

K Means Clustering for Imagery Analysis



Sajjad Salaria

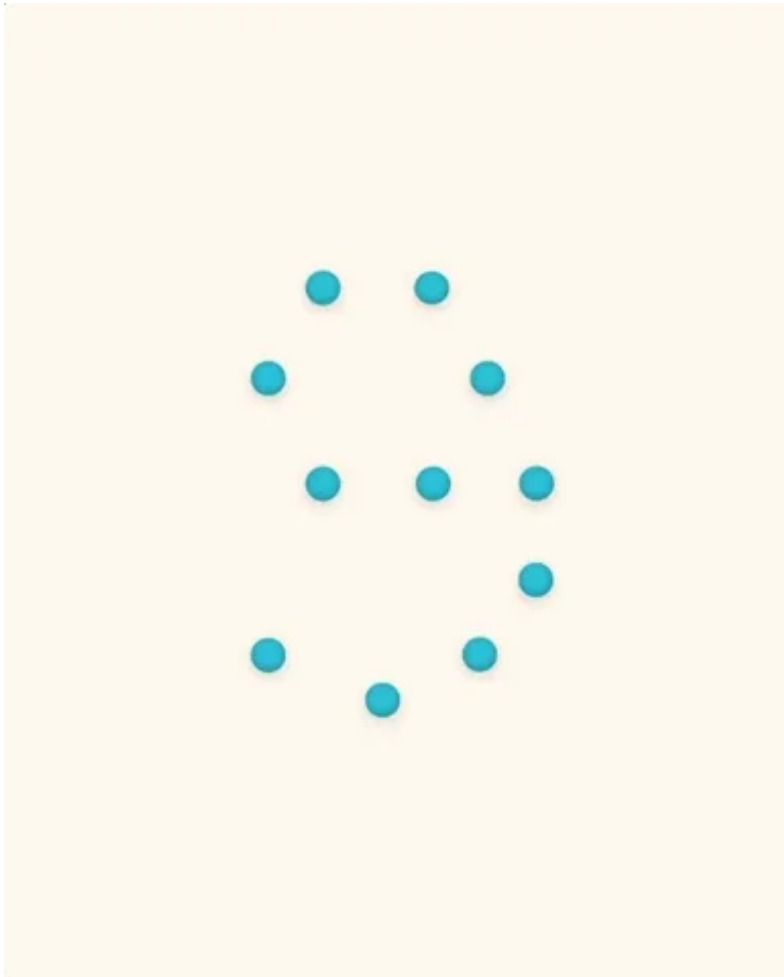
Follow

Aug 22, 2019 · 4 min read

Let's learn about K-Means by doing a mini-project.

In this project, we will use a K-means algorithm to perform image classification.

Clustering isn't limited to the consumer information and population sciences, it can be used for imagery analysis as well. Leveraging Scikit-learn and the MNIST dataset, we will investigate the use of K-means clustering for computer vision.



<https://giphy.com/gifs/l0HlNBe9Z3Xz9x1Yc/html5>

In this project, we will learn how to:

DDI Editor's Pick: 5 Machine Learning Books That Turn You from Novice to Expert | Data Driven...

The booming growth in the Machine Learning industry has brought renewed interest in people about Artificial...

www.datadriveninvestor.com

- Preprocess images for clustering
- Deploy K-means clustering algorithms
- Use common metrics to evaluate cluster performance
- Visualize high-dimensional cluster centroids

Let's get started by importing a few of the libraries we will use in this project.

```
1  import sys
2  import sklearn
3  import matplotlib
4  import numpy as np
5  import matplotlib.pyplot as plt
```

```
6 %matplotlib inline
7
8 print('Python: {}'.format(sys.version))
9 print('Sklearn: {}'.format(sklearn.version))
10 print('Matplotlib: {}'.format(matplotlib.version))
11 print('NumPy: {}'.format(np.version))
```

import.py hosted with ❤ by GitHub

[view raw](#)

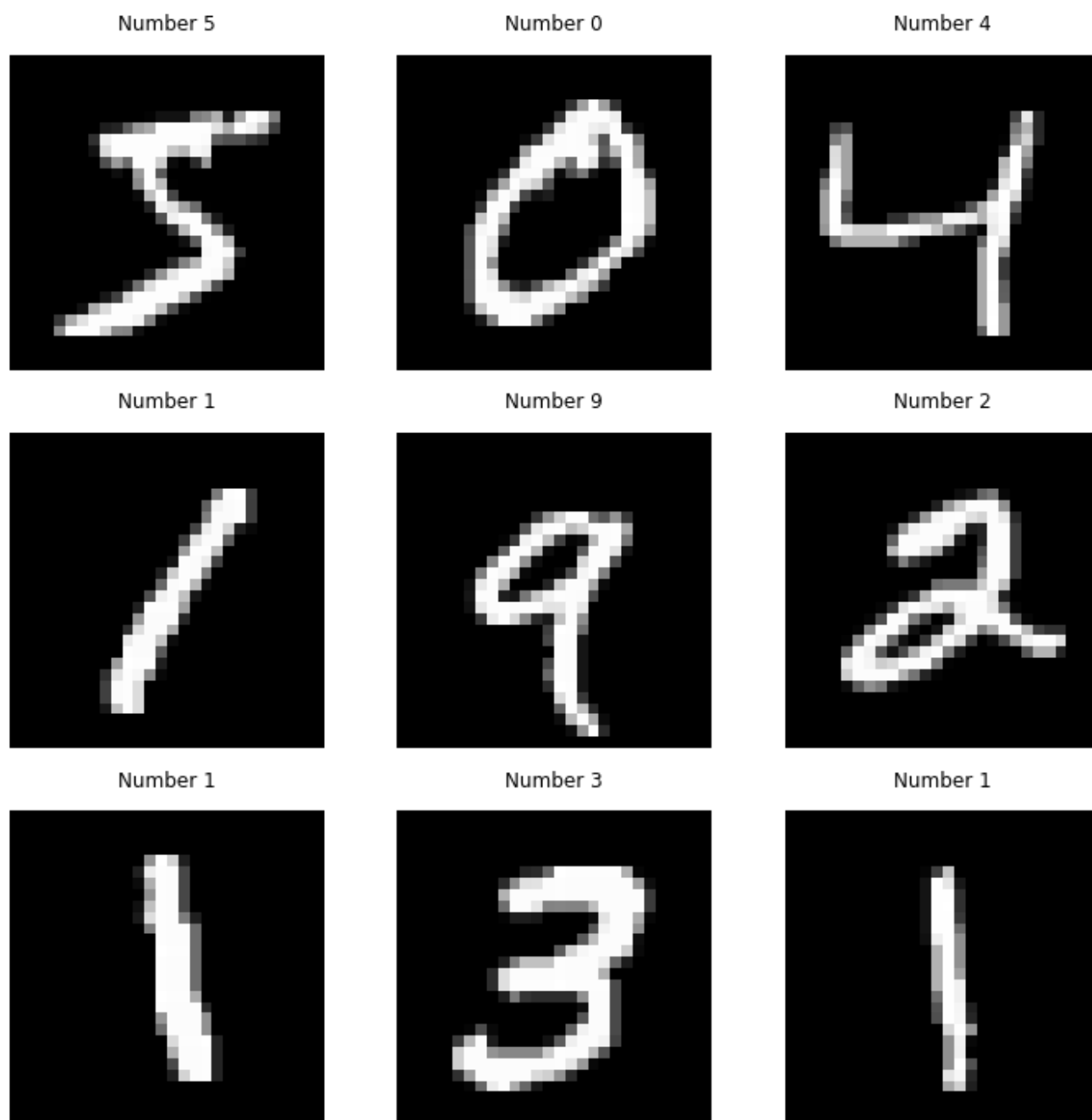
1. Import the MNIST dataset

For this project, we will be using the MNIST dataset. It is available through keras, a deep learning library we have used in previous tutorials. Although we won't be using other features of keras today, it will save us time to import mnist from this library. It is also available through the tensorflow library or for download at <http://yann.lecun.com/exdb/mnist/>.

```
1 from keras.datasets import mnist
2
3 (x_train, y_train), (x_test, y_test) = mnist.load_data()
4
5 print('Training Data: {}'.format(x_train.shape))
6 print('Training Labels: {}'.format(y_train.shape))
7
8 Training Data: (60000L, 28L, 28L)
9 Training Labels: (60000L,)
10
11 print('Testing Data: {}'.format(x_test.shape))
12 print('Testing Labels: {}'.format(y_test.shape))
13
14 Testing Data: (10000L, 28L, 28L)
15 Testing Labels: (10000L,)
16
17 # EDA
18
19 fig, axs = plt.subplots(3, 3, figsize = (12, 12))
20 plt.gray()
21
22 for i, ax in enumerate(axs.flat):
23     ax.matshow(x_train[i])
24     ax.axis('off')
25     ax.set_title('Number {}'.format(y_train[i]))
26
27 # Save the figure
```

```
fig.show()
```

Dataset.py hosted with ❤ by GitHub

[view raw](#)

2. Preprocessing the MNIST images

Images stored as NumPy arrays are 2-dimensional arrays. However, the K-means clustering algorithm provided by scikit-learn ingests 1-dimensional arrays; as a result, we will need to reshape each image.

Clustering algorithms almost always use 1-dimensional data. For example, if you were clustering a set of X, Y coordinates, each point would be passed to the clustering algorithm as a 1-dimensional array with a length of two (example: [2,4] or [-1, 4]). If you were using 3-dimensional data, the array would have a length of 3 (example: [2, 4, 1] or [-1, 4, 5]).

MNIST contains images that are 28 by 28 pixels; as a result, they will have a length of 784 once we reshape them into a 1-dimensional array.

```
1  # convert each image to 1 dimensional array
2
3  X = x_train.reshape(len(x_train),-1)
4  Y = y_train
5
6  # normalize the data to 0 - 1
7
8  X = X.astype(float) / 255.
9
10 print(X.shape)
11 print(X[0].shape)
12
13 (60000L, 784L)
14 (784L,)
```

Pre-Processing.py hosted with ❤ by GitHub

[view raw](#)

3. K-Means Clustering

Time to start clustering! Due to the size of the MNIST dataset, we will use the mini-batch implementation of k-means clustering provided by scikit-learn. This will dramatically reduce the amount of time it takes to fit the algorithm to the data.

The MNIST dataset contains images of the integers 0 to 9. Because of this, let's start by setting the number of clusters to 10, one for each digit.

```
1  from sklearn.cluster import MiniBatchKMeans
2
3  n_digits = len(np.unique(y_test))
4  print(n_digits)
5
6  # Initialize KMeans model
```

```
7
8 kmeans = MiniBatchKMeans(n_clusters = n_digits)
9
10 # Fit the model to the training data
11
12 kmeans.fit(X)
13
14 kmeans.labels_
```

K-means.py hosted with ❤ by GitHub

[view raw](#)

4. Assigning Cluster Labels

K-means clustering is an unsupervised machine learning method; consequently, the labels assigned by our KMeans algorithm refer to the cluster each array was assigned to, not the actual target integer. To fix this, let's define a few functions that will predict which integer corresponds to each cluster.

```
1 def infer_cluster_labels(kmeans, actual_labels):
2     inferred_labels = {}
3
4     for i in range(kmeans.n_clusters):
5
6         # find index of points in cluster
7         labels = []
8         index = np.where(kmeans.labels_ == i)
9
10        # append actual labels for each point in cluster
11        labels.append(actual_labels[index])
12
13        # determine most common label
14        if len(labels[0]) == 1:
15            counts = np.bincount(labels[0])
16        else:
17            counts = np.bincount(np.squeeze(labels))
18
19        # assign the cluster to a value in the inferred_labels dictionary
20        if np.argmax(counts) in inferred_labels:
21            # append the new number to the existing array at this slot
22            inferred_labels[np.argmax(counts)].append(i)
23        else:
24            # create a new array in this slot
25            inferred_labels[np.argmax(counts)] = [i]
```

```

26
27     #print(labels)
28     #print('Cluster: {}, label: {}'.format(i, np.argmax(counts)))
29
30     return inferred_labels
31
32     def infer_data_labels(X_labels, cluster_labels):
33         # empty array of len(X)
34         predicted_labels = np.zeros(len(X_labels)).astype(np.uint8)
35
36         for i, cluster in enumerate(X_labels):
37             for key, value in cluster_labels.items():
38                 if cluster in value:
39                     predicted_labels[i] = key
40
41         return predicted_labels
42
43     # test the infer_cluster_labels() and infer_data_labels() functions
44
45     cluster_labels = infer_cluster_labels(kmeans, Y)
46     X_clusters = kmeans.predict(X)
47     predicted_labels = infer_data_labels(X_clusters, cluster_labels)
48     print(predicted_labels[:20])
49     print(Y[:20])
50
51
52

```

5. Optimizing and Evaluating the Clustering Algorithm

With the functions defined above, we can now determine the accuracy of our algorithms. Since we are using this clustering algorithm for classification, accuracy is ultimately the most important metric; however, there are other metrics out there that can be applied directly to the clusters themselves, regardless of the associated labels. Two of these metrics that we will use are inertia and homogeneity.

```

1  # Initialize and fit KMeans algorithm
2  kmeans = MiniBatchKMeans(n_clusters = 36)
3  kmeans.fit(X)
4
5  # record centroid values
6  centroids = kmeans.cluster_centers_

```

```
7
8 # reshape centroids into images
9 images = centroids.reshape(36, 28, 28)
10 images *= 255
11 images = images.astype(np.uint8)
12
13 # determine cluster labels
14 cluster_labels = infer_cluster_labels(kmeans, Y)
15
16 # create figure with subplots using matplotlib.pyplot
17 fig, axs = plt.subplots(6, 6, figsize = (20, 20))
18 plt.gray()
19
20 # loop through subplots and add centroid images
21 for i, ax in enumerate(axs.flat):
22
23     # determine inferred label using cluster_labels dictionary
24     for key, value in cluster_labels.items():
25         if i in value:
26             ax.set_title('Inferred Label: {}'.format(key))
27
28     # add image to subplot
29     ax.matshow(images[i])
30     ax.axis('off')
31
32 # display the figure
33 fig.show()
```

ClusteringAlgorithmVI.py hosted with ❤ by GitHub

[view raw](#)

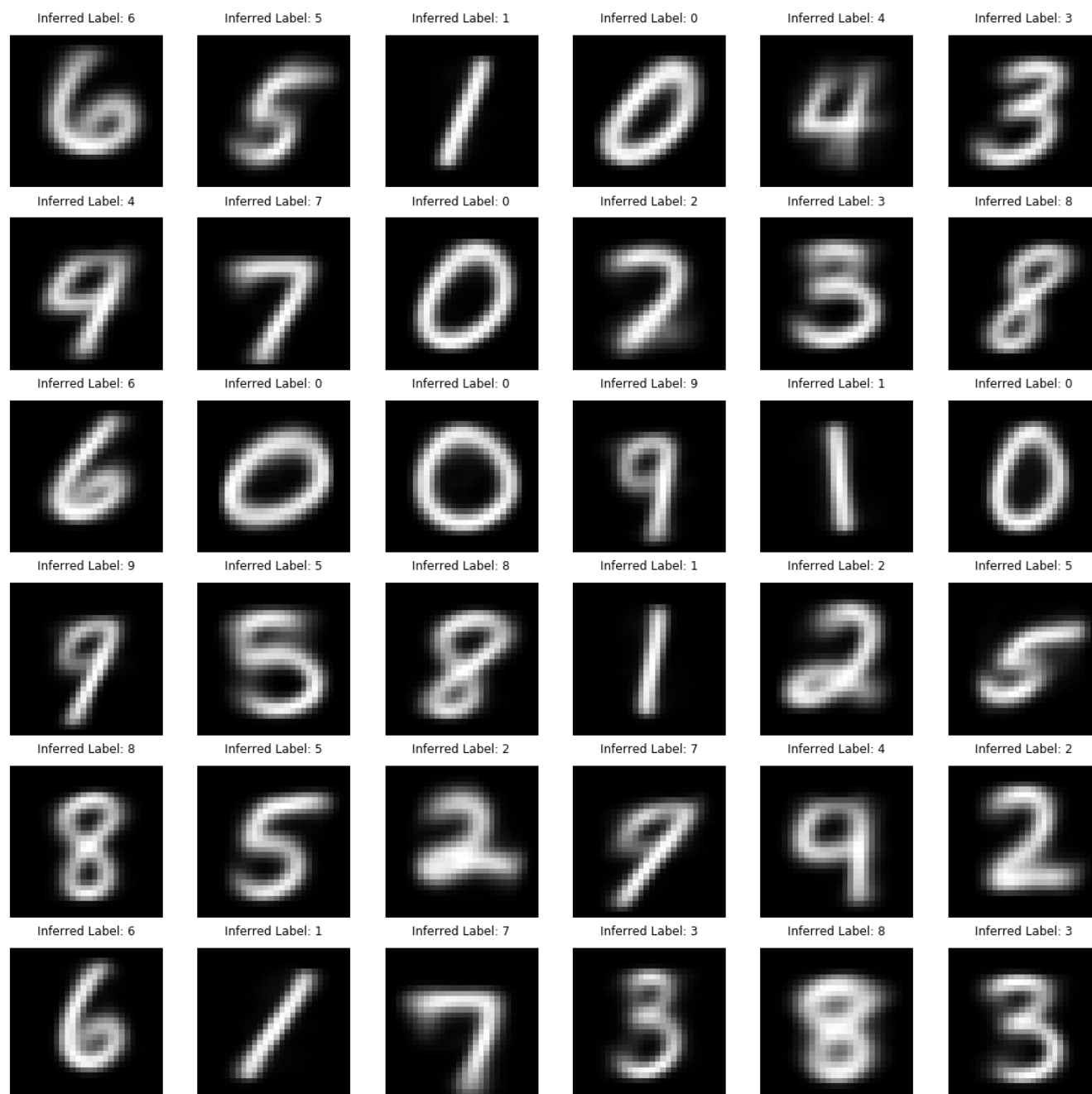
Furthermore, earlier we made the assumption that $K = 10$ was the appropriate number of clusters; however, this might not be the case. Let's fit the K-means clustering algorithm with several different values of K , then evaluate the performance using our metrics.

6. Visualizing Cluster Centroids

The most representative point within each cluster is called the centroid. If we were dealing with X, Y points, the centroid would simply be a point on the graph. However, since we are using arrays of length 784, our centroid is also going to be an array of length 784. We can reshape this array back into a 28 by 28-pixel image and plot it.

These graphs will display the most representative image for each cluster.

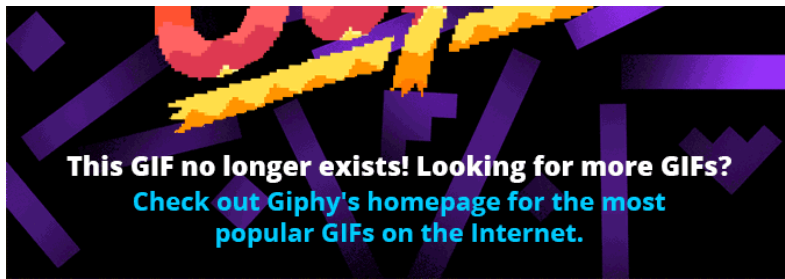

```
print('Accuracy: {}'.format(metrics.accuracy_score(y_test, predicted_labels)))
```



Predictions — Accuracy: 0.9023

For Code: refer -> <https://github.com/xoraus/K-Means-Clustering-for-Imagery-Analysis>





<https://giphy.com/gifs/ncK1aJlwmpHNu/html5>

References:

[1]: https://en.wikipedia.org/wiki/K-means_clustering

[2]: <http://yann.lecun.com/exdb/mnist/>

[3]: https://en.wikipedia.org/wiki/MNIST_database

[4]: <https://www.kaggle.com/ngbolin/mnist-dataset-digit-recognizer>

Gain Access to Expert Views

Give me access!



I agree to leave Medium.com and submit this information, which will be collected and used according to Unsubscribe's privacy policy

Sign up for DDI Highlights

By Data Driven Investor

In each issue we cover all things awesome in the markets, economy, crypto, tech, and more! [Take a look](#)

Get this newsletter

Emails will be sent to lfulton159@gmail.com.
[Not you?](#)

[Machine Learning](#)

[Computer Science](#)

[Deep learning](#)

[Digits Recognition](#)

[K Means Clustering](#)

[About](#) [Help](#) [Legal](#)

Get the Medium app

