

[Live Webinar] Gain a Competitive Edge with Graph Analytics

Sign Up Now ▶

DZone > AI Zone > 10 Interesting Use Cases for the K-Means Algorithm

10 Interesting Use Cases for the K-Means Algorithm

by Kaushik Raghupathi · Mar. 27, 18 · AI Zone · Analysis



We at DZone want to help developers be better developers, but we need to know more about you to do that. Please take our 8-minute Community Survey and you could be one of two people to win \$250. [Take the survey](#) ▶



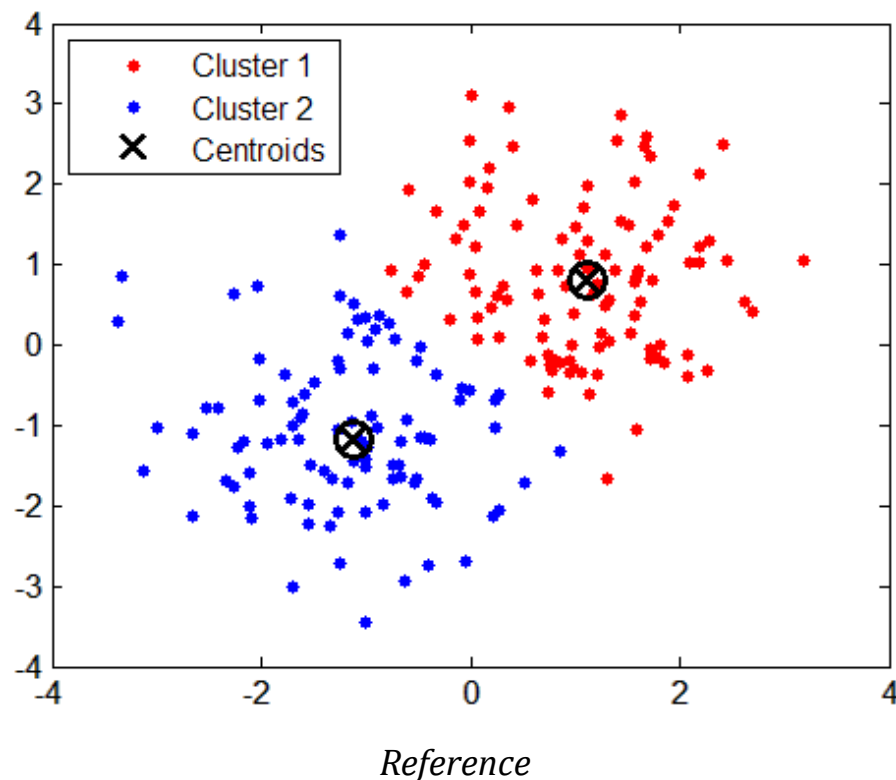
The k-means algorithm is one of the oldest and most commonly used clustering algorithms. It is a great starting point for new ML enthusiasts to pick up, given the simplicity of its implementation. As part of this post, we will review the origins of this algorithm and typical usage scenarios.

The History

The term "k-means" was first used by James MacQueen in 1967 as part of his paper on "Some methods for classification and analysis of multivariate observations". The standard algorithm was also used in Bell Labs as part of a technique in pulse code modulation in 1957. It was also published by In 1965 by E. W. Forgy and typically is also known as the Lloyd-Forgy method.

What Is K-Means?

Clustering is the task of dividing the population or data points into a number of groups such that data points in the same groups are more similar to other data points in the same group than those in other groups. In simple words, the aim is to segregate groups with similar traits and assign them into clusters. The goal of the k-means algorithm is to find groups in the data, with the number of groups represented by the variable K. The algorithm works iteratively to assign each data point to one of K groups based on the features that are provided. In the reference image below, $K=2$, and there are two clusters identified from the source dataset.



The outputs of executing a k-means on a dataset are:

- K centroids: Centroids for each of the k clusters identified from the dataset.
- Complete dataset labeled to ensure each data point is assigned to one of the clusters.

Where Can I Apply K-Means?

k-means can typically be applied to data that has a smaller number of dimensions, is numeric, and is continuous. Think of a scenario in which you want to make groups of similar things from a randomly distributed collection of things; k-means is very suitable for such scenarios.

Here is a list of ten interesting use cases for k-means.

1. Document Classification

Cluster documents in multiple categories based on tags, topics, and the content of the document. This is a very standard classification problem and k-means is a highly suitable algorithm for this purpose. The initial processing of the documents is needed to represent each document as a vector and uses term frequency to identify commonly used terms that help classify the document. The document vectors are then clustered to help identify similarity in document groups. Here is a sample implementation of the k-means for document clustering.

2. Delivery Store Optimization

Optimize the process of good delivery using truck drones by using a combination of k-means to find the optimal number of launch locations and a genetic algorithm to solve the truck route as a traveling salesman problem. Here is a whitepaper on the same topic.

3. Identifying Crime Localities

With data related to crimes available in specific localities in a city, the category of crime, the area of the crime, and the association between the two can give quality insight into crime-prone areas within a city or a locality. Here is an interesting paper based on crime data from Delhi FIRs.

4. Customer Segmentation

Clustering helps marketers improve their customer base, work on target areas, and segment customers based on purchase history, interests, or activity monitoring. Here is a white paper on how telecom providers can cluster pre-paid customers to identify patterns in terms of money spent in recharging, sending SMS, and browsing the internet. The classification would help the company target specific clusters of customers for specific campaigns.

5. Fantasy League Stat Analysis

Analyzing player stats has always been a critical element of the sporting world, and with increasing competition, machine learning has a critical role to play here. As an interesting exercise, if you would like to create a fantasy draft team and like to identify similar players based on player stats, k-means can be a useful option. Check out this article for details and a sample implementation.

6. Insurance Fraud Detection

Machine learning has a critical role to play in fraud detection and has numerous applications in automobile, healthcare, and insurance fraud detection. Utilizing past historical data on fraudulent claims, it is possible to isolate new claims based on its proximity to clusters that indicate fraudulent patterns. Since insurance fraud can potentially have a multi-million dollar impact on a company, the ability to detect frauds is crucial. Check out this white paper on using clustering in automobile insurance to detect frauds.

7. Rideshare Data Analysis

The publicly available Uber ride information dataset provides a large amount of valuable data around traffic, transit time, peak pickup localities, and more. Analyzing this data is useful not just in the context of Uber but also in providing insight into urban traffic patterns and helping us plan for the cities of the future. Here is an article with links to a sample dataset and a process for analyzing Uber data.

8. Cyber-Profiling Criminals

Cyber-profiling is the process of collecting data from individuals and groups to identify significant co-relations. The idea of cyber profiling is derived from criminal profiles, which provide information on the investigation division to classify the types of criminals who were at the crime scene. Here is an interesting white paper on how to cyber-profile users in an academic environment based on user data preferences.

9. Call Record Detail Analysis

A call detail record (CDR) is the information captured by telecom companies during the call, SMS, and internet activity of a customer. This information provides greater insights about the customer's needs when used with customer demographics. In this article, you will understand how you can cluster customer activities for 24 hours by using the unsupervised k-means clustering algorithm. It is used to understand segments of customers with respect to their usage by hours.

10. Automatic Clustering of IT Alerts

Large enterprise IT infrastructure technology components such as network, storage, or database generate large volumes of alert messages. Because alert messages potentially point to operational issues, they must be manually screened for prioritization for downstream processes. Clustering of data can provide insight into categories of alerts and mean time to repair, and help in failure predictions.



The content team at DZone is passionate about helping developers. We need to know more about you in order to do that. Tell us who you are in our 8-minute Community Survey for a chance to win \$250. [Take the survey](#) ►



Like This Article? Read More From DZone



Solving a Clustering Problem Using the k-Means Algorithm With Oracle



Homemade Machine Learning in Python




Monitoring Real-Time Uber Data Using Apache APIs, Part 3: Real-Time Dashboard Using Vert.x



**Free DZone Refcard
Introduction to TensorFlow**

Topics: MACHINE LEARNING , AI , ALGORITHM , K-MEANS , CLUSTERING

Published at DZone with permission of Kaushik Raghupathi . [See the original article here.](#) 
Opinions expressed by DZone contributors are their own.
