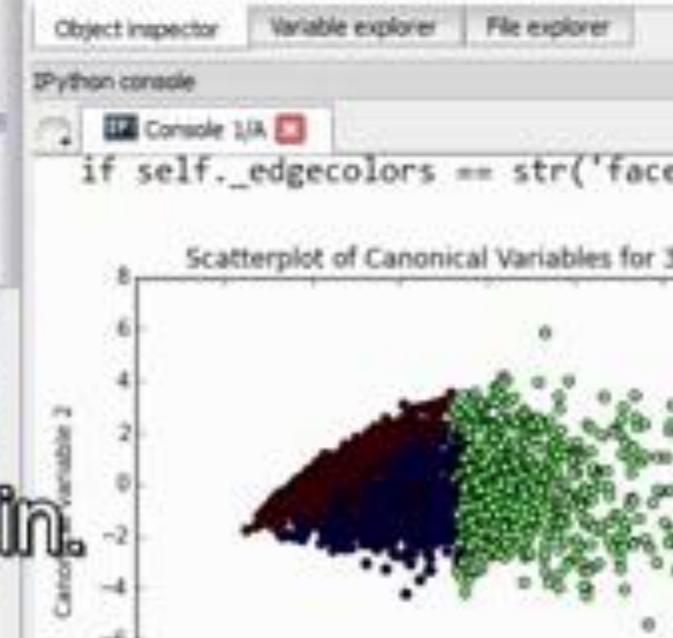
```
plt.show()
1
2 """
3 BEGIN multiple steps to merge cluster assignment with clustering variables to examine
4 cluster variable means by cluster
5 """
6 # create a unique identifier variable from the index for the
7 # cluster training data to merge with the cluster assignment variable
8 clus_train.reset_index(level=0, inplace=True)
9 # create a list that has the new index variable
10 cluslist=list(clus_train['index'])
11 # create a list of cluster assignments
12 labels=list(model3.labels_)
13 # combine index variable list with cluster assignment list into a dictionary
14 newlist=dict(zip(cluslist, labels))
15 newlist
16 # convert newlist dictionary to a dataframe
17 newclus=DataFrame.from_dict(newlist, orient='index')
18 newclus
19
20 # now do the same for the cluster assignment variable
21 # rename the cluster assignment column
22 newclus.columns = ['cluster']
23 # create a unique identifier variable from the index for the
24 # cluster assignment dataframe
25 # to merge with cluster training data
26 newclus.reset_index(level=0, inplace=True)
27 # merge the cluster assignment dataframe with the cluster training variable dataframe
28 # by the index variable
29 merged_train=pd.merge(clus_train, newclus, on='index')
30 merged_train.head(n=100)
31 # cluster frequencies
32 merged_train.cluster.value_counts(),
33 """
the newclus data frame using
the model.vista_index function again.
```

Usage

Here you can get help on it, either on the Editor or

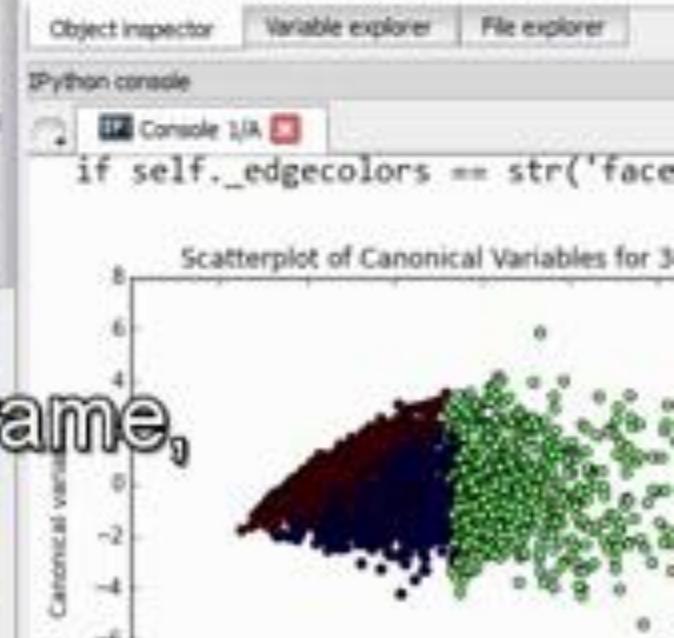
Help can also be shown next to an object. You can use Object Inspector.

New to



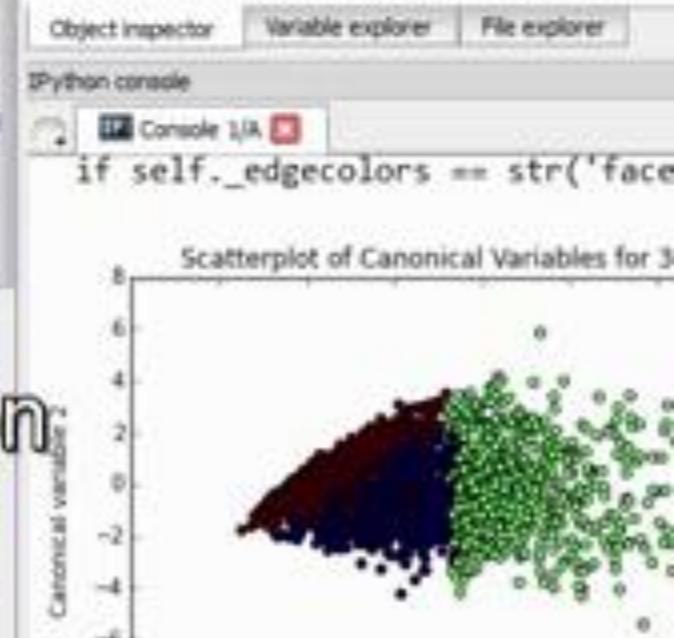
```
plt.show()
1
2 """
3 BEGIN multiple steps to merge cluster assignment with clustering variables to examine
4 cluster variable means by cluster
5 """
6 # create a unique identifier variable from the index for the
7 # cluster training data to merge with the cluster assignment variable
8 clus_train.reset_index(level=0, inplace=True)
9 # create a list that has the new index variable
10 cluslist=list(clus_train['index'])
11 # create a list of cluster assignments
12 labels=list(model3.labels_)
13 # combine index variable list with cluster assignment list into a dictionary
14 newlist=dict(zip(cluslist, labels))
15 newlist
16 # convert newlist dictionary to a dataframe
17 newclus=DataFrame.from_dict(newlist, orient='index')
18 newclus
19
20 # now do the same for the cluster assignment variable
21 # rename the cluster assignment column
22 newclus.columns = ['cluster']
23 # create a unique identifier variable from the index for the
24 # cluster assignment dataframe
25 # to merge with cluster training data
26 newclus.reset_index(level=0, inplace=True)
27 # merge the cluster assignment dataframe with the cluster training variable dataframe
28 # by the index variable
29 merged_train=pd.merge(clus_train, newclus, on='index')
30 merged_train.head(n=100)
31 # cluster frequencies
32 merged_train.cluster.value_counts()
33 """
```

We do this by creating a new data frame, called merged_train, and



```
plt.show()
1
2 """
3 BEGIN multiple steps to merge cluster assignment with clustering variables to examine
4 cluster variable means by cluster
5 """
6 # create a unique identifier variable from the index for the
7 # cluster training data to merge with the cluster assignment variable
8 clus_train.reset_index(level=0, inplace=True)
9 # create a list that has the new index variable
10 cluslist=list(clus_train['index'])
11 # create a list of cluster assignments
12 labels=list(model3.labels_)
13 # combine index variable list with cluster assignment list into a dictionary
14 newlist=dict(zip(cluslist, labels))
15 newlist
16 # convert newlist dictionary to a dataframe
17 newclus=DataFrame.from_dict(newlist, orient='index')
18 newclus
19
20 # now do the same for the cluster assignment variable
21 # rename the cluster assignment column
22 newclus.columns = ['cluster']
23 # create a unique identifier variable from the index for the
24 # cluster assignment dataframe
25 # to merge with cluster training data
26 newclus.reset_index(level=0, inplace=True)
27 # merge the cluster assignment dataframe with the cluster training variable dataframe
28 # by the index variable
29 merged_train=pd.merge(clus_t
30 merged_train.head(n=100)
31 # cluster frequencies
32 merged_train.cluster.value_counts()
33 """
```

we use the pandas merge function to merge the two data sets.



Usage

Here you can get help on it, either on the Editor or

Help can also be shown next to an object. You can use Object Inspector.

New to

```
plt.show()
1
2 """
3 BEGIN multiple steps to merge cluster assignment with clustering variables to examine
4 cluster variable means by cluster
5 """
6 # create a unique identifier variable from the index for the
7 # cluster training data to merge with the cluster assignment variable
8 clus_train.reset_index(level=0, inplace=True)
9 # create a list that has the new index variable
10 cluslist=list(clus_train['index'])
11 # create a list of cluster assignments
12 labels=list(model3.labels_)
13 # combine index variable list with cluster assignment list into a dictionary
14 newlist=dict(zip(cluslist, labels))
15 newlist
16 # convert newlist dictionary to a dataframe
17 newclus=DataFrame.from_dict(newlist, orient='index')
18 newclus
19
20 # now do the same for the cluster assignment variable
21 # rename the cluster assignment column
22 newclus.columns = ['cluster']
23 # create a unique identifier variable from the index for the
24 # cluster assignment dataframe
25 # to merge with cluster training data
26 newclus.reset_index(level=0, inplace=True)
27 # merge the cluster assignment dataframe with the cluster training variable dataframe
28 # by the index variable
29 merged_train=pd.merge(clus_train, newclus, on='index')
30 merged_train.head(n=100)
31 # cluster frequencies
32 merged_train.cluster.value_counts()
33 """
```

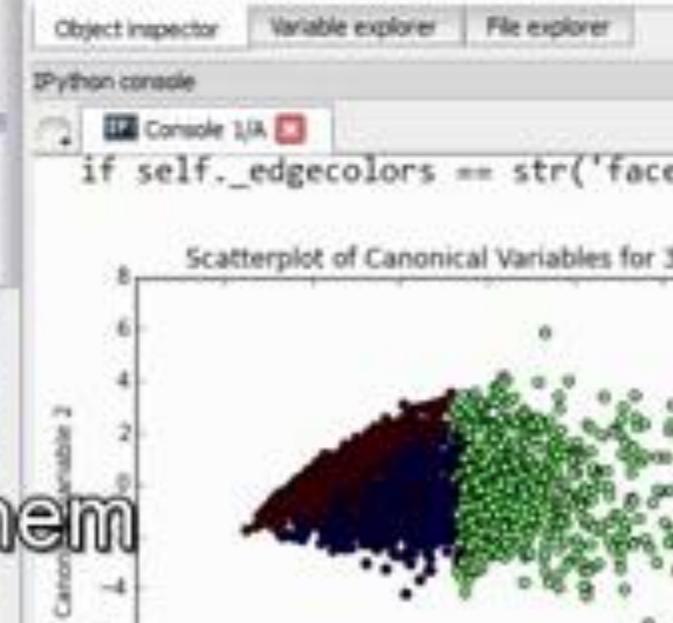
On="index" tells Python to merge them

Usage

Here you can get help on it, either on the Editor or

Help can also be shown next to an object. You can use the Object Inspector.

New to



```
1 # cluster frequencies
2 merged_train.cluster.value_counts()
3
4 """
5 END multiple steps to merge cluster assignment with clustering variables to e
6 cluster variable means by cluster
7 """
8
9 # FINALLY calculate clustering variable means by cluster
10 clustergrp = merged_train.groupby('cluster').mean()
11 print ("Clustering variable means by cluster")
12 print(clustergrp)
13
14
15 # validate clusters in training data by examining cluster differences in GPA
16 # first have to merge GPA with clustering variables and cluster assignment da
17 gpa_data=data_clean[["GPA","cluster"]].dropna() Finally, we can get the means on
18 # split GPA data into train and test sets
19 gpa_train, gpa_test=gpa_data.sample(frac=0.85, random_state=123), random_state=123
20 gpa_train1=pd.DataFrame(gpa_train)
```

```
1 # cluster frequencies
2 merged_train.cluster.value_counts()
3
4 """
5 END multiple steps to merge cluster assignment with clustering variables to e
6 cluster variable means by cluster
7 """
8
9 # FINALLY calculate clustering variable means by cluster
10 clustergrp = merged_train.groupby('cluster').mean()
11 print ("Clustering variable means by cluster")
12 print(clustergrp)
13
14
15 # validate clusters in training data by examining cluster differences in GPA
16 # first have to merge GPA with clustering variables and cluster assignment da
17 gpa_data=data_clean['GPA1']
18 # split GPA data into train and test sets
19 gpa_train, gpa_all=all.the clustering.variables.by cluster. 3, random_state=1
20 gpa_train1=pd.DataFrame(gpa_train)
```

```
....  
....:  
....: # FINALLY calculate clustergrp = merger  
....: print ("Clustering")  
....: print(clustergrp)  
....:
```

Clustering variable means by cluster

| | index | ALCEVR1 | MAREVER1 | ALCPROBS1 | WT1 | VIOL1 | \ |
|---------|-------------|-----------|-----------|-----------|-----------|-----------|---|
| cluster | | | | | | | |
| 0 | 3328.930505 | 0.946562 | -0.057109 | -0.058550 | -6...1192 | -0.168252 | |
| 1 | 3312.255193 | 0.661266 | 1.097891 | 0.906 | 1.106490 | 0.791967 | |
| 2 | 3239.829577 | -1.056455 | -0.474543 | -0.412562 | -0.451110 | -0.264092 | |

| | DEP1 | ESTEEM1 | SCHCONN1 | PARACTV | PARPRES | FAMCONCT | |
|---------|-----------|-----------|-----------|-----------|-----------|-----------|--|
| cluster | | | | | | | |
| 0 | -0.196987 | 0.188367 | 0.131356 | 0.151205 | 0.099521 | 0.226842 | |
| 1 | 0.853587 | -0.655070 | -0.927499 | -0.409328 | -0.473980 | -0.955398 | |
| 2 | -0.292352 | 0.206912 | 0.341590 | 0.091216 | 0.156604 | 0.298288 | |

In [14]:

Remember that the variables are standardized to be on the same scale with an overall sample mean of zero, and a standard deviation of 1

```
....:  
....: # FINALLY calculate clustering variable means by cluster  
....: clustergrp = merged_train.groupby('cluster').mean()  
....: print ("Clustering variable means by cluster")  
....: print(clustergrp)  
....:
```

```
Clustering variable means by cluster
```

| | index | ALCEVR1 | MAREVER1 | ALCPROBS1 | DEVIANT1 | VIOL1 | \ |
|---------|-------------|-----------|-----------|-----------|-----------|-----------|---|
| cluster | | | | | | | |
| 0 | 3328.930505 | 0.946562 | -0.057109 | -0.058550 | -0.121192 | -0.168252 | |
| 1 | 3312.255193 | 0.661266 | 1.097891 | 0.900592 | 1.106490 | 0.791967 | |
| 2 | 3239.829577 | -1.056455 | -0.474543 | -0.412562 | -0.451110 | -0.264092 | |

| | DEP1 | ESTEEM1 | SCHCONN1 | PARACTV | PARPRES | FAMCONCT | |
|---------|-----------|-----------|-----------|-----------|-----------|-----------|--|
| cluster | | | | | | | |
| 0 | -0.196987 | 0.188367 | 0.131356 | 0.151205 | 0.099521 | 0.226842 | |
| 1 | 0.853587 | -0.655070 | -0.927499 | -0.409328 | -0.473980 | -0.955398 | |
| 2 | -0.292352 | 0.206912 | 0.341590 | 0.091216 | 0.156604 | 0.298288 | |

```
In [14]:
```

```
....:  
....: # FINALLY calculate clustering variable means by cluster  
....: clustergrp = merged_train.groupby('cluster').mean()  
....: print ("Clustering variable means by cluster")  
....: print(clustergrp)  
....:
```

```
Clustering variable means by cluster
```

| | index | ALCEVR1 | MAREVER1 | ALCPROBS1 | DEVIANT1 | VIOL1 | \ |
|---------|-------------|-----------|-----------|-----------|-----------|-----------|---|
| cluster | | | | | | | |
| 0 | 3328.930505 | 0.946562 | -0.057109 | -0.058550 | -0.121192 | -0.168252 | |
| 1 | 3312.255193 | 0.661266 | 1.097891 | 0.900592 | 1.106490 | 0.791967 | |
| 2 | 3239.829577 | -1.056455 | -0.474543 | -0.412562 | -0.451110 | -0.264092 | |

| | DEP1 | ESTEEM1 | SCHCONN1 | PARACTV | PARPRES | FAMCONCT | |
|---------|-----------|-----------|-----------|-----------|-----------|-----------|--|
| cluster | | | | | | | |
| 0 | -0.196987 | 0.188367 | 0.131356 | 0.151205 | 0.099521 | 0.226842 | |
| 1 | 0.853587 | -0.655070 | -0.927499 | -0.409328 | -0.473980 | -0.955398 | |
| 2 | -0.292352 | 0.206912 | 0.341590 | 0.091216 | 0.156604 | 0.298288 | |

```
In [14]:
```

```
123
124 """
125 END multiple steps to merge cluster assignment with clustering variables to examine
126 cluster variable means by cluster
127 """
128
129 # FINALLY calculate clustering variable means by cluster
130 clustergrp = merged_train.groupby('cluster').mean()
131 print ("Clustering variable means by cluster")
132 print(clustergrp)
133
134
135 # validate clusters in training data by examining cluster differences in GPA using ANOVA
136 # first have to merge GPA with clustering variables and cluster assignment data
137 gpa_data=data_clean['GPA1']
138 # split GPA data into train and test sets
139 gpa_train, gpa_test = train_test_split(gpa_data, test_size=.3, random_state=123)
140 gpa_train1=pd.DataFrame(gpa_train)
141 gpa_train1.reset_index(level=0, inplace=True)
142 merged_train_all=pd.merge(gpa_train1, merged_train, on='index')
143 sub1 = merged_train_all[['GPA1', 'cluster']].dropna()
144
145 import statsmodels.formula.api as smf
146 import statsmodels.stats.multicomp as multi
147
148 gpamod = smf.ols(formula='GPA1 ~ C(cluster)', data=sub1).fit()
149 print (gpamod.summary())
150
151 print ('means for GPA by cluster')
```

Finally, let's see how
the clusters differ on GPA.

Object inspector

iPython console

Console 1

packages\pa

return

File "C:\

packages\pa

code, i

File "C:\

```
123
124 """
125 END multiple steps to merge cluster assignment with clustering variables to examine
126 cluster variable means by cluster
127 """
128
129 # FINALLY calculate clustering variable means by cluster
130 clustergrp = merged_train.groupby('cluster').mean()
131 print ("Clustering variable means by cluster")
132 print(clustergrp)
133
134
135 # validate clusters in training data by examining cluster differences in GPA using ANOVA
136 # first have to merge GPA with clustering variables and cluster assignment data
137 gpa_data=data_clean['GPA1']
138 # split GPA data into train and test sets
139 gpa_train, gpa_test = train_test_split(gpa_data, test_size=.3, random_state=123)
140 gpa_train1=pd.DataFrame(gpa_train)
141 gpa_train1.reset_index(level=0, inplace=True)
142 merged_train_all=pd.merge(gpa_train1, merged_train, on='index')
143 sub1 = merged_train_all[['GPA1', 'cluster']].dropna()
144
145 import statsmodels.formula.api as smf
146 import statsmodels.stats.multicomp as multi
147
148 gpamod = smf.ols(formula='GPA1 ~ C(cluster)', data=sub1).fit()
149 print (gpamod.summary())
150
151 print ('means for GPA by cluster')
```

We first have to extract the GPA variable from the original data set,

Object inspector

iPython console

Console 1

packages\pa

return

File "C:\

packages\pa

code, i

File "C:\

```
123
124 """
125 END multiple steps to merge cluster assignment with clustering variables to examine
126 cluster variable means by cluster
127 """
128
129 # FINALLY calculate clustering variable means by cluster
130 clustergrp = merged_train.groupby('cluster').mean()
131 print ("Clustering variable means by cluster")
132 print(clustergrp)
133
134
135 # validate clusters in training data by examining cluster differences in GPA using ANOVA
136 # first have to merge GPA with clustering variables and cluster assignment data
137 gpa_data=data_clean['GPA1']
138 # split GPA data into train and test sets
139 gpa_train, gpa_test = train_test_split(gpa_data, test_size=.3, random_state=123)
140 gpa_train1=pd.DataFrame(gpa_train)
141 gpa_train1.reset_index(level=0, inplace=True)
142 merged_train_all=pd.merge(gpa_train1, merged_train, on='index')
143 sub1 = merged_train_all[['GPA1', 'cluster']].dropna()
144
145 import statsmodels.formula.api as smf
146 import statsmodels.stats.multicomp as multi
147
148 gpamod = smf.ols(formula='GPA1 ~ C(cluster)', data=sub1).fit()
149 print (gpamod.summary())
150
151 print ('means for GPA by cluster')
```

Before merging,
we need to reset the index for

Object inspector
Python console
Console 1
packages\pa
return
File "C:\
packages\pa
code, i
File "C:\

```
123
124 """
125 END multiple steps to merge cluster assignment with clustering variables to examine
126 cluster variable means by cluster
127 """
128
129 # FINALLY calculate clustering variable means by cluster
130 clustergrp = merged_train.groupby('cluster').mean()
131 print ("Clustering variable means by cluster")
132 print(clustergrp)
133
134
135 # validate clusters in training data by examining cluster differences in GPA using ANOVA
136 # first have to merge GPA with clustering variables and cluster assignment data
137 gpa_data=data_clean['GPA1']
138 # split GPA data into train and test sets
139 gpa_train, gpa_test = train_test_split(gpa_data, test_size=.3, random_state=123)
140 gpa_train1=pd.DataFrame(gpa_train)
141 gpa_train1.reset_index(level=0, inplace=True)
142 merged_train_all=pd.merge(gpa_train1, merged_train, on='index')
143 sub1 = merged_train_all[['GPA1', 'cluster']].dropna()
144
145 import statsmodels.formula.api as smf
146 import statsmodels.stats.multicomp as multi
147
148 gpamod = smf.ols(formula='GPA1 ~ C(cluster)', data=sub1).fit()
149 print (gpamod.summary)
150
151 print ('means for GPA by cluster')
```

the GPA variable data frame to create the unique identifier to link the datasets.

Object inspector

Python console

Console 1

packages\pa

return

File "C:\

packages\pa

code, i

File "C:\

```
123
124 """
125 END multiple steps to merge cluster assignment with clustering variables to examine
126 cluster variable means by cluster
127 """
128
129 # FINALLY calculate clustering variable means by cluster
130 clustergrp = merged_train.groupby('cluster').mean()
131 print ("Clustering variable means by cluster")
132 print(clustergrp)
133
134
135 # validate clusters in training data by examining cluster differences in GPA using ANOVA
136 # first have to merge GPA with clustering variables and cluster assignment data
137 gpa_data=data_clean['GPA1']
138 # split GPA data into train and test sets
139 gpa_train, gpa_test = train_test_split(gpa_data, test_size=.3, random_state=123)
140 gpa_train1=pd.DataFrame(gpa_train)
141 gpa_train1.reset_index(level=0, inplace=True)
142 merged_train_all=pd.merge(gpa_train1, merged_train, on='index')
143 sub1 = merged_train_all[['GPA1', 'cluster']].dropna()
144
145 import statsmodels.formula.api as smf
146 import statsmodels.stats.multicomp as multi
147
148 gpamod = smf.ols(formula="GPA1 ~ C(CLUSTER)", data=sub1).fit()
149 print (gpamod.summary())
150
151 print ('means for GPA by cluster')
```

We then merge the datasets by the unique identifier index into a data frame

Object inspector
Python console
Console 1
packages\pa
return
File "C:\
packages\pa
code, i
File "C:\

```
29 # FINALLY calculate clustering variable means by cluster
30 clustergrp = merged_train.groupby('cluster').mean()
31 print ("Clustering variable means by cluster")
32 print(clustergrp)
33
34
35 # validate clusters in training data by examining cluster differences in GPA using ANOVA
36 # first have to merge GPA with clustering variables and cluster assignment data
37 gpa_data=data_clean['GPA1']
38 # split GPA data into train and test sets
39 gpa_train, gpa_test = train_test_split(gpa_data, test_size=.3, random_state=123)
40 gpa_trainl=pd.DataFrame(gpa_train)
41 gpa_trainl.reset_index(level=0, inplace=True)
42 merged_train_all=pd.merge(gpa_trainl, merged_train, on='index')
43 sub1 = merged_train_all[['GPA1', 'cluster']].dropna()
44
45 import statsmodels.formula.api as smf
46 import statsmodels.stats.multicomp as multi
47
48 gpamod = smf.ols(formula='GPA1 ~ C(cluster)', data=sub1).fit()
49 print (gpamod.summary())
50
51 print ('means for GPA by cluster')
52 m1= sub1.groupby('cluster').mean()
53 print (m1)
54
55 print ('standard deviations for GPA by cluster')
```

We'll use analysis of variance to test whether there are significant differences

```
137 gpa_data=data_clean['GPA1']
138 # split GPA data into train and test sets
139 gpa_train, gpa_test = train_test_split(gpa_data, test_size=.3, random_state=123)
140 gpa_train1=pd.DataFrame(gpa_train)
141 gpa_train1.reset_index(level=0, inplace=True)
142 merged_train_all=pd.merge(gpa_train1, merged_train, on='index')
143 sub1 = merged_train_all[['GPA1', 'cluster']].dropna()
144
145 import statsmodels.formula.api as smf
146 import statsmodels.stats.multicomp as multi
147
148 gpamod = smf.ols(formula='GPA1 ~ C(cluster)', data=sub1).fit()
149 print (gpamod.summary())
150
151 print ('means for GPA by cluster')
152 m1= sub1.groupby('cluster').mean()
153 print (m1)
154
155 print ('standard deviations for GPA by cluster')
156 m2= sub1.groupby('cluster').std()
157 print (m2)
158
159 mcl = multi.MultiComparison(sub1['GPA1'], sub1['cluster'])
160 res1 = mcl.tukeyhsd()
161 print(res1.summary())
162
163
```

between clusters on
the quantitative GPA variable.

Object inspect
IPython console
Console
...
...
...
...
...
=====

Dep. Var.:
Model:
Method:
Date:
Time:
No. Observations: 161
Df Residuals: 159
Df Model: 2
Covariance Type: opg
=====

=

Int.]

```
138 # split GPA data into train and test sets
139 gpa_train, gpa_test = train_test_split(gpa_data, test_size=.3, random_state=123)
140 gpa_train1=pd.DataFrame(gpa_train)
141 gpa_train1.reset_index(level=0, inplace=True)
142 merged_train_all=pd.merge(gpa_train1, merged_train, on='index')
143 sub1 = merged_train_all[['GPA1', 'cluster']].dropna()
144
145 import statsmodels.formula.api as smf
146 import statsmodels.stats.multicomp as multi
147
148 gpamod = smf.ols(formula='GPA1 ~ C(cluster)', data=sub1).fit()
149 print (gpamod.summary())
150
151 print ('means for GPA by cluster')
152 m1= sub1.groupby('cluster').mean()
153 print (m1)
154
155 print ('standard deviations for GPA by cluster')
156 m2= sub1.groupby('cluster').std()
157 print (m2)
158
159 mc1 = multi.MultiComparison(sub1['GPA1'], sub1['cluster'])
160 res1 = mc1.tukeyhsd()
161 print(res1.summary())
162
```

We use the ols function to
test the analysis of variance.

```
138 # split GPA data into train and test sets
139 gpa_train, gpa_test = train_test_split(gpa_data, test_size=.3, random_state=123)
140 gpa_train1=pd.DataFrame(gpa_train)
141 gpa_train1.reset_index(level=0, inplace=True)
142 merged_train_all=pd.merge(gpa_train1, merged_train, on='index')
143 sub1 = merged_train_all[['GPA1', 'cluster']].dropna()
144
145 import statsmodels.formula.api as smf
146 import statsmodels.stats.multicomp as multi
147
148 gpamod = smf.ols(formula='GPA1 ~ C(cluster)', data=sub1).fit()
149 print (gpamod.summary())
150
151 print ('means for GPA by cluster')
152 m1= sub1.groupby('cluster').mean()
153 print (m1)
154
155 print ('standard deviations for GPA by cluster')
156 m2= sub1.groupby('cluster').std()
157 print (m2)
158
159 mc1 = multi.MultiComparison(sub1['GPA1'], sub1['cluster'])
160 res1 = mc1.tukey()
161 print(res1.summary())
162
```

The formula specifies the model,
with GPA as the response variable and

```
138 # split GPA data into train and test sets
139 gpa_train, gpa_test = train_test_split(gpa_data, test_size=.3, random_state=123)
140 gpa_train1=pd.DataFrame(gpa_train)
141 gpa_train1.reset_index(level=0, inplace=True)
142 merged_train_all=pd.merge(gpa_train1, merged_train, on='index')
143 sub1 = merged_train_all[['GPA1', 'cluster']].dropna()
144
145 import statsmodels.formula.api as smf
146 import statsmodels.stats.multicomp as multi
147
148 gpamod = smf.ols(formula='GPA1 ~ C(cluster)', data=sub1).fit()
149 print (gpamod.summary())
150
151 print ('means for GPA by cluster')
152 m1= sub1.groupby('cluster').mean()
153 print (m1)
154
155 print ('standard deviations for GPA by cluster')
156 m2= sub1.groupby('cluster').std()
157 print (m2)
158
159 mc1 = multi.MultiComparison(sub1['GPA1'], sub1['cluster'])
160 res1 = mc1.tukeyhsd()
161 print(res1.summary())
162
```

The capital C tells Python that
the cluster assignment variable is

```
138 # split GPA data into train and test sets
139 gpa_train, gpa_test = train_test_split(gpa_data, test_size=.3, random_state=123)
140 gpa_train1=pd.DataFrame(gpa_train)
141 gpa_train1.reset_index(level=0, inplace=True)
142 merged_train_all=pd.merge(gpa_train1, merged_train, on='index')
143 sub1 = merged_train_all[['GPA1', 'cluster']].dropna()
144
145 import statsmodels.formula.api as smf
146 import statsmodels.stats.multicomp as multi
147
148 gpamod = smf.ols(formula='GPA1 ~ C(cluster)', data=sub1).fit()
149 print (gpamod.summary())
150
151 print ('means for GPA by cluster')
152 m1= sub1.groupby('cluster').mean()
153 print (m1)
154
155 print ('standard deviations for GPA by cluster')
156 m2= sub1.groupby('cluster').std()
157 print (m2)
158
159 mc1 = multi.MultiComparison(sub1['GPA1'], sub1['cluster'])
160 res1 = mc1.tukeyhsd()
161 print(res1.summary())
162
```

categorical.

```
138 # split GPA data into train and test sets
139 gpa_train, gpa_test = train_test_split(gpa_data, test_size=.3, random_state=123)
140 gpa_train1=pd.DataFrame(gpa_train)
141 gpa_train1.reset_index(level=0, inplace=True)
142 merged_train_all=pd.merge(gpa_train1, merged_train, on='index')
143 sub1 = merged_train_all[['GPA1', 'cluster']].dropna()
144
145 import statsmodels.formula.api as smf
146 import statsmodels.stats.multicomp as multi
147
148 gpamod = smf.ols(formula='GPA1 ~ C(cluster)', data=sub1).fit()
149 print (gpamod.summary())
150
151 print ('means for GPA by cluster')
152 m1= sub1.groupby('cluster').mean()
153 print (m1)
154
155 print ('standard deviations for GPA by cluster')
156 m2= sub1.groupby('cluster').std()
157 print (m2)
158
159 mcl = multi.MultiComparison(sub1['GPA1'], sub1['cluster'])
160 res1 = mcl.tukeyhsd()
161 print(res1.summary())
162
```

we will request a tukey test to evaluate post hoc comparisons

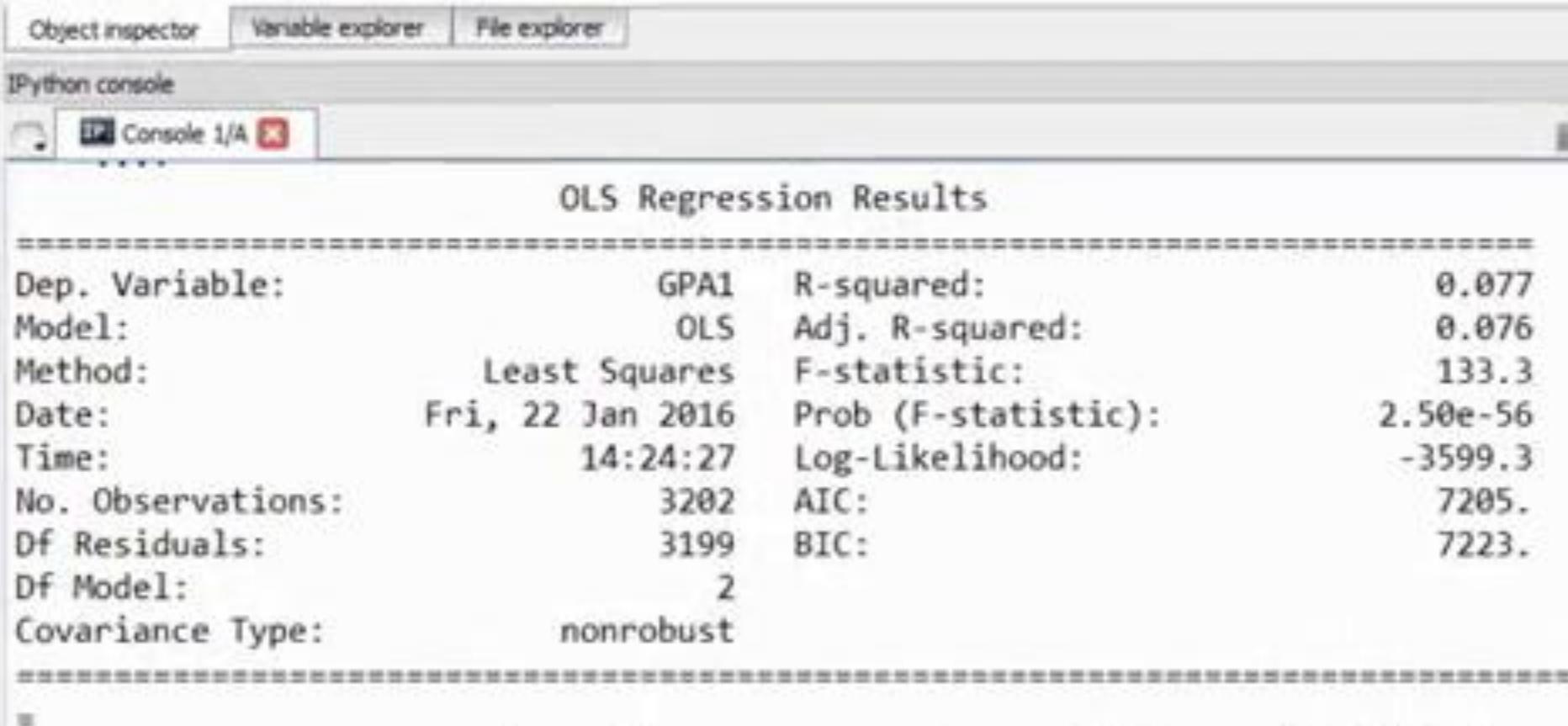
```
138 # split GPA data into train and test sets
139 gpa_train, gpa_test = train_test_split(gpa_data, test_size=.3, random_state=123)
140 gpa_train1=pd.DataFrame(gpa_train)
141 gpa_train1.reset_index(level=0, inplace=True)
142 merged_train_all=pd.merge(gpa_train1, merged_train, on='index')
143 sub1 = merged_train_all[['GPA1', 'cluster']].dropna()
144
145 import statsmodels.formula.api as smf
146 import statsmodels.stats.multicomp as multi
147
148 gpamod = smf.ols(formula='GPA1 ~ C(cluster)', data=sub1).fit()
149 print (gpamod.summary())
150
151 print ('means for GPA by cluster')
152 m1= sub1.groupby('cluster').mean()
153 print (m1)
154
155 print ('standard deviations for GPA by cluster')
156 m2= sub1.groupby('cluster').std()
157 print (m2)
158
159 mcl = multi.MultiComparison(sub1['GPA1'], sub1['cluster'])
160 res1 = mcl.tukeyhsd()
161 print(res1.summary())
162
```

```
138 # split GPA data into train and test sets
139 gpa_train, gpa_test = train_test_split(gpa_data, test_size=.3, random_state=123)
140 gpa_train1=pd.DataFrame(gpa_train)
141 gpa_train1.reset_index(level=0, inplace=True)
142 merged_train_all=pd.merge(gpa_train1, merged_train, on='index')
143 sub1 = merged_train_all[['GPA1', 'cluster']].dropna()
144
145 import statsmodels.formula.api as smf
146 import statsmodels.stats.multicomp as multi
147
148 gpamod = smf.ols(formula='GPA1 ~ C(cluster)', data=sub1).fit()
149 print (gpamod.summary())
150
151 print ('means for GPA by cluster')
152 m1= sub1.groupby('cluster').mean()
153 print (m1)
154
155 print ('standard deviations for GPA by cluster')
156 m2= sub1.groupby('cluster').std()
157 print (m2)
158
159 mcl = multi.MultiComparison(sub1['GPA1'], sub1['cluster'])
160 res1 = mcl.tukeyh()
161 print(res1.summary())
162
```

between the clusters using
the multi comparison function from

ences in GPA using ANOVA
assignment data

random_state=123)



The screenshot shows a Jupyter Notebook interface with the following details:

- Header: Object inspector, Variable explorer, File explorer.
- Section: IPython console, Console 1/A.
- Output:

```
OLS Regression Results
=====
Dep. Variable: GPA1 R-squared: 0.077
Model: OLS Adj. R-squared: 0.076
Method: Least Squares F-statistic: 133.3
Date: Fri, 22 Jan 2016 Prob (F-statistic): 2.50e-56
Time: 14:24:27 Log-Likelihood: -3599.3
No. Observations: 3202 AIC: 7205.
Df Residuals: 3199 BIC: 7223.
Df Model: 2
Covariance Type: nonrobust
```

The analysis of variance summary table

ences in GPA using ANOVA
assignment data

random_state=123)

Object inspector Variable explorer File explorer

IPython console

Console 1/A

OLS Regression Results

| Dep. Variable: | GPA1 | R-squared: | 0.077 |
|-------------------|------------------|---------------------|----------|
| Model: | OLS | Adj. R-squared: | 0.076 |
| Method: | Least Squares | F-statistic: | 133.3 |
| Date: | Fri, 22 Jan 2016 | Prob (F-statistic): | 2.50e-56 |
| Time: | 14:24:27 | Log-Likelihood: | -3599.3 |
| No. Observations: | 3202 | AIC: | 7205. |
| Df Residuals: | 3199 | BIC: | 7223. |
| Df Model: | 2 | | |
| Covariance Type: | nonrobust | | |

indicates that the clusters
differed significantly on GPA.

ences in GPA using ANOVA
ssignment data

random_state=123)

Object inspector Variable explorer File explorer

Jupyter console

Console 1/A 

means for GPA by cluster
GPA1

cluster

0 2.829949

1 2.426063

2 2.994542

standard deviations for GPA by cluster
GPA1

cluster

0 0.727230

1 0.738174

2 0.740505

adolescents in cluster 1, the most
troubled group, had the lowest GPA, .05

ences in GPA using ANOVA
ssignment data

random_state=123)

Object inspector Variable explorer File explorer

IPython console

Console 1/A

means for GPA by cluster

GPA1

| cluster | mean |
|---------|----------|
| 0 | 2.829949 |
| 1 | 2.426063 |
| 2 | 2.994542 |

standard deviations for GPA by cluster

GPA1

| cluster | std |
|---------|----------|
| 0 | 0.727230 |
| 1 | 0.738174 |
| 2 | 0.740000 |

and adolescents in cluster 2, the least
troubled group, had the highest GPA.

cluster differences in GPA using ANOVA
and cluster assignment data

```
size=.3, random_state=123)
```

```
= 'index')
```

```
1).fit()
```

```
r'])
```

Object inspector Variable explorer File explorer

Python console

means for GPA by cluster

GPA1

cluster

| cluster | GPA1 |
|---------|----------|
| 0 | 2.829949 |
| 1 | 2.426063 |
| 2 | 2.994542 |

standard deviations for GPA by cluster

GPA1

cluster

| cluster | GPA1 |
|---------|----------|
| 0 | 0.727230 |
| 1 | 0.786903 |
| 2 | 0.729174 |

Multiple Comparison of Means - Tukey HSD, FWER=0.05

| group1 | group2 | meandiff | lower | upper | reject |
|--------|--------|----------|--------|--------|--------|
| 0 | 2 | 0.1646 | 0.0946 | 0.2346 | True |

The tukey test shows that the clusters
differed significantly in mean GPA,

```
nter differences in GPA using ANOVA  
nd cluster assignment data
```

```
, size=.3, random_state=123)
```

```
= 'index')
```

```
1).fit()
```

```
r'])
```

Object inspector Variable explorer File explorer

Python console

Console 1/A

means for GPA by cluster

GPA1

cluster

| | |
|---|----------|
| 0 | 2.829949 |
| 1 | 2.426063 |
| 2 | 2.994542 |

standard deviations for GPA by cluster

GPA1

cluster

| | |
|---|----------|
| 0 | 0.727230 |
| 1 | 0.786903 |
| 2 | 0.729174 |

Multiple Comparison of Means - Tukey HSD, FWER=0.05

| group1 | group2 | meandiff | lower | upper | reject |
|--------|--------|----------|-------|-------|--------|
|--------|--------|----------|-------|-------|--------|

| | | | | | |
|---|---|---------|---------|---------|------|
| 0 | 1 | -0.4039 | -0.4892 | -0.3186 | True |
| 0 | 2 | 0.1646 | 0.0946 | 0.2346 | True |
| 1 | 2 | 0.5685 | 0.4868 | 0.6502 | True |

cluster differences in GPA using ANOVA
and cluster assignment data

```
size=.3, random_state=123)
```

```
= 'index')
```

```
1).fit()
```

```
r'])
```

Object inspector Variable explorer File explorer

Python console

Console 1/A

means for GPA by cluster

GPA1

cluster

| cluster | GPA1 |
|---------|----------|
| 0 | 2.829949 |
| 1 | 2.426063 |
| 2 | 2.994542 |

standard deviations for GPA by cluster

GPA1

cluster

| cluster | GPA1 |
|---------|----------|
| 0 | 0.727230 |
| 1 | 0.786903 |
| 2 | 0.738174 |

Multiple Comparison of Means - Tukey HSD, FWER=0.05

group1 group2 meandiff lower upper reject

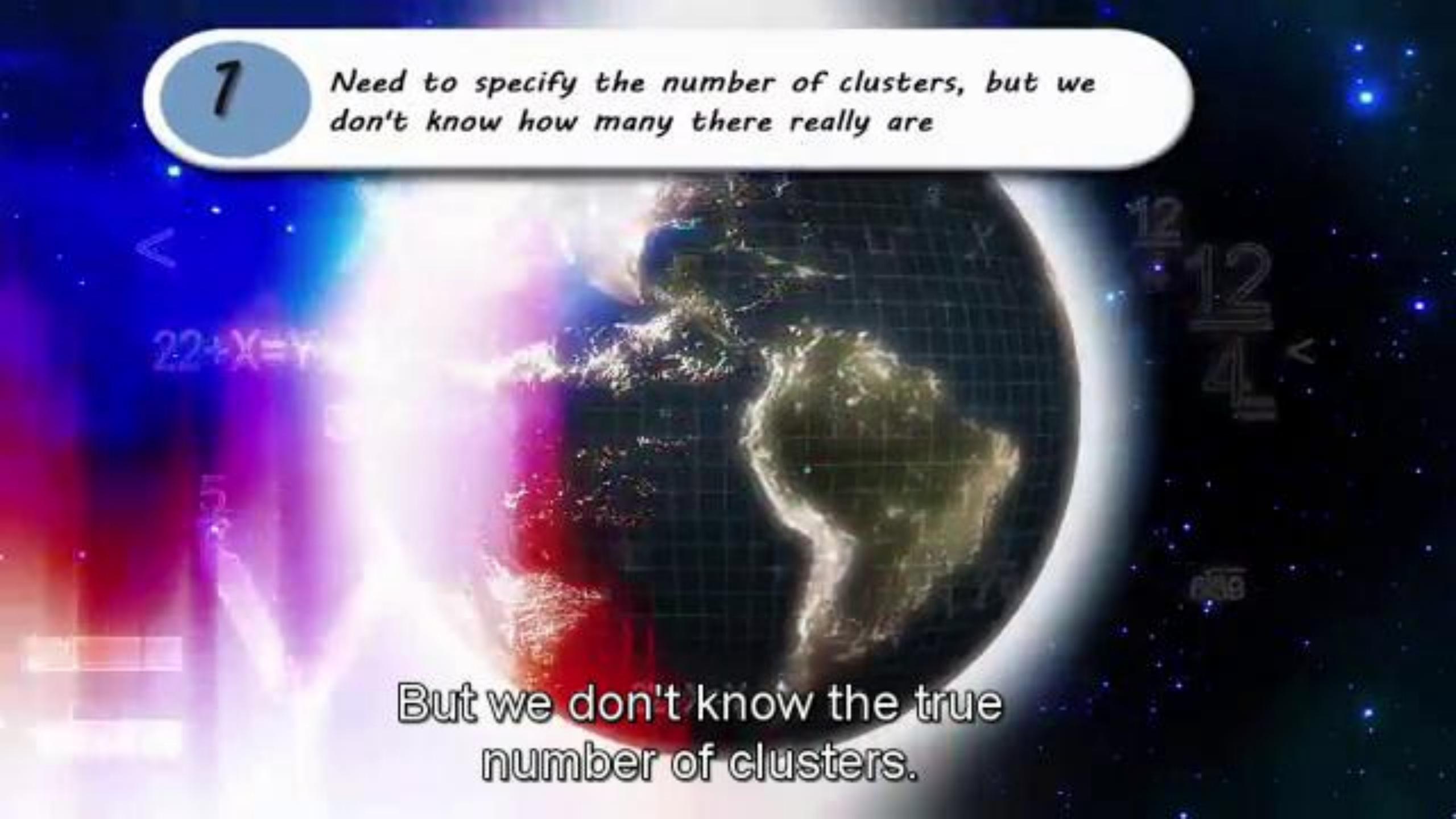
although the difference between
cluster 0 and cluster 2 were smaller.

Module 4

Lesson 4 - k-Means Cluster Analysis Limitations

1

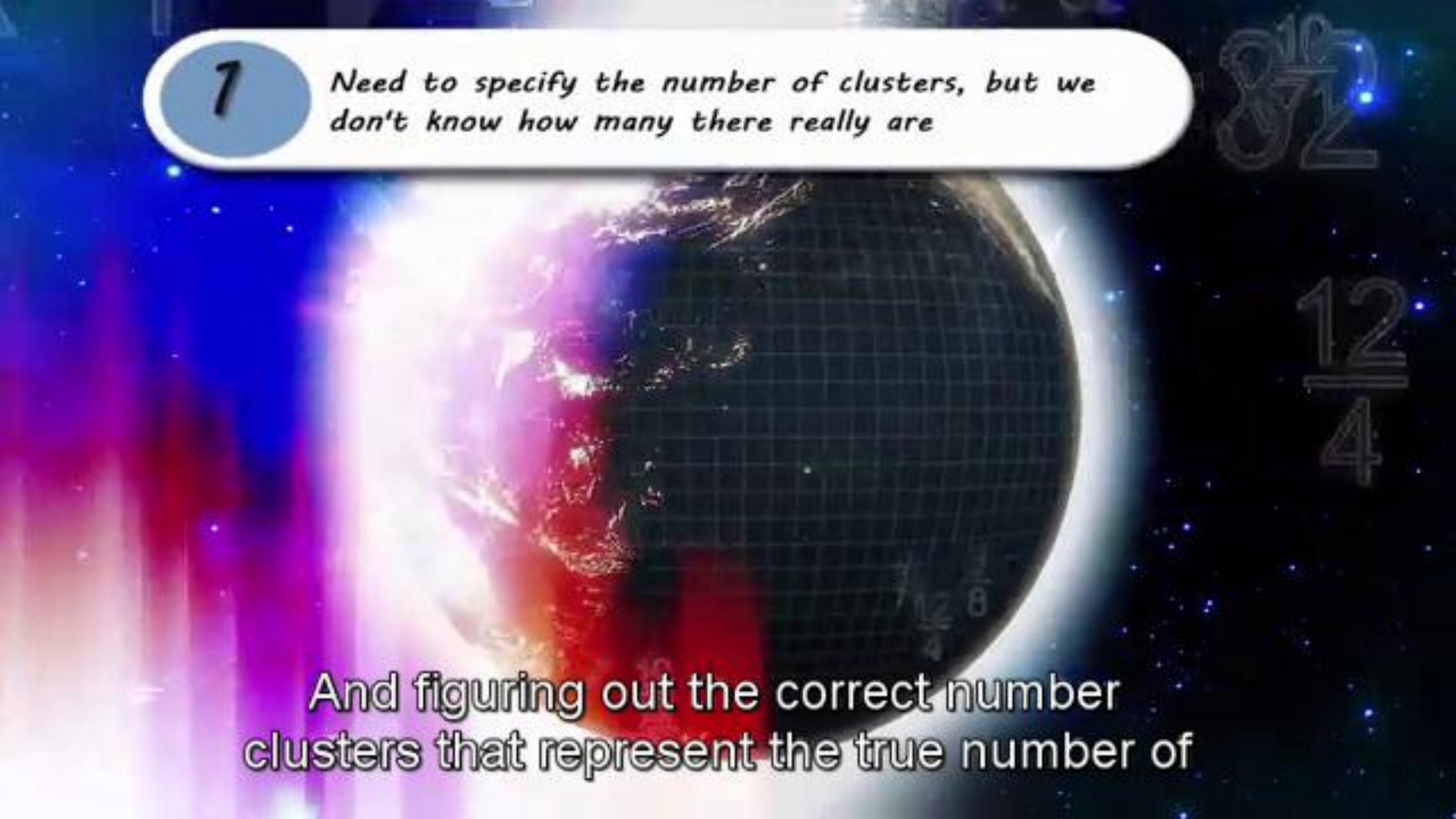
Need to specify the number of clusters, but we
don't know how many there really are



But we don't know the true
number of clusters.

1

Need to specify the number of clusters, but we
don't know how many there really are



And figuring out the correct number
clusters that represent the true number of

1

Need to specify the number of clusters, but we
don't know how many there really are

22+X



clusters in the population
is pretty subjective.

1

Need to specify the number of clusters, but we don't know how many there really are

2

Results can change depending on the location of the initial centroids

On top of that, your results can change depending on the location of

1

Need to specify the number of clusters, but we don't know how many there really are

2

Results can change depending on the location of the initial centroids

3

K-means cluster analysis isn't recommended if you have a lot of categorical variables

**K-means cluster analysis
is not recommended if**

1

Need to specify the number of clusters, but we don't know how many there really are

2

Results can change depending on the location of the initial centroids

3

K-means cluster analysis isn't recommended if you have a lot of categorical variables

If you have a lot of categorical variables, then you need to use

1

Need to specify the number of clusters, but we don't know how many there really are

2

Results can change depending on the location of the initial centroids

3

K-means cluster analysis isn't recommended if you have a lot of categorical variables

4

Assumes that clusters are spherical, distinct, and approximately equal in size

1 Need to specify the number of clusters, but we don't know how many there really are

2 Results can change depending on the location of the initial centroids

3 K-means cluster analysis isn't recommended if you have a lot of categorical variables

4 Assumes that clusters are spherical, distinct, and as a result tends to identify clusters with these characteristics.

1

Need to specify the number of clusters, but we don't know how many there really are

2

Results can change depending on the location of the initial centroids

3

K-means cluster analysis isn't recommended if you have a lot of categorical variables

4

Assumes that clusters are spherical, distinct, and it won't work as well if clusters are elongated or not equal in size.

A photograph of a two-lane asphalt road curving through a hilly, coastal landscape. The road has a double yellow line in the center and white lines on the edges. It's surrounded by green grass and shrubs. In the background, there are more hills and a clear blue sky.

There are a few steps you can take to
help you feel more confident about

A photograph of a two-lane asphalt road curving through a hilly, coastal landscape. The road has a double yellow line in the center and white lines on the edges. It's surrounded by green grass and shrubs, with a large, rocky hillside to the left and a dense forest of tall trees on a hill to the right. A small yellow sign is visible on the right side of the road.

the reliability and
validity of your clusters.

- *Conduct k-means cluster analysis using a range of values of k*

First, conduct the k-means cluster analysis using a range of values of k.

- Conduct k-means cluster analysis using a range of values of k

This helps, but doesn't completely solve the cluster instability problem

- Conduct k-means cluster analysis using a range of values of k
- Split data into training and test sets to run multiple samples through the algorithm

Splitting your data into training and test data sets,

- Conduct k-means cluster analysis using a range of values of k
- Split data into training and test sets to run multiple samples through the algorithm

will allow you to run more than one sample through your algorithm, and

- Conduct k-means cluster analysis using a range of values of k
- Split data into training and test sets to run multiple samples through the algorithm

can be helpful in determining whether the clusters you find are reliable.

- Conduct k-means cluster analysis using a range of values of k
- Split data into training and test sets to run multiple samples through the algorithm

If you get the same results in different samples, you can be more confident that

- Conduct k-means cluster analysis using a range of values of k
- Split data into training and test sets to run multiple samples through the algorithm

the clusters are reasonably catching the underlying subgroups in your population.

- Conduct k-means cluster analysis using a range of values of k
- Split data into training and test sets to run multiple samples through the algorithm
- Validate the clusters

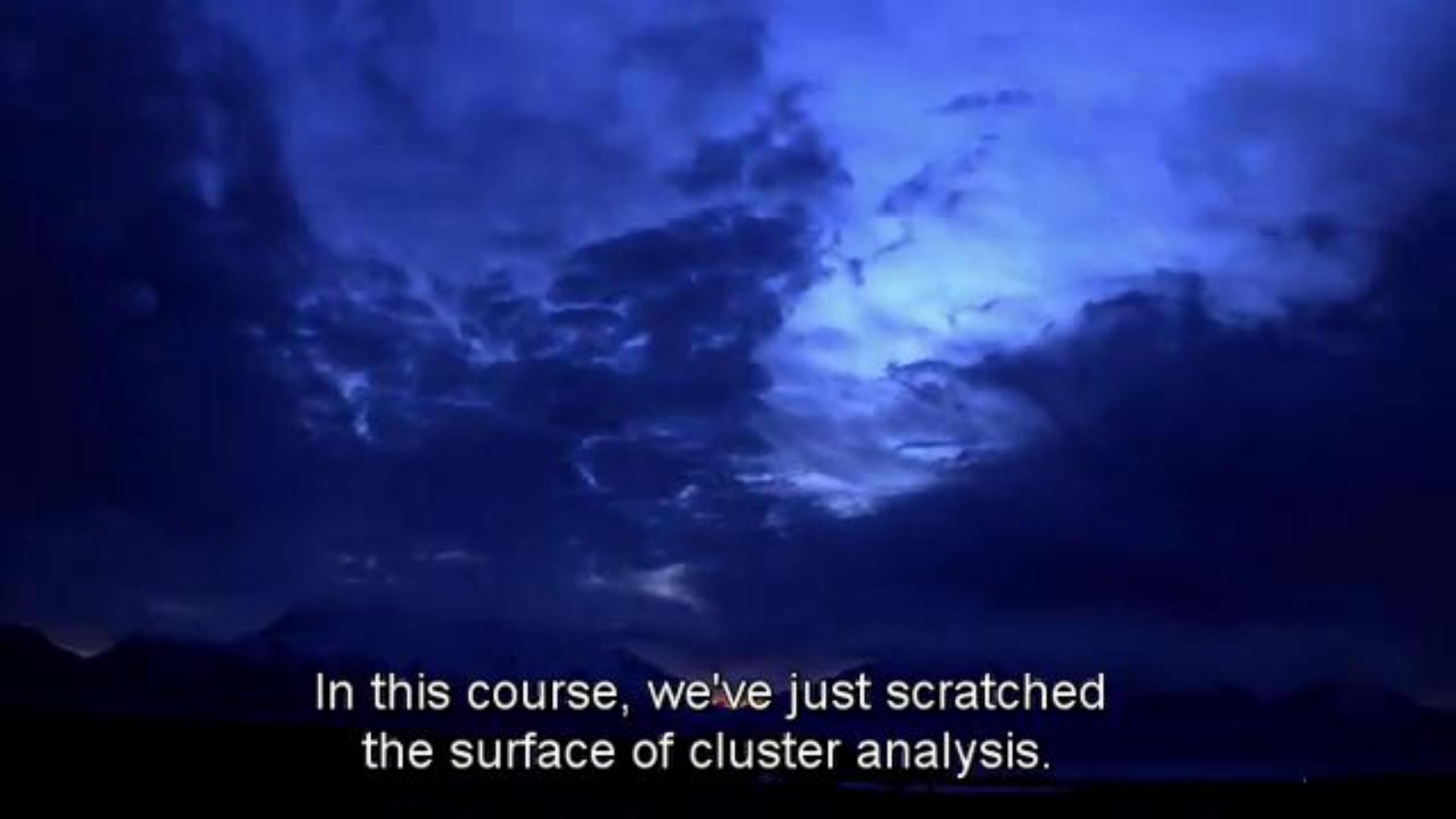
validating the clusters by determining whether they are interpretable, and

- Conduct k-means cluster analysis using a range of values of k
- Split data into training and test sets to run multiple samples through the algorithm
- Validate the clusters

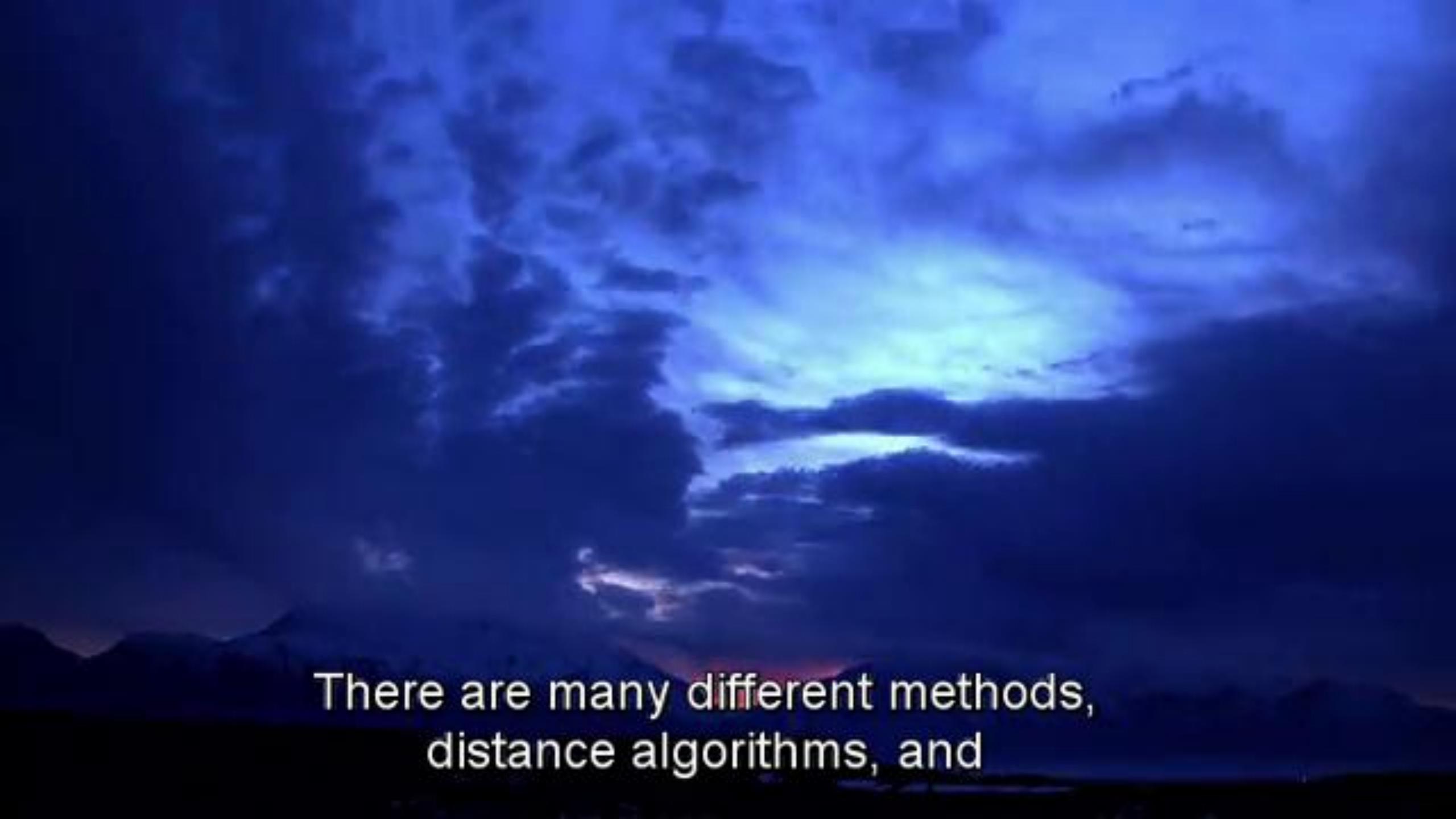
whether they differ from each other on other variables not used in the cluster

- Conduct k-means cluster analysis using a range of values of k
- Split data into training and test sets to run multiple samples through the algorithm
- Validate the clusters

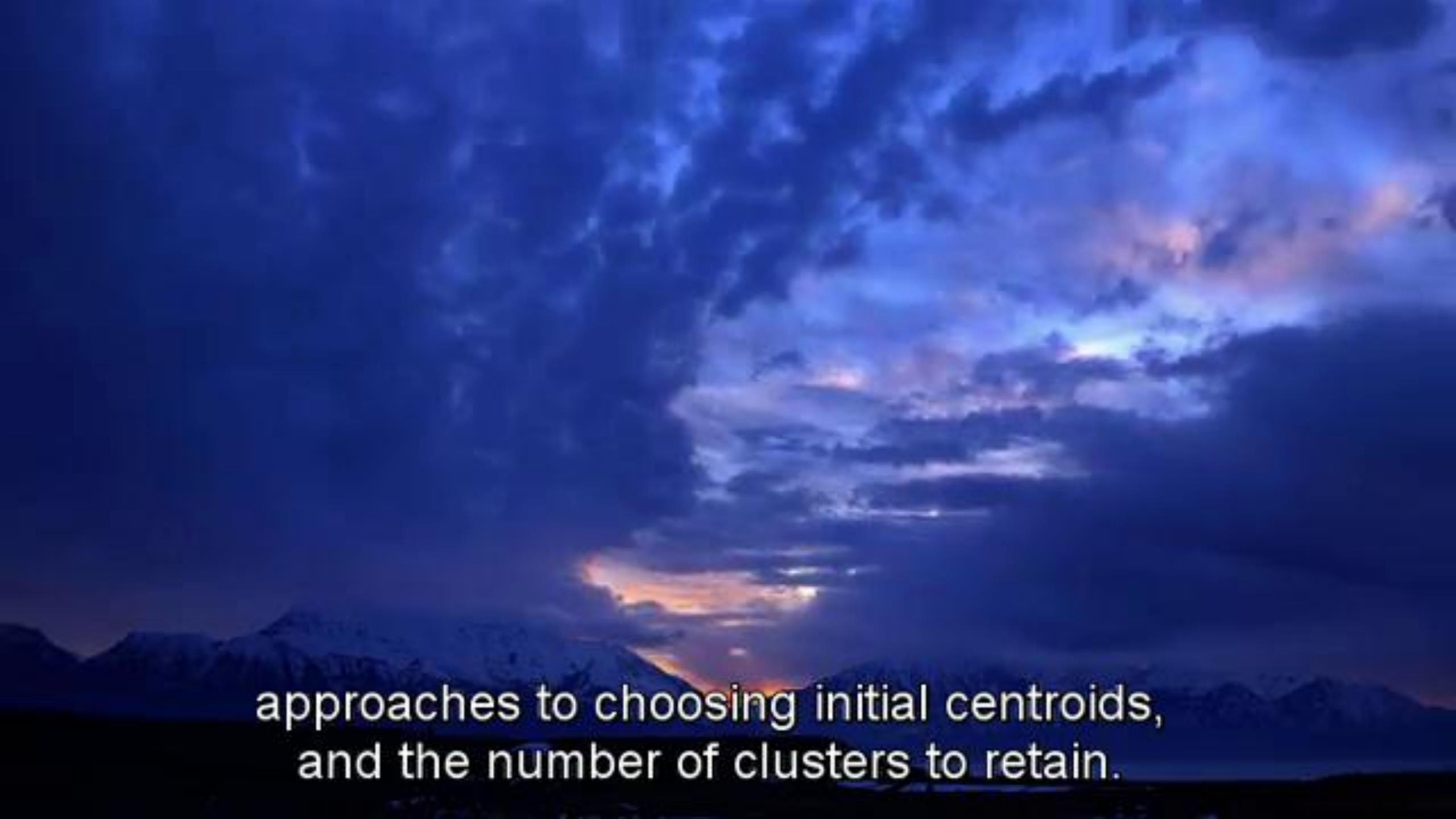
analysis, can increase your confidence in the cluster solution that you choose.



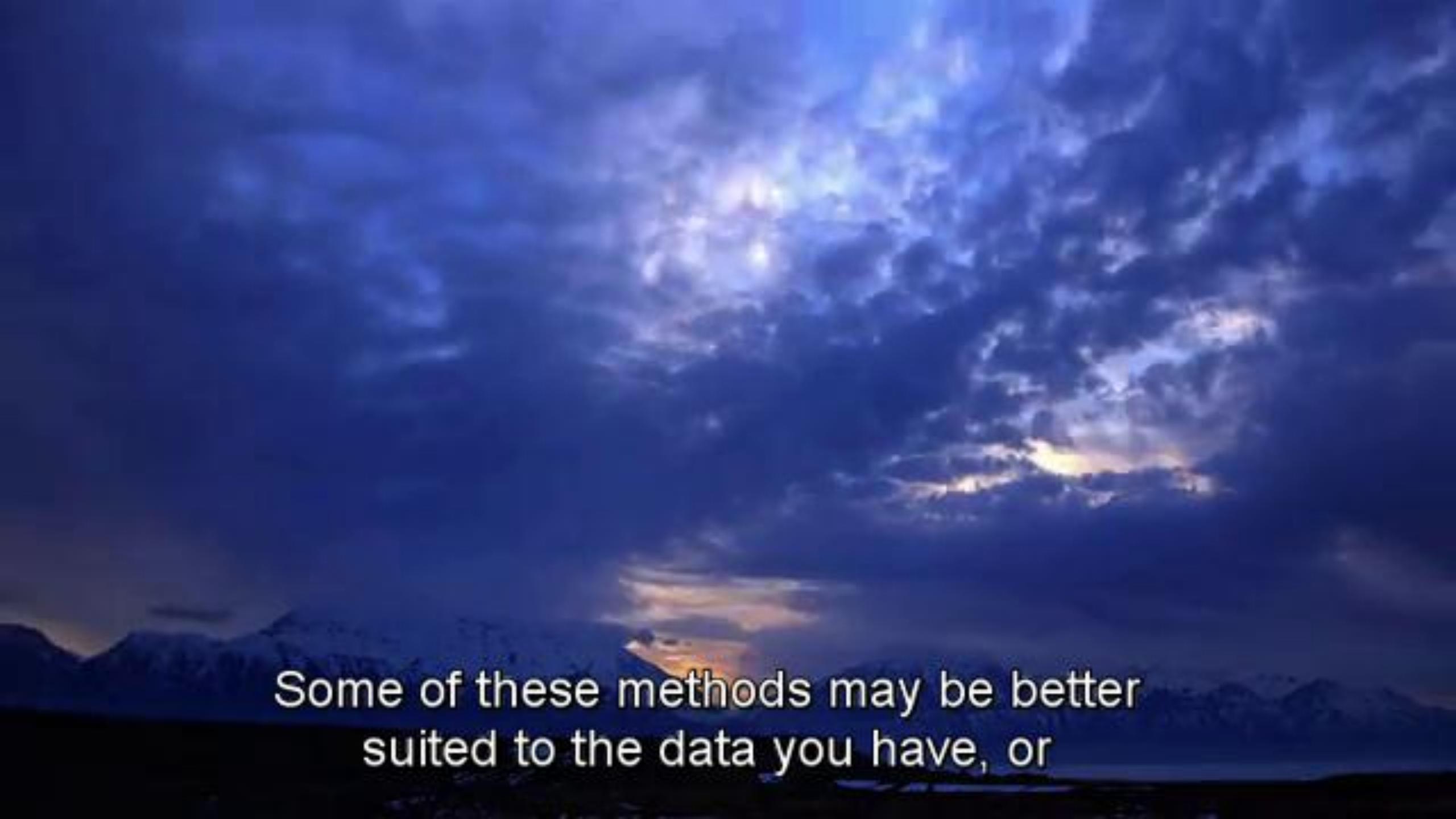
In this course, we've just scratched
the surface of cluster analysis.



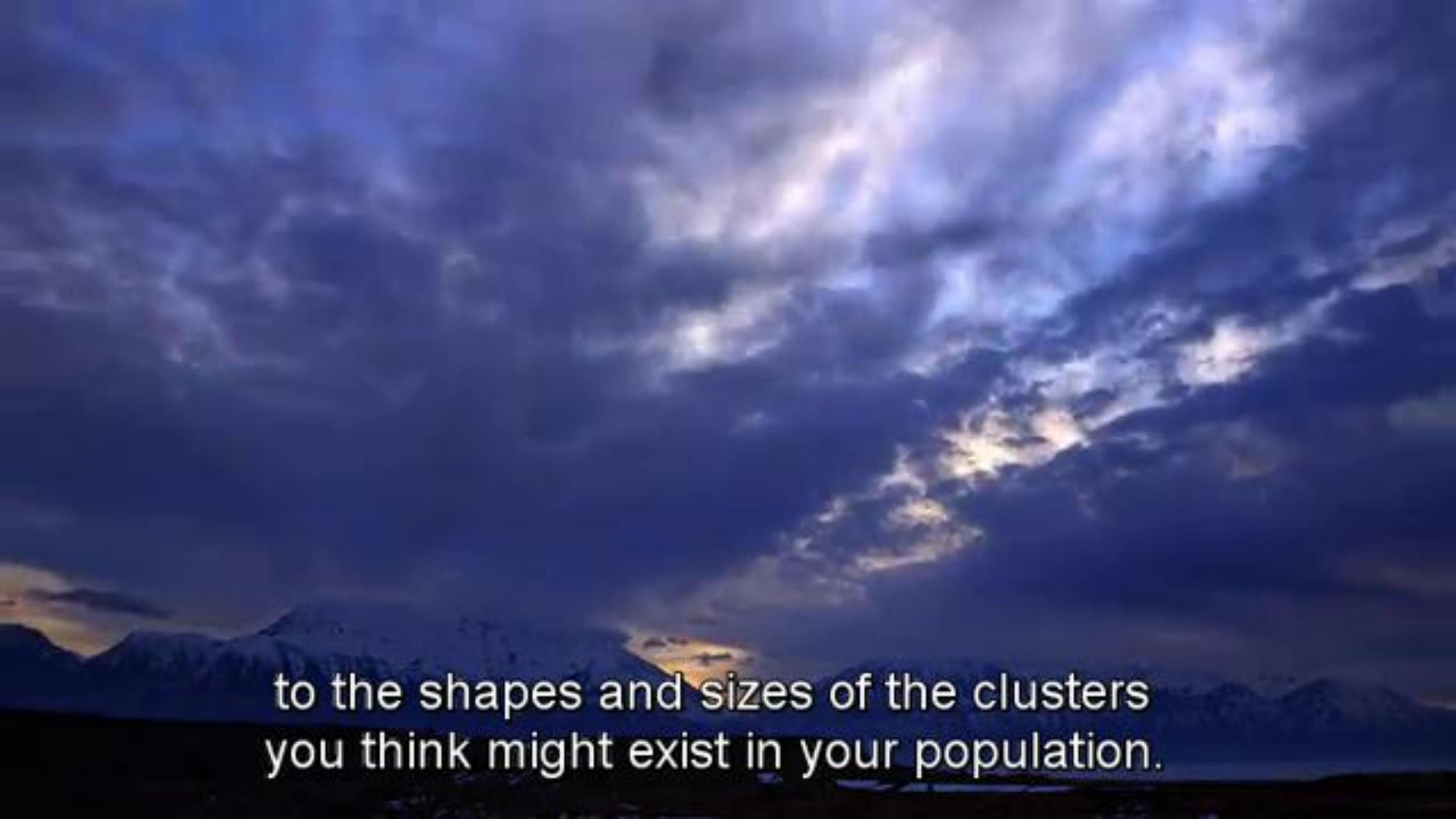
There are many different methods,
distance algorithms, and

The background of the slide features a wide-angle photograph of a mountainous landscape at dusk or dawn. The sky is filled with large, billowing clouds that are illuminated from below by the setting or rising sun, creating a vibrant palette of orange, yellow, and purple hues. The mountains in the foreground are dark and silhouetted against the bright sky.

approaches to choosing initial centroids,
and the number of clusters to retain.

The background of the slide features a wide-angle photograph of a natural landscape at dusk or dawn. The sky is filled with large, billowing clouds in shades of deep blue, purple, and hints of orange and yellow where the sun is low on the horizon. Below the clouds, a dark silhouette of mountains or hills is visible against the lighter sky. The overall mood is serene and vast.

Some of these methods may be better
suited to the data you have, or



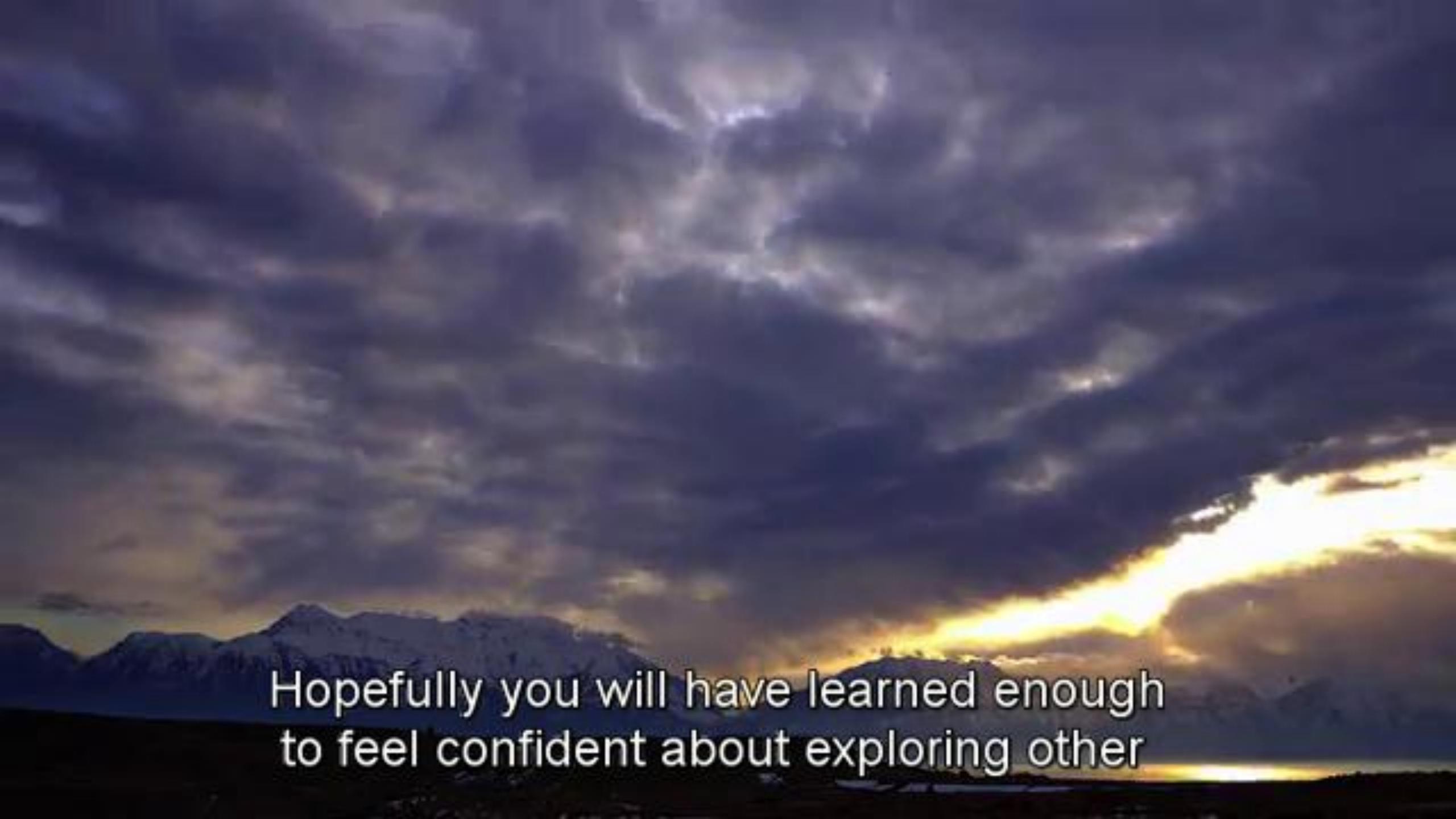
to the shapes and sizes of the clusters
you think might exist in your population.



K-means cluster analysis
is a good starting point

A wide-angle photograph of a dramatic sky at dusk or dawn. The upper half of the image is filled with heavy, dark, textured clouds. In the lower half, a bright, glowing horizon line cuts across the frame, with the light reflecting off the water below. The colors range from deep blues and purples to bright yellows and oranges near the horizon.

because its simplicity makes it
easier to convey the concepts.



Hopefully you will have learned enough
to feel confident about exploring other

methods.