

Introduction to Machine Learning

K-Means and *K*-Centers

Andres Mendez-Vazquez

February 25, 2019

Outline

1

K-Means Clustering

- The NP-Hard Problem
- *K*-Means Clustering Heuristic
- Convergence Criterion
- The Distance Function
- Example
- Properties of *K*-Means

2

The *K*-Center Criterion Clustering

- Introduction
- Re-Stating the *K*-center as a Clustering Problem
- Comparison with *K*-means
- The Greedy *K*-Center Algorithm
- Pseudo-Code
- The *K*-Center Algorithm
- Notes in Implementation
- Examples
- *K*-Center Algorithm Properties

Outline

1

K-Means Clustering

- The NP-Hard Problem
- *K*-Means Clustering Heuristic
- Convergence Criterion
- The Distance Function
- Example
- Properties of *K*-Means

2

The *K*-Center Criterion Clustering

- Introduction
- Re-Statting the *K*-center as a Clustering Problem
- Comparison with *K*-means
- The Greedy *K*-Center Algorithm
- Pseudo-Code
- The *K*-Center Algorithm
- Notes in Implementation
- Examples
- *K*-Center Algorithm Properties

The Hardness of K -means clustering

Definition

- Given a multiset $S \subseteq \mathbb{R}^d$, an integer k and $L \in \mathbb{R}$, is there a subset $T \subset \mathbb{R}^d$ with $|T| = k$ such that

$$\sum_{x \in S} \min_{t \in T} \|x - t\|^2 \leq L?$$

- The k -means clustering problem is NP-complete even for $d = 2$.

The Hardness of K -means clustering

Definition

- Given a multiset $S \subseteq \mathbb{R}^d$, an integer k and $L \in \mathbb{R}$, is there a subset $T \subset \mathbb{R}^d$ with $|T| = k$ such that

$$\sum_{x \in S} \min_{t \in T} \|x - t\|^2 \leq L?$$

Theorem

- The k -means clustering problem is NP-complete even for $d = 2$.

Reduction

The reduction to an NP-Complete problem

- Exact Cover by 3-Sets problem

Definition

- Given a finite set U containing exactly $3n$ elements and a collection $C = \{S_1, S_2, \dots, S_l\}$ of subsets of U each of which contains exactly 3 elements, Are there n sets in C such that their union is U ?

Reduction

The reduction to an NP-Complete problem

- Exact Cover by 3-Sets problem

Definition

- Given a finite set U containing exactly $3n$ elements and a collection $\mathcal{C} = \{S_1, S_2, \dots, S_l\}$ of subsets of U each of which contains exactly 3 elements, Are there n sets in \mathcal{C} such that their union is U ?

However

There are efficient heuristic and approximation algorithms

- Which can solve this problem

Outline

1

K-Means Clustering

- The NP-Hard Problem
- ***K*-Means Clustering Heuristic**
- Convergence Criterion
- The Distance Function
- Example
- Properties of *K*-Means

2

The *K*-Center Criterion Clustering

- Introduction
- Re-Stating the *K*-center as a Clustering Problem
- Comparison with *K*-means
- The Greedy *K*-Center Algorithm
- Pseudo-Code
- The *K*-Center Algorithm
- Notes in Implementation
- Examples
- *K*-Center Algorithm Properties

K-Means - Stuart Lloyd(Circa 1957)

History

Invented by Stuart Loyd in Bell Labs to obtain the best quantization in a signal data set.

Something More

The paper was published until 1982

Basically, what does it do?

It tries to find k points $\mu_1, \dots, \mu_k \in \mathbb{R}^d$ that minimize the expression (i.e. a partition S of the vector points):

$$\sum_{k=1}^K \sum_{i: x_i \in C_k} \|x_i - \mu_k\|^2 = \sum_{k=1}^K \sum_{i: x_i \in C_k} (x_i - \mu_k)^T (x_i - \mu_k)$$

K-Means - Stuart Lloyd(Circa 1957)

History

Invented by Stuart Loyd in Bell Labs to obtain the best quantization in a signal data set.

Something Notable

The paper was published until 1982

Basically, it's an iterative algorithm

It tries to find k points $\mu_1, \dots, \mu_k \in \mathbb{R}^d$ that minimize the expression (i.e. a partition S of the vector points):

$$\sum_{k=1}^K \sum_{i: x_i \in C_k} \|x_i - \mu_k\|^2 = \sum_{k=1}^K \sum_{i: x_i \in C_k} (x_i - \mu_k)^T (x_i - \mu_k)$$

K -Means - Stuart Lloyd(Circa 1957)

History

Invented by Stuart Loyd in Bell Labs to obtain the best quantization in a signal data set.

Something Notable

The paper was published until 1982

Basically given N vectors $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^d$

It tries to find k points $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k \in \mathbb{R}^d$ that minimize the expression (i.e. a partition S of the vector points):

$$\sum_{k=1}^N \sum_{i:\mathbf{x}_i \in C_k} \|\mathbf{x}_i - \boldsymbol{\mu}_k\|^2 = \sum_{k=1}^N \sum_{i:\mathbf{x}_i \in C_k} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T (\mathbf{x}_i - \boldsymbol{\mu}_k)$$

K-means clustering

K-means

It is a partitional clustering algorithm.

K-means clustering

K-means

It is a partitional clustering algorithm.

Definition

Let the set of data points (or instances) D be $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ where $\mathbf{x}_i = (x_{i1}, \dots, x_{ir})^T$:

- The K -means algorithm partitions the given data into K clusters.
- Each cluster has a cluster center, called centroid.
- K is specified by the user.

K-means clustering

K-means

It is a partitional clustering algorithm.

Definition

Let the set of data points (or instances) D be $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ where $\mathbf{x}_i = (x_{i1}, \dots, x_{ir})^T$:

- The K -means algorithm partitions the given data into K clusters.
 - Each cluster has a cluster center, called centroid.
 - K is specified by the user.

K-means clustering

K-means

It is a partitional clustering algorithm.

Definition

Let the set of data points (or instances) D be $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ where $\mathbf{x}_i = (x_{i1}, \dots, x_{ir})^T$:

- The K -means algorithm partitions the given data into K clusters.
- Each cluster has a cluster center, called centroid.

• K is specified by the user

K-means clustering

K-means

It is a partitional clustering algorithm.

Definition

Let the set of data points (or instances) D be $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ where $\mathbf{x}_i = (x_{i1}, \dots, x_{ir})^T$:

- The K -means algorithm partitions the given data into K clusters.
- Each cluster has a cluster center, called centroid.
- K is specified by the user.

K-means algorithm

The *K*-means algorithm works as follows

Given k as the possible number of cluster:

- ➊ Randomly choose K data points (seeds) to be the initial centroids, cluster centers,

- ▶ $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$

- ➋ Assign each data point to the closest centroid

- ▶ $c_i = \arg \min_j \{dist(\mathbf{x}_i - \mathbf{v}_j)\}$

- ➌ Re-compute the centroids using the current cluster memberships.

$$\mathbf{v}_j = \frac{\sum_{i=1}^n I(c_i = j) \mathbf{x}_i}{\sum_{i=1}^n I(c_i = j)}$$

- ➍ If a convergence criterion is not met, go to 2

K-means algorithm

The *K*-means algorithm works as follows

Given k as the possible number of cluster:

- ① Randomly choose K data points (seeds) to be the initial **centroids**, cluster centers,

- ▶ $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$

- ② Assign each data point to the closest centroid

- ▶ $c_j = \arg \min_j \{dist(\mathbf{x}_i - \mathbf{v}_j)\}$

- ③ Re-compute the **centroids** using the current cluster memberships.

$$\mathbf{v}_j = \frac{\sum_{i=1}^n I(c_i = j) \mathbf{x}_i}{\sum_{i=1}^n I(c_i = j)}$$

- ④ If a convergence criterion is not met, go to 2

K-means algorithm

The *K*-means algorithm works as follows

Given k as the possible number of cluster:

- ① Randomly choose K data points (seeds) to be the initial **centroids**, cluster centers,

- ▶ $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$

- ② Assign each data point to the closest **centroid**

- ▶ $c_i = \arg \min_j \{dist(\mathbf{x}_i - \mathbf{v}_j)\}$

- ③ Re-compute the centroids using the current cluster memberships.

$$\mathbf{v}_j = \frac{\sum_{i=1}^n I(c_i = j) \mathbf{x}_i}{\sum_{i=1}^n I(c_i = j)}$$

- ④ If a convergence criterion is not met, go to 2

K -means algorithm

The K -means algorithm works as follows

Given k as the possible number of cluster:

- ① Randomly choose K data points (seeds) to be the initial **centroids**, cluster centers,
 - ▶ $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$
- ② Assign each data point to the closest **centroid**
 - ▶ $c_i = \arg \min_j \{dist(\mathbf{x}_i - \mathbf{v}_j)\}$
- ③ Re-compute the **centroids** using the current cluster memberships.

$$\mathbf{v}_j = \frac{\sum_{i=1}^n I(c_i = j) \mathbf{x}_i}{\sum_{i=1}^n I(c_i = j)}$$

④ If a convergence criterion is not met, go to 2

K-means algorithm

The *K*-means algorithm works as follows

Given k as the possible number of cluster:

- ① Randomly choose K data points (seeds) to be the initial **centroids**, cluster centers,
 - ▶ $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$
- ② Assign each data point to the closest **centroid**
 - ▶ $c_i = \arg \min_j \{dist(\mathbf{x}_i - \mathbf{v}_j)\}$
- ③ Re-compute the **centroids** using the current cluster memberships.

$$\mathbf{v}_j = \frac{\sum_{i=1}^n I(c_i = j) \mathbf{x}_i}{\sum_{i=1}^n I(c_i = j)}$$

- ④ If a convergence criterion is not met, go to 2.

What is the code trying to do?

It is trying to find a partition S

K -means tries to find a partition S such that it minimizes the cost function:

$$\min_S \sum_{k=1}^N \sum_{i: x_i \in C_k} (x_i - \mu_k)^T (x_i - \mu_k) \quad (1)$$

Where μ_k is the centroid for cluster C_k .

$$\mu_k = \frac{1}{N_k} \sum_{i: x_i \in C_k} x_i \quad (2)$$

Where N_k is the number of samples in the cluster C_k .

What is the code trying to do?

It is trying to find a partition S

K -means tries to find a partition S such that it minimizes the cost function:

$$\min_S \sum_{k=1}^N \sum_{i: x_i \in C_k} (x_i - \mu_k)^T (x_i - \mu_k) \quad (1)$$

Where μ_k is the centroid for cluster C_k .

$$\mu_k = \frac{1}{N_k} \sum_{i: x_i \in C_k} x_i \quad (2)$$

Where N_k is the number of samples in the cluster C_k .

What is the code trying to do?

It is trying to find a partition S

K -means tries to find a partition S such that it minimizes the cost function:

$$\min_S \sum_{k=1}^N \sum_{i:x_i \in C_k} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T (\mathbf{x}_i - \boldsymbol{\mu}_k) \quad (1)$$

Where $\boldsymbol{\mu}_k$ is the centroid for cluster C_k

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{i:x_i \in C_k} \mathbf{x}_i \quad (2)$$

Where N_k is the number of samples in the cluster C_k .

Outline

1

K-Means Clustering

- The NP-Hard Problem
- *K*-Means Clustering Heuristic
- Convergence Criterion**
- The Distance Function
- Example
- Properties of *K*-Means

2

The *K*-Center Criterion Clustering

- Introduction
- Re-Stating the *K*-center as a Clustering Problem
- Comparison with *K*-means
- The Greedy *K*-Center Algorithm
- Pseudo-Code
- The *K*-Center Algorithm
- Notes in Implementation
- Examples
- *K*-Center Algorithm Properties

What Stopping/convergence criterion should we use?

First

No (or minimum) re-assignments of data points to different clusters.

What Stopping/convergence criterion should we use?

First

No (or minimum) re-assignments of data points to different clusters.

Second

No (or minimum) change of centroids.

What Stopping/convergence criterion should we use?

First

No (or minimum) re-assignments of data points to different clusters.

Second

No (or minimum) change of centroids.

Third

Minimum decrease in the sum of squared error (SSE),

- C_k is cluster k .

- v_k is the centroid of cluster C_k .

$$SSE = \sum_{k=1}^K \sum_{x \in C_k} \text{dist}(x, v_k)^2$$

What Stopping/convergence criterion should we use?

First

No (or minimum) re-assignments of data points to different clusters.

Second

No (or minimum) change of centroids.

Third

Minimum decrease in the sum of squared error (SSE),

- C_k is cluster k .

• v_k is the centroid of cluster C_k .

$$SSE = \sum_{k=1}^K \sum_{x \in C_k} \text{dist}(x, v_k)^2$$

What Stopping/convergence criterion should we use?

First

No (or minimum) re-assignments of data points to different clusters.

Second

No (or minimum) change of centroids.

Third

Minimum decrease in the sum of squared error (SSE),

- C_k is cluster k .
- \mathbf{v}_k is the centroid of cluster C_k .

$$SSE = \sum_{k=1}^K \sum_{x \in c_k} dist(\mathbf{x}, \mathbf{v}_k)^2$$

Outline

1

K-Means Clustering

- The NP-Hard Problem
- *K*-Means Clustering Heuristic
- Convergence Criterion
- The Distance Function**
- Example
- Properties of *K*-Means

2

The *K*-Center Criterion Clustering

- Introduction
- Re-Stating the *K*-center as a Clustering Problem
- Comparison with *K*-means
- The Greedy *K*-Center Algorithm
- Pseudo-Code
- The *K*-Center Algorithm
- Notes in Implementation
- Examples
- *K*-Center Algorithm Properties

The distance function

Actually, we have the following distance functions:

Euclidean

$$dist(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\| = \sqrt{(\mathbf{x} - \mathbf{y})^T (\mathbf{x} - \mathbf{y})}$$

Manhattan

$$dist(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_1 = \sum_{i=1}^n |x_i - y_i|$$

Minkowski

$$dist(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_A = \sqrt{(\mathbf{x} - \mathbf{y})^T A (\mathbf{x} - \mathbf{y})}$$

The distance function

Actually, we have the following distance functions:

Euclidean

$$dist(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\| = \sqrt{(\mathbf{x} - \mathbf{y})^T (\mathbf{x} - \mathbf{y})}$$

Manhattan

$$dist(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_1 = \sum_{i=1}^n |x_i - y_i|$$

Minkowski

$$dist(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_A = \sqrt{(\mathbf{x} - \mathbf{y})^T A (\mathbf{x} - \mathbf{y})}$$

The distance function

Actually, we have the following distance functions:

Euclidean

$$dist(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\| = \sqrt{(\mathbf{x} - \mathbf{y})^T (\mathbf{x} - \mathbf{y})}$$

Manhattan

$$dist(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_1 = \sum_{i=1}^n |x_i - y_i|$$

Mahalanobis

$$dist(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_A = \sqrt{(\mathbf{x} - \mathbf{y})^T A (\mathbf{x} - \mathbf{y})}$$

Outline

1

K-Means Clustering

- The NP-Hard Problem
- *K*-Means Clustering Heuristic
- Convergence Criterion
- The Distance Function
- Example
- Properties of *K*-Means

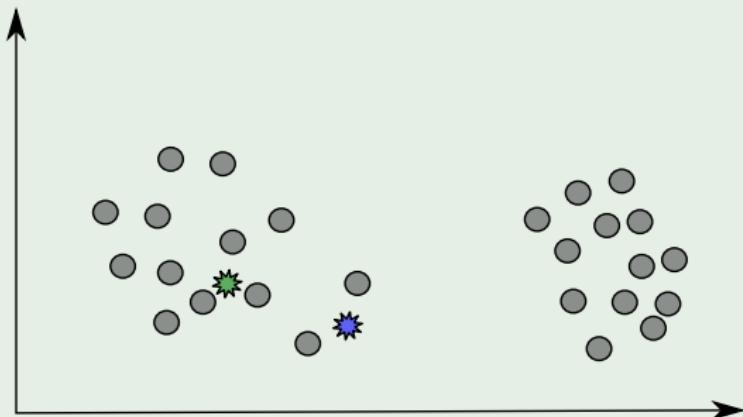
2

The *K*-Center Criterion Clustering

- Introduction
- Re-Stating the *K*-center as a Clustering Problem
- Comparison with *K*-means
- The Greedy *K*-Center Algorithm
- Pseudo-Code
- The *K*-Center Algorithm
- Notes in Implementation
- Examples
- *K*-Center Algorithm Properties

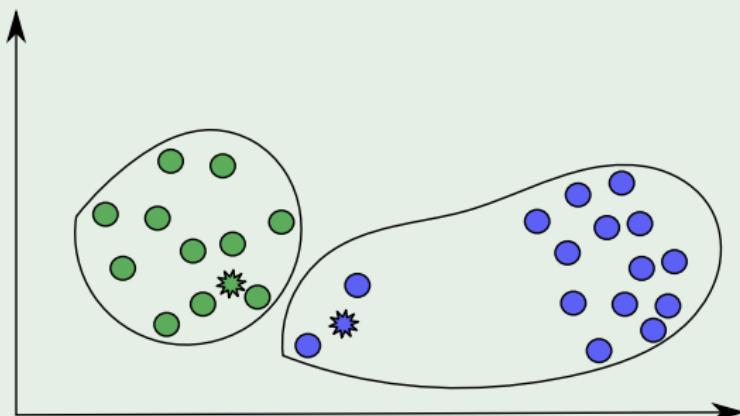
An example

Dropping two possible centroids



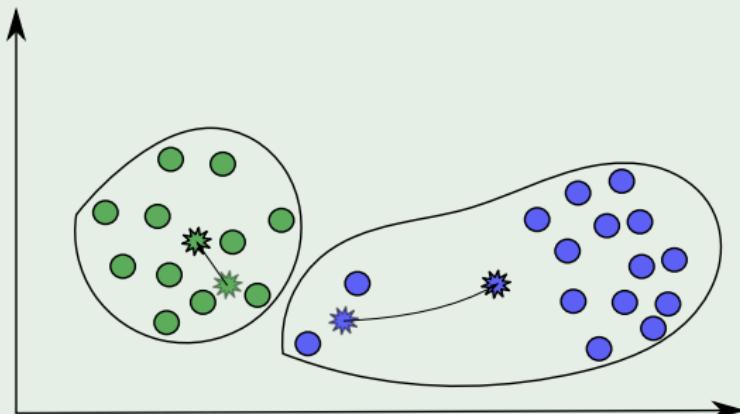
An example

Calculate the memberships



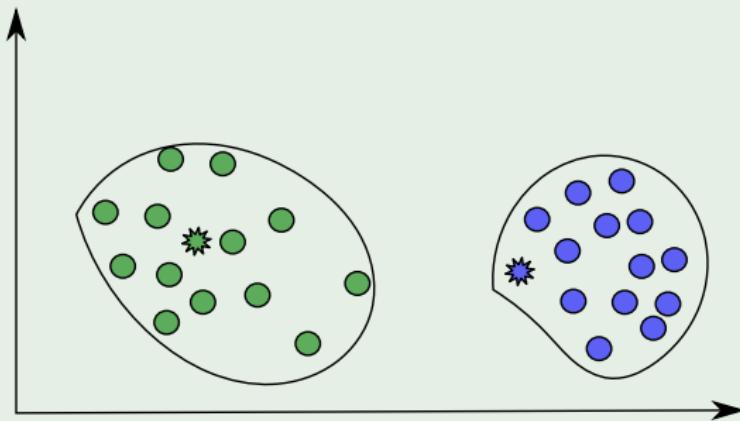
An example

We re-calculate centroids



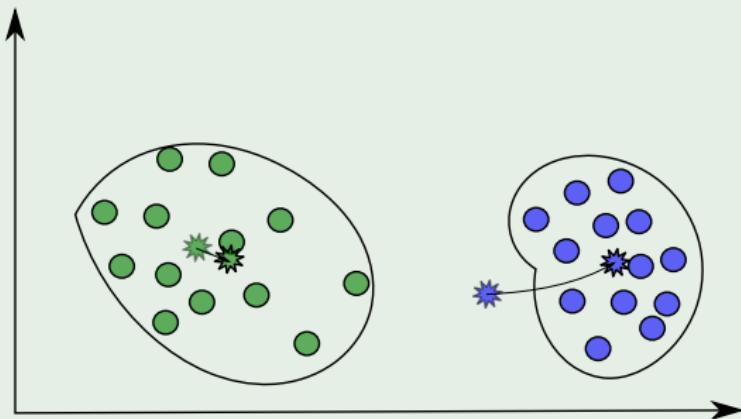
An example

We re-calculate memberships



An example

We re-calculate centroids and keep going



Outline

1

K-Means Clustering

- The NP-Hard Problem
- *K*-Means Clustering Heuristic
- Convergence Criterion
- The Distance Function
- Example
- Properties of *K*-Means

2

The *K*-Center Criterion Clustering

- Introduction
- Re-Stating the *K*-center as a Clustering Problem
- Comparison with *K*-means
- The Greedy *K*-Center Algorithm
- Pseudo-Code
- The *K*-Center Algorithm
- Notes in Implementation
- Examples
- *K*-Center Algorithm Properties

Strengths of K -means

Strengths

- Simple: easy to understand and to implement
- Efficient: Time complexity: $O(tKN)$, where N is the number of data points, K is the number of clusters, and t is the number of iterations.
- Since both K and t are small, K -means is considered a linear algorithm.

Strengths of K -means

Strengths

- Simple: easy to understand and to implement
- Efficient: Time complexity: $O(tKN)$, where N is the number of data points, K is the number of clusters, and t is the number of iterations.
 - Since both K and t are small, K -means is considered a linear algorithm.

Popularity

K -means is the most popular clustering algorithm.

Strengths of K -means

Strengths

- Simple: easy to understand and to implement
- Efficient: Time complexity: $O(tKN)$, where N is the number of data points, K is the number of clusters, and t is the number of iterations.
- Since both K and t are small. K -means is considered a linear algorithm.

Popularity

K -means is the most popular clustering algorithm.

More than

It terminates at a local optimum if SSE is used. The global optimum is hard to find due to complexity.

Strengths of K -means

Strengths

- Simple: easy to understand and to implement
- Efficient: Time complexity: $O(tKN)$, where N is the number of data points, K is the number of clusters, and t is the number of iterations.
- Since both K and t are small. K -means is considered a linear algorithm.

Popularity

K -means is the most popular clustering algorithm.

More Info

It terminates at a local optimum if SSE is used. The global optimum is hard to find due to complexity.

Strengths of K -means

Strengths

- Simple: easy to understand and to implement
- Efficient: Time complexity: $O(tKN)$, where N is the number of data points, K is the number of clusters, and t is the number of iterations.
- Since both K and t are small. K -means is considered a linear algorithm.

Popularity

K -means is the most popular clustering algorithm.

Note that

It terminates at a **local optimum** if SSE is used. The **global optimum** is hard to find due to complexity.

Weaknesses of K -means

Important

The algorithm is only applicable if the mean is defined.

- For categorical data, K -mode – the centroid is represented by most frequent values

Weaknesses of K -means

Important

The algorithm is only applicable if the mean is defined.

- For categorical data, K -mode - the centroid is represented by most frequent values.

In addition

The user needs to specify K .

Weaknesses of K -means

Important

The algorithm is only applicable if the mean is defined.

- For categorical data, K -mode - the centroid is represented by most frequent values.

In addition

The user needs to specify K .

Outliers

The algorithm is sensitive to outliers.

- Outliers are data points that are very far away from other data points.
- Outliers could be errors in the data recording or some special data points with very different values.

Weaknesses of K -means

Important

The algorithm is only applicable if the mean is defined.

- For categorical data, K -mode - the centroid is represented by most frequent values.

In addition

The user needs to specify K .

Outliers

The algorithm is sensitive to **outliers**.

- Outliers are data points that are very far away from other data points.
- Outliers could be errors in the data recording or some special data points with very different values.

Weaknesses of K -means

Important

The algorithm is only applicable if the mean is defined.

- For categorical data, K -mode - the centroid is represented by most frequent values.

In addition

The user needs to specify K .

Outliers

The algorithm is sensitive to **outliers**.

- Outliers are data points that are very far away from other data points.
 - Outliers could be errors in the data recording or some special data points with very different values.

Weaknesses of K -means

Important

The algorithm is only applicable if the mean is defined.

- For categorical data, K -mode - the centroid is represented by most frequent values.

In addition

The user needs to specify K .

Outliers

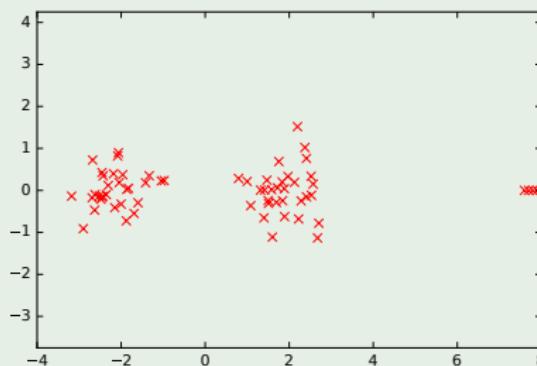
The algorithm is sensitive to **outliers**.

- Outliers are data points that are very far away from other data points.
- Outliers could be errors in the data recording or some special data points with very different values.

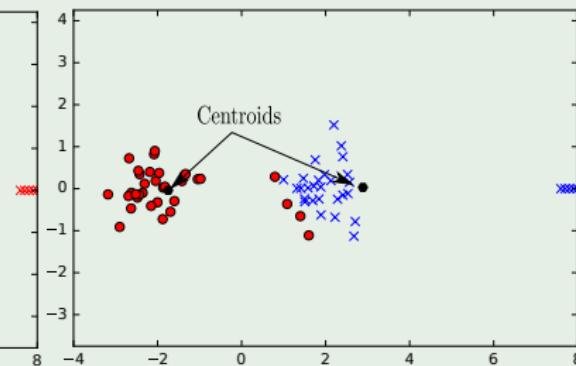
Weaknesses of K -means: Problems with outliers

A series of outliers

Initial Data

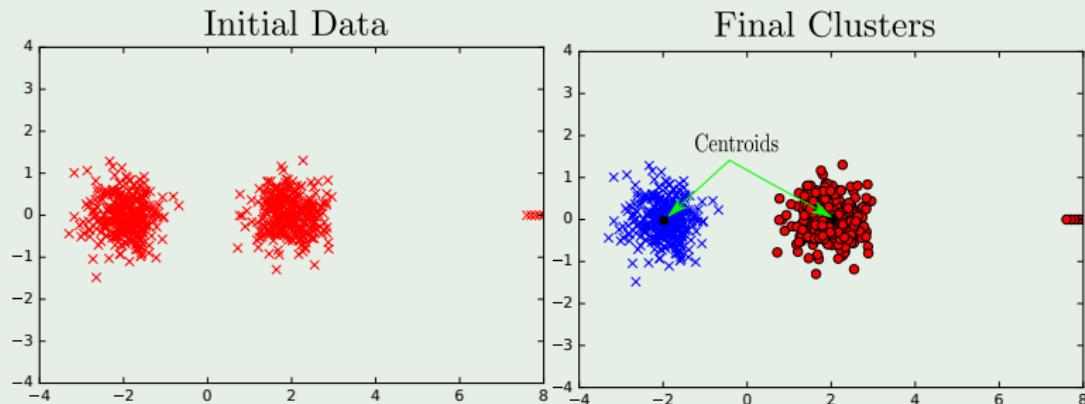


Final Clusters



Weaknesses of K -means: Problems with outliers

Nevertheless, if you have more dense clusters



Weaknesses of K -means: How to deal with outliers

One method

To remove some data points in the clustering process that are much further away from the centroids than other data points.

- To be safe, we may want to monitor these possible outliers over a few iterations and then decide to remove them.

Weaknesses of K -means: How to deal with outliers

One method

To remove some data points in the clustering process that are much further away from the centroids than other data points.

- To be safe, we may want to monitor these possible outliers over a few iterations and then decide to remove them.

Another method

To perform random sampling.

- Since in sampling we only choose a small subset of the data points, the chance of selecting an outlier is very small.
- Assign the rest of the data points to the clusters by distance or similarity comparison, or classification.

Weaknesses of K -means: How to deal with outliers

One method

To remove some data points in the clustering process that are much further away from the centroids than other data points.

- To be safe, we may want to monitor these possible outliers over a few iterations and then decide to remove them.

Another method

To perform random sampling.

- Since in sampling we only choose a small subset of the data points, the chance of selecting an outlier is very small.
- Assign the rest of the data points to the clusters by distance or similarity comparison, or classification.

Weaknesses of K -means: How to deal with outliers

One method

To remove some data points in the clustering process that are much further away from the centroids than other data points.

- To be safe, we may want to monitor these possible outliers over a few iterations and then decide to remove them.

Another method

To perform random sampling.

- Since in sampling we only choose a small subset of the data points, the chance of selecting an outlier is very small.
 - Assign the rest of the data points to the clusters by distance or similarity comparison, or classification.

Weaknesses of K -means: How to deal with outliers

One method

To remove some data points in the clustering process that are much further away from the centroids than other data points.

- To be safe, we may want to monitor these possible outliers over a few iterations and then decide to remove them.

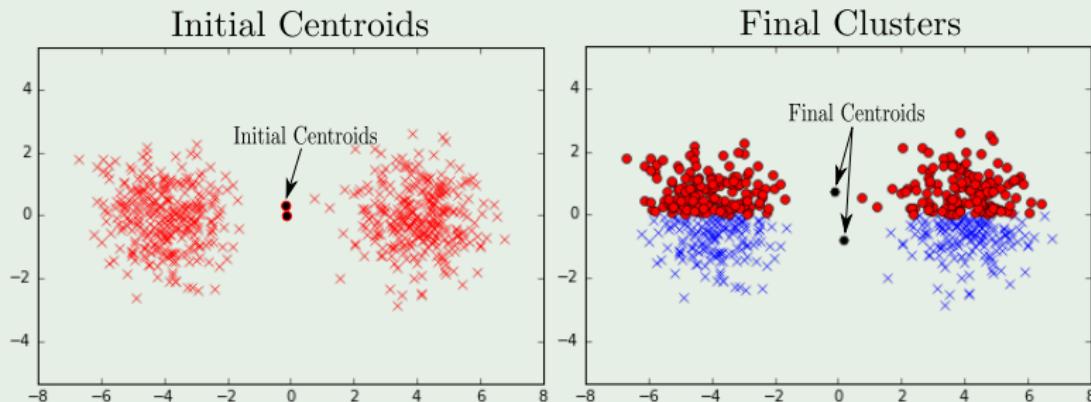
Another method

To perform random sampling.

- Since in sampling we only choose a small subset of the data points, the chance of selecting an outlier is very small.
- Assign the rest of the data points to the clusters by distance or similarity comparison, or classification.

Weaknesses of K -means (cont...)

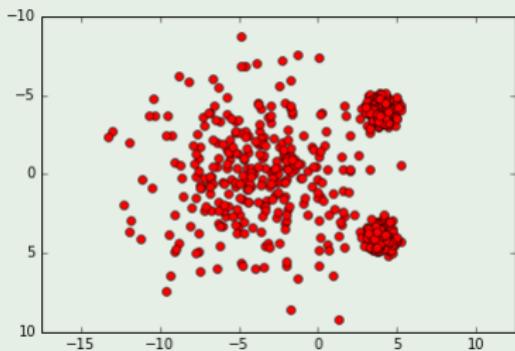
The algorithm is sensitive to **initial seeds**



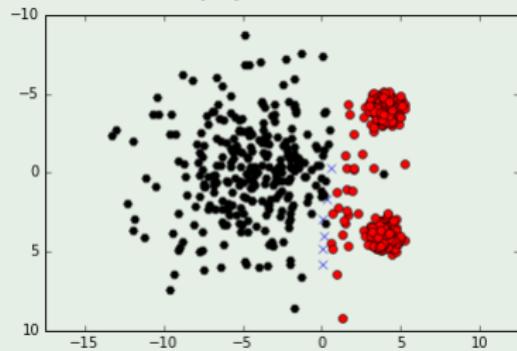
Weaknesses of K -means : Different Densities

We have three cluster nevertheless

DATA

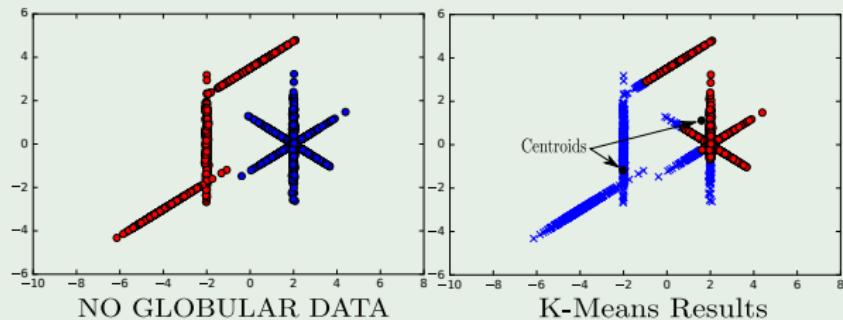


3 Clusters



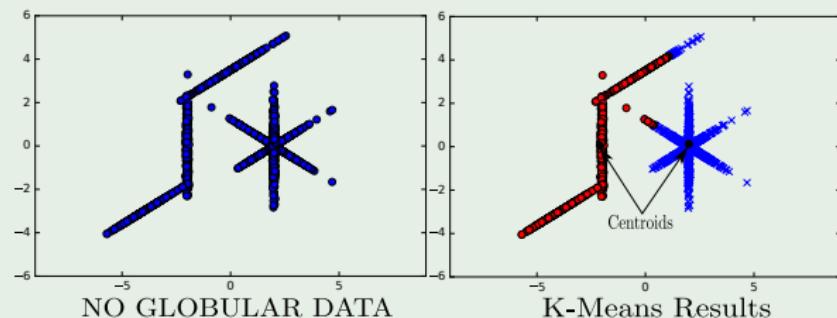
Weaknesses of K -means: Non-globular Shapes

Here, we notice that K -means may only detect globular shapes



Weaknesses of K -means: Non-globular Shapes

However, it sometimes work better than expected



Outline

1

K-Means Clustering

- The NP-Hard Problem
- *K*-Means Clustering Heuristic
- Convergence Criterion
- The Distance Function
- Example
- Properties of *K*-Means

2

The *K*-Center Criterion Clustering

● Introduction

- Re-Stating the *K*-center as a Clustering Problem
- Comparison with *K*-means
- The Greedy *K*-Center Algorithm
- Pseudo-Code
- The *K*-Center Algorithm
- Notes in Implementation
- Examples
- *K*-Center Algorithm Properties

The K -center Problem

The input

It is a set of points with distances represented by a weighted graph
 $G = (V, V \times V)$.

The K -center Problem

The input

It is a set of points with distances represented by a weighted graph
 $G = (V, V \times V)$.

Output

- Select K centroids in G .
 - such that minimize maximum distance of every point to a centroid

The K -center Problem

The input

It is a set of points with distances represented by a weighted graph
 $G = (V, V \times V)$.

Output

- Select K centroids in G .
- Such that minimize maximum distance of every point to a centroid.

Theorem (In the general case for any distance)

It is NP-hard to approximate the general K -center problem within any factor α .

The K -center Problem

The input

It is a set of points with distances represented by a weighted graph
 $G = (V, V \times V)$.

Output

- Select K centroids in G .
- Such that minimize maximum distance of every point to a centroid.

Theorem (In the general case for any distance)

It is NP-hard to approximate the general K -center problem within any factor α .

Therefore

We change the distance to be constrained by

The Triangle Inequality

Given x, y and z

$$L(x, z) \leq L(x, y) + L(y, z)$$

Therefore

We change the distance to be constrained by

The Triangle Inequality

Given x, y and z

$$L(x, z) \leq L(x, y) + L(y, z)$$

Outline

1

K-Means Clustering

- The NP-Hard Problem
- *K*-Means Clustering Heuristic
- Convergence Criterion
- The Distance Function
- Example
- Properties of *K*-Means

2

The *K*-Center Criterion Clustering

- Introduction
- Re-Stating the *K*-center as a Clustering Problem
- Comparison with *K*-means
- The Greedy *K*-Center Algorithm
- Pseudo-Code
- The *K*-Center Algorithm
- Notes in Implementation
- Examples
- *K*-Center Algorithm Properties

We have a new criterion

Instead of using the K -mean criterion

We use the K -center criterion under the triangle inequality.

New criterion

The K -center criterion partitions the points into K clusters so as to minimize the maximum distance of any point to its cluster center.

We have a new criterion

Instead of using the K -mean criterion

We use the K -center criterion under the triangle inequality.

New Criterion

The K -center criterion partitions the points into K clusters so as to minimize the maximum distance of any point to its cluster center.

Explanation

Setup

Suppose we have a data set A that contains N objects.

Explanation

Setup

Suppose we have a data set A that contains N objects.

We want the following

We want to partition A into K sets labeled C_1, C_2, \dots, C_K .

Explanation

Setup

Suppose we have a data set A that contains N objects.

We want the following

We want to partition A into K sets labeled C_1, C_2, \dots, C_K .

Now

Define a cluster size for any cluster C_k as follows.

- The smallest value D for which all points in this cluster C_k are:
 - Within distance D of each other.
 - Or within distance $\frac{D}{2}$ of some point called the cluster center.

Explanation

Setup

Suppose we have a data set A that contains N objects.

We want the following

We want to partition A into K sets labeled C_1, C_2, \dots, C_K .

Now

Define a cluster size for any cluster C_k as follows.

- The smallest value D for which all points in this cluster C_k are:
 - Within distance D of each other.
 - Or within distance $\frac{D}{2}$ of some point called the cluster center.

Explanation

Setup

Suppose we have a data set A that contains N objects.

We want the following

We want to partition A into K sets labeled C_1, C_2, \dots, C_K .

Now

Define a cluster size for any cluster C_k as follows.

- The smallest value D for which all points in this cluster C_k are:
 - ▶ Within distance D of each other.
 - ▶ Or within distance $\frac{D}{2}$ of some point called the cluster center.

Explanation

Setup

Suppose we have a data set A that contains N objects.

We want the following

We want to partition A into K sets labeled C_1, C_2, \dots, C_K .

Now

Define a cluster size for any cluster C_k as follows.

- The smallest value D for which all points in this cluster C_k are:
 - ▶ Within distance D of each other.
 - ▶ Or within distance $\frac{D}{2}$ of some point called the cluster center.

Another Way

We have

Another way to define the cluster size:

- D is the maximum pairwise distance between an arbitrary pair of points in the cluster.
- Or twice the maximum distance between a data point and a chosen centroid.

This

Denoting the cluster size of C_k by D_k , we have that the cluster size of partition (the way the points are grouped) S by:

$$D = \max_{k=1,\dots,K} D_k \quad (3)$$

In another word,

The cluster size of a partition composed of multiple clusters is the maximum size of these clusters.

Another Way

We have

Another way to define the cluster size:

- D is the maximum pairwise distance between an arbitrary pair of points in the cluster.
- Or twice the maximum distance between a data point and a chosen centroid.

Thus

Denoting the cluster size of C_k by D_k , we have that the cluster size of partition (the way the points are grouped) S by:

$$D = \max_{k=1,\dots,K} D_k \quad (3)$$

The maximum cluster size

The cluster size of a partition composed of multiple clusters is the maximum size of these clusters.

Another Way

We have

Another way to define the cluster size:

- D is the maximum pairwise distance between an arbitrary pair of points in the cluster.
- Or twice the maximum distance between a data point and a chosen centroid.

Thus

Denoting the cluster size of C_k by D_k , we have that the cluster size of partition (the way the points are grouped) S by:

$$D = \max_{k=1,\dots,K} D_k \quad (3)$$

In another words

The cluster size of a partition composed of multiple clusters is the maximum size of these clusters.

Outline

1

K-Means Clustering

- The NP-Hard Problem
- *K*-Means Clustering Heuristic
- Convergence Criterion
- The Distance Function
- Example
- Properties of *K*-Means

2

The *K*-Center Criterion Clustering

- Introduction
- Re-Stating the *K*-center as a Clustering Problem
- Comparison with *K*-means
- The Greedy *K*-Center Algorithm
- Pseudo-Code
- The *K*-Center Algorithm
- Notes in Implementation
- Examples
- *K*-Center Algorithm Properties

Comparison with K -means

We use the following distance for comparison

In order to compare these methods, we use Euclidean distance.

In K -means we assume the distance between vectors is the squared Euclidean distance

K -means tries to find a partition S that minimizes:

$$\min_S \sum_{k=1}^N \sum_{x_i \in C_k} (x_i - \mu_k)^T (x_i - \mu_k) \quad (4)$$

What is the centroid of cluster k ?

$$\mu_k = \frac{1}{N_k} \sum_{x_i \in C_k} x_i \quad (5)$$

Comparison with K -means

We use the following distance for comparison

In order to compare these methods, we use Euclidean distance.

In K -means we assume the distance between vectors is the squared Euclidean distance

K -means tries to find a partition S that minimizes:

$$\min_S \sum_{k=1}^N \sum_{i: x_i \in C_k} (x_i - \mu_k)^T (x_i - \mu_k) \quad (4)$$

What is the centroid of a cluster?

$$\mu_k = \frac{1}{N_k} \sum_{x_i \in C_k} x_i \quad (5)$$

Comparison with K -means

We use the following distance for comparison

In order to compare these methods, we use Euclidean distance.

In K -means we assume the distance between vectors is the squared Euclidean distance

K -means tries to find a partition S that minimizes:

$$\min_S \sum_{k=1}^N \sum_{i:x_i \in C_k} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T (\mathbf{x}_i - \boldsymbol{\mu}_k) \quad (4)$$

Where $\boldsymbol{\mu}_k$ is the centroid for cluster C_k

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{i:x_i \in C_k} \mathbf{x}_i \quad (5)$$

Now, what about the K -centers

K -center, on the other hand, minimizes the worst case distance to these centroids

$$\min_S \max_{k=1,\dots,K} \max_{i:x_i \in C_k} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T (\mathbf{x}_i - \boldsymbol{\mu}_k) \quad (6)$$

Now, what about the K -centers

K -center, on the other hand, minimizes the worst case distance to these centroids

$$\min_S \max_{k=1,\dots,K} \max_{i:x_i \in C_k} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T (\mathbf{x}_i - \boldsymbol{\mu}_k) \quad (6)$$

Where

$\boldsymbol{\mu}_k$ is called the “centroid”, but may not be the mean vector.

Now, what about the K -centers

K -center, on the other hand, minimizes the worst case distance to these centroids

$$\min_S \max_{k=1,\dots,K} \max_{i:x_i \in C_k} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T (\mathbf{x}_i - \boldsymbol{\mu}_k) \quad (6)$$

Where

$\boldsymbol{\mu}_k$ is called the “centroid”, but may not be the mean vector.

Properties

- The above objective function shows that for each cluster, only the worst scenario matters, that is, the farthest data point to the centroid.
- Moreover, among the clusters, only the worst cluster matters, whose farthest data point yields the maximum distance to the centroid comparing with the farthest data points of the other clusters.

Now, what about the K -centers

K -center, on the other hand, minimizes the worst case distance to these centroids

$$\min_S \max_{k=1,\dots,K} \max_{i:x_i \in C_k} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T (\mathbf{x}_i - \boldsymbol{\mu}_k) \quad (6)$$

Where

$\boldsymbol{\mu}_k$ is called the “centroid”, but may not be the mean vector.

Properties

- The above objective function shows that for each cluster, only the worst scenario matters, that is, the farthest data point to the centroid.
- Moreover, among the clusters, only the worst cluster matters, whose farthest data point yields the maximum distance to the centroid comparing with the farthest data points of the other clusters.

In other words

First

This minimax type of problem is much harder to solve than solving the objective function of k-means.

Second, centroids are replaced by cluster boundaries. We can calculate max distance instead of using centroids.

$$\min_S \max_{k=1,\dots,K} \max_{x_i, x_j \in C_k} L(x_i, x_j) \quad (7)$$

Why?

$L(x_i, x_j)$ denotes any distance between a pair of objects in the same cluster.

In other words

First

This minimax type of problem is much harder to solve than solving the objective function of k-means.

Another formulation of K -center minimizes the worst case pairwise distance instead of using centroids

$$\min_S \max_{k=1, \dots, K} \max_{i, j: x_i, x_j \in C_k} L(x_i, x_j) \quad (7)$$

What?

$L(x_i, x_j)$ denotes any distance between a pair of objects in the same cluster.

In other words

First

This minimax type of problem is much harder to solve than solving the objective function of k-means.

Another formulation of K -center minimizes the worst case pairwise distance instead of using centroids

$$\min_S \max_{k=1, \dots, K} \max_{i, j: x_i, x_j \in C_k} L(x_i, x_j) \quad (7)$$

With

$L(x_i, x_j)$ denotes any distance between a pair of objects in the same cluster.

Outline

1

K-Means Clustering

- The NP-Hard Problem
- *K*-Means Clustering Heuristic
- Convergence Criterion
- The Distance Function
- Example
- Properties of *K*-Means

2

The *K*-Center Criterion Clustering

- Introduction
- Re-Stating the *K*-center as a Clustering Problem
- Comparison with *K*-means
- The Greedy *K*-Center Algorithm**
- Pseudo-Code
- The *K*-Center Algorithm
- Notes in Implementation
- Examples
- *K*-Center Algorithm Properties

Greedy Algorithm for K -Center

Main Idea

The idea behind the Greedy Algorithm is to choose a subset H from the original dataset S consisting of K points that are farthest apart from each other.

Intuition

Since the points in set H are far apart then the worst-case scenario has been taken care of and hopefully the cluster size for the partition is small.

Implementation

Each point $h_k \in H$ represents one cluster or subset of points C_k .

Greedy Algorithm for K -Center

Main Idea

The idea behind the Greedy Algorithm is to choose a subset H from the original dataset S consisting of K points that are farthest apart from each other.

Intuition

Since the points in set H are far apart then the worst-case scenario has been taken care of and hopefully the cluster size for the partition is small.



Each point $h_k \in H$ represents one cluster or subset of points C_k .

Greedy Algorithm for K -Center

Main Idea

The idea behind the Greedy Algorithm is to choose a subset H from the original dataset S consisting of K points that are farthest apart from each other.

Intuition

Since the points in set H are far apart then the worst-case scenario has been taken care of and hopefully the cluster size for the partition is small.

Thus

Each point $h_k \in H$ represents one cluster or subset of points C_k .

Then

Something Notable

We can think of it as a centroid.

Important

Technically it is not a centroid because it tends to be at the boundary of a cluster, but conceptually we can think of it as a centroid.

Important

The way that we partition these points given the centroids is the same as in K -means, that is, the nearest-neighbor rule.

Then

Something Notable

We can think of it as a centroid.

However

Technically it is not a centroid because it tends to be at the boundary of a cluster, but conceptually we can think of it as a centroid.

Decision Rule

The way that we partition these points given the centroids is the same as in K -means, that is, the nearest-neighbor rule.

Then

Something Notable

We can think of it as a centroid.

However

Technically it is not a centroid because it tends to be at the boundary of a cluster, but conceptually we can think of it as a centroid.

Important

The way that we partition these points given the centroids is the same as in K -means, that is, the nearest-neighbor rule.

Specifically

We do the following

For every point x_i , in order to see which cluster C_k it is partitioned into, we compute its distance to each cluster centroid as follows, and find out which centroid is the closest:

$$L(x_i, \mathbf{h}_k) = \min_{k'=1, \dots, K} L(x_i, \mathbf{h}_{k'}) \quad (8)$$



Whichever centroid with the minimum distance is selected as the cluster for x_i .

Specifically

We do the following

For every point x_i , in order to see which cluster C_k it is partitioned into, we compute its distance to each cluster centroid as follows, and find out which centroid is the closest:

$$L(x_i, \mathbf{h}_k) = \min_{k'=1, \dots, K} L(x_i, \mathbf{h}_{k'}) \quad (8)$$

Thus

Whichever centroid with the minimum distance is selected as the cluster for x_i .

Important

For K -center clustering

We only need pairwise distance $L(\mathbf{x}_i, \mathbf{x}_j)$ for any $\mathbf{x}_i, \mathbf{x}_j \in S$.

Where

\mathbf{x}_i can be a non-vector representation of the objects.

Allowing vector or non-vector

$L(\mathbf{x}_i, \mathbf{x}_j)$ which makes the K -center more general than the K -means.

Important

For K -center clustering

We only need pairwise distance $L(\mathbf{x}_i, \mathbf{x}_j)$ for any $\mathbf{x}_i, \mathbf{x}_j \in S$.

Where

\mathbf{x}_i can be a non-vector representation of the objects.

Assuming we can calculate

$L(\mathbf{x}_i, \mathbf{x}_j)$ which makes the K -center more general than the K -means.

Important

For K -center clustering

We only need pairwise distance $L(\mathbf{x}_i, \mathbf{x}_j)$ for any $\mathbf{x}_i, \mathbf{x}_j \in S$.

Where

\mathbf{x}_i can be a non-vector representation of the objects.

As long we can calculate

$L(\mathbf{x}_i, \mathbf{x}_j)$ which makes the K -center more general than the K -means.

Properties

Something Notable

The greedy algorithm achieves an approximation factor of 2 as the distance measure L satisfies the triangle inequality.

That is, we have that

$$D^* = \min_S \max_{k=1, \dots, K} \max_{i, j: x_i, x_j \in C_k} L(x_i, x_j) \quad (9)$$

Then we have the following lemma about the greedy algorithm

$$D \leq 2D^* \quad (10)$$

Properties

Something Notable

The greedy algorithm achieves an approximation factor of 2 as the distance measure L satisfies the triangle inequality.

Thus, we have that

$$D^* = \min_S \max_{k=1,\dots,K} \max_{i,j:x_i,x_j \in C_k} L(\mathbf{x}_i, \mathbf{x}_j) \quad (9)$$

$$D \leq 2D^* \quad (10)$$

Properties

Something Notable

The greedy algorithm achieves an approximation factor of 2 as the distance measure L satisfies the triangle inequality.

Thus, we have that

$$D^* = \min_S \max_{k=1,\dots,K} \max_{i,j:x_i,x_j \in C_k} L(\mathbf{x}_i, \mathbf{x}_j) \quad (9)$$

Then, we have the following guarantee for the greedy algorithm

$$D \leq 2D^* \quad (10)$$

Nevertheless

We have that

K -center does not provide a locally optimal solution.

We can get only a solution

Guaranteeing to be within a certain performance range of the theoretical optimal solution.

Nevertheless

We have that

K -center does not provide a locally optimal solution.

We can get only a solution

Guaranteeing to be within a certain performance range of the theoretical optimal solution.

Outline

1

K-Means Clustering

- The NP-Hard Problem
- *K*-Means Clustering Heuristic
- Convergence Criterion
- The Distance Function
- Example
- Properties of *K*-Means

2

The *K*-Center Criterion Clustering

- Introduction
- Re-Stating the *K*-center as a Clustering Problem
- Comparison with *K*-means
- The Greedy *K*-Center Algorithm
- Pseudo-Code**
- The *K*-Center Algorithm
- Notes in Implementation
- Examples
- *K*-Center Algorithm Properties

Setup

First

Set H denotes the set of cluster centroids or cluster of representative objects $\{h_1, \dots, h_k\} \subset S$.

Second

Let $cluster(x_i)$ be the identity of the cluster $x_i \in S$ belongs to.

Third

The distance $dist(x_i)$ is the distance between x_i and its closest cluster representative object (centroid).

$$dist(x_i) = \min_{h_j \in H} L(x_i, h_j) \quad (11)$$

Setup

First

Set H denotes the set of cluster centroids or cluster of representative objects $\{h_1, \dots, h_k\} \subset S$.

Second

Let $cluster(x_i)$ be the identity of the cluster $x_i \in S$ belongs to.

Third

The distance $dist(x_i)$ is the distance between x_i and its closest cluster representative object (centroid).

$$dist(x_i) = \min_{h_j \in H} L(x_i, h_j) \quad (11)$$

Setup

First

Set H denotes the set of cluster centroids or cluster of representative objects $\{\mathbf{h}_1, \dots, \mathbf{h}_k\} \subset S$.

Second

Let $cluster(\mathbf{x}_i)$ be the identity of the cluster $\mathbf{x}_i \in S$ belongs to.

Third

The distance $dist(\mathbf{x}_i)$ is the distance between \mathbf{x}_i and its closest cluster representative object (centroid):

$$dist(\mathbf{x}_i) = \min_{\mathbf{h}_j \in H} L(\mathbf{x}_i, \mathbf{h}_j) \quad (11)$$

The Main Idea

Something Notable

We always assign x_i to the closest centroid. Therefore $dist(x_i)$ is the minimum distance between x_i and any centroid.

The Main Idea

Something Notable

We always assign x_i to the closest centroid. Therefore $dist(x_i)$ is the minimum distance between x_i and any centroid.

The Algorithm is Iterative

- We generate one cluster centroid first and then add others one by one until we get K clusters.

• The set of centroids H starts with only a single centroid, h_1 .

The Main Idea

Something Notable

We always assign x_i to the closest centroid. Therefore $dist(x_i)$ is the minimum distance between x_i and any centroid.

The Algorithm is Iterative

- We generate one cluster centroid first and then add others one by one until we get K clusters.
- The set of centroids H starts with only a single centroid, h_1 .

Outline

1

K-Means Clustering

- The NP-Hard Problem
- *K*-Means Clustering Heuristic
- Convergence Criterion
- The Distance Function
- Example
- Properties of *K*-Means

2

The *K*-Center Criterion Clustering

- Introduction
- Re-Stating the *K*-center as a Clustering Problem
- Comparison with *K*-means
- The Greedy *K*-Center Algorithm
- Pseudo-Code
- The *K*-Center Algorithm**
- Notes in Implementation
- Examples
- *K*-Center Algorithm Properties

Algorithm

Step 1

- Randomly select an object x_j from S , let $\mathbf{h}_1 = x_j$, $H = \{\mathbf{h}_1\}$.
 - It does not matter how x_j is selected.

Algorithm

Step 1

- Randomly select an object x_j from S , let $\mathbf{h}_1 = x_j$, $H = \{\mathbf{h}_1\}$.
- It does not matter how x_j is selected.

Step 2

For $j = 1$ to n :

- $dist(x_i) = L(x_i, h_1)$.
- $cluster(x_i)$.

Algorithm

Step 1

- Randomly select an object x_j from S , let $\mathbf{h}_1 = x_j$, $H = \{\mathbf{h}_1\}$.
- It does not matter how x_j is selected.

Step 2

For $j = 1$ to n :

- ➊ $dist(x_i) = L(x_i, h_1)$.
- ➋ $cluster(x_i)$.

In other words:

- The $dist(x_i)$ is the computed distance between x_i and h_1 .
- Because so far we only have one cluster, we will assign a cluster label 1 to every x_j .

Algorithm

Step 1

- Randomly select an object x_j from S , let $\mathbf{h}_1 = x_j$, $H = \{\mathbf{h}_1\}$.
- It does not matter how x_j is selected.

Step 2

For $j = 1$ to n :

- ➊ $dist(\mathbf{x}_i) = L(\mathbf{x}_i, \mathbf{h}_1)$.

➋ $cluster(\mathbf{x}_j)$.

In other words:

- The $dist(\mathbf{x}_i)$ is the computed distance between \mathbf{x}_j and \mathbf{h}_1 .
- Because so far we only have one cluster, we will assign a cluster label 1 to every \mathbf{x}_j .

Algorithm

Step 1

- Randomly select an object x_j from S , let $\mathbf{h}_1 = x_j$, $H = \{\mathbf{h}_1\}$.
- It does not matter how x_j is selected.

Step 2

For $j = 1$ to n :

- ➊ $dist(\mathbf{x}_i) = L(\mathbf{x}_i, \mathbf{h}_1)$.
- ➋ $cluster(\mathbf{x}_i)$.

In other words:

- The $dist(\mathbf{x}_i)$ is the computed distance between \mathbf{x}_j and \mathbf{h}_1 .
- Because so far we only have one cluster, we will assign a cluster label 1 to every \mathbf{x}_j .

Algorithm

Step 1

- Randomly select an object x_j from S , let $\mathbf{h}_1 = x_j$, $H = \{\mathbf{h}_1\}$.
- It does not matter how x_j is selected.

Step 2

For $j = 1$ to n :

- ① $dist(\mathbf{x}_i) = L(\mathbf{x}_i, \mathbf{h}_1)$.
- ② $cluster(\mathbf{x}_i)$.

In other words

- The $dist(x_i)$ is the computed distance between x_j and \mathbf{h}_1 .
 - Because so far we only have one cluster, we will assign a cluster label 1 to every x_i .

Algorithm

Step 1

- Randomly select an object x_j from S , let $\mathbf{h}_1 = x_j$, $H = \{\mathbf{h}_1\}$.
- It does not matter how x_j is selected.

Step 2

For $j = 1$ to n :

- ① $dist(\mathbf{x}_i) = L(\mathbf{x}_i, \mathbf{h}_1)$.
- ② $cluster(\mathbf{x}_i)$.

In other words

- The $dist(x_i)$ is the computed distance between x_j and \mathbf{h}_1 .
- Because so far we only have one cluster, we will assign a cluster label 1 to every x_j .

Algorithm

Step 3

For $i = 2$ to K

① $D = \max_{x_j: x_j \in S - H} dist(x_j)$

- ② Choose $h_i \in S - H$ such that $dist(h_i) == D$
- ③ $H = H \cup \{h_i\}$
- ④ for $j = 1$ to N
 - ⑤ if $L(x_j, h_i) \leq dist(x_j)$
 - ⑥ $dist(x_j) = L(x_j, h_i)$
 - ⑦ $cluster(x_j) = i$

Algorithm

Step 3

For $i = 2$ to K

① $D = \max_{x_j: x_j \in S - H} dist(\mathbf{x}_j)$

② Choose $\mathbf{h}_i \in S - H$ such that $dist(\mathbf{h}_i) == D$

③ $H = H \cup \{\mathbf{h}_i\}$

④ for $j = 1$ to N

⑤ if $L(\mathbf{x}_j, \mathbf{h}_i) \leq dist(\mathbf{x}_j)$

⑥ $dist(\mathbf{x}_j) = L(\mathbf{x}_j, \mathbf{h}_i)$

⑦ $cluster(\mathbf{x}_j) = i$

Algorithm

Step 3

For $i = 2$ to K

- ① $D = \max_{x_j: x_j \in S - H} dist(\mathbf{x}_j)$
- ② Choose $\mathbf{h}_i \in S - H$ such that $dist(\mathbf{h}_i) == D$
- ③ $H = H \cup \{\mathbf{h}_i\}$
- ④ for $j = 1$ to N
 - ⑤ if $L(\mathbf{x}_j, \mathbf{h}_i) \leq dist(\mathbf{x}_j)$
 - ⑥ $dist(\mathbf{x}_j) = L(\mathbf{x}_j, \mathbf{h}_i)$
 - ⑦ $cluster(\mathbf{x}_j) = i$

Algorithm

Step 3

For $i = 2$ to K

- ① $D = \max_{x_j: x_j \in S - H} dist(\mathbf{x}_j)$
- ② Choose $\mathbf{h}_i \in S - H$ such that $dist(\mathbf{h}_i) == D$
- ③ $H = H \cup \{\mathbf{h}_i\}$
- ④ for $j = 1$ to N
 - ⑤ if $L(x_j, \mathbf{h}_i) \leq dist(\mathbf{x}_j)$
 - ⑥ $dist(\mathbf{x}_j) = L(x_j, \mathbf{h}_i)$
 - ⑦ $cluster(\mathbf{x}_j) = i$

Algorithm

Step 3

For $i = 2$ to K

- ① $D = \max_{x_j: x_j \in S - H} dist(\mathbf{x}_j)$
- ② Choose $\mathbf{h}_i \in S - H$ such that $dist(\mathbf{h}_i) == D$
- ③ $H = H \cup \{\mathbf{h}_i\}$
- ④ for $j = 1$ to N
 - ⑤ if $L(\mathbf{x}_j, \mathbf{h}_i) \leq dist(\mathbf{x}_j)$
 - ⑥ $dist(\mathbf{x}_j) = L(\mathbf{x}_j, \mathbf{h}_i)$
 - ⑦ $cluster(\mathbf{x}_j) = i$

Algorithm

Step 3

For $i = 2$ to K

- ① $D = \max_{x_j: x_j \in S - H} dist(\mathbf{x}_j)$
- ② Choose $\mathbf{h}_i \in S - H$ such that $dist(\mathbf{h}_i) == D$
- ③ $H = H \cup \{\mathbf{h}_i\}$
- ④ for $j = 1$ to N
- ⑤ if $L(\mathbf{x}_j, \mathbf{h}_i) \leq dist(\mathbf{x}_j)$
- ⑥ $dist(\mathbf{x}_j) = L(\mathbf{x}_j, \mathbf{h}_i)$

Distance (\mathbf{x}_j) = ?

Algorithm

Step 3

For $i = 2$ to K

- ① $D = \max_{x_j: x_j \in S - H} dist(\mathbf{x}_j)$
- ② Choose $\mathbf{h}_i \in S - H$ such that $dist(\mathbf{h}_i) == D$
- ③ $H = H \cup \{\mathbf{h}_i\}$
- ④ for $j = 1$ to N
- ⑤ if $L(\mathbf{x}_j, \mathbf{h}_i) \leq dist(\mathbf{x}_j)$
- ⑥ $dist(\mathbf{x}_j) = L(\mathbf{x}_j, \mathbf{h}_i)$
- ⑦ $cluster(\mathbf{x}_j) = i$

Thus

As the algorithm progresses

- We gradually add more and more cluster centroids, beginning with 2 until we get to K .
- At each iteration, we find among all of the points which are not yet included in the set, a worst point:
 - ▶ Worst in the sense that this point has maximum distance to its corresponding centroid.

Thus

As the algorithm progresses

- We gradually add more and more cluster centroids, beginning with 2 until we get to K .
- At each iteration, we find among all of the points which are not yet included in the set, a worst point:
 - ▶ Worst in the sense that this point has maximum distance to its corresponding centroid.

What does it do?

- It is added to the set H .
- To stress gain, points already included in H are not among the consideration.

Thus

As the algorithm progresses

- We gradually add more and more cluster centroids, beginning with 2 until we get to K .
- At each iteration, we find among all of the points which are not yet included in the set, a worst point:
 - ▶ Worst in the sense that this point has maximum distance to its corresponding centroid.

This worst point

- It is added to the set H .
 - ▶ To stress gain, points already included in H are not among the consideration.

Thus

As the algorithm progresses

- We gradually add more and more cluster centroids, beginning with 2 until we get to K .
- At each iteration, we find among all of the points which are not yet included in the set, a worst point:
 - ▶ Worst in the sense that this point has maximum distance to its corresponding centroid.

This worst point

- It is added to the set H .
- To stress gain, points already included in H are not among the consideration.

Outline

1

K-Means Clustering

- The NP-Hard Problem
- *K*-Means Clustering Heuristic
- Convergence Criterion
- The Distance Function
- Example
- Properties of *K*-Means

2

The *K*-Center Criterion Clustering

- Introduction
- Re-Stating the *K*-center as a Clustering Problem
- Comparison with *K*-means
- The Greedy *K*-Center Algorithm
- Pseudo-Code
- The *K*-Center Algorithm
- Notes in Implementation
- Examples
- *K*-Center Algorithm Properties

Implementation

We can use the following for implementation

- The disjoint-set data structure with the following operations:
 - ▶ MakeSet
 - ▶ Find
 - ▶ Union
 - ▶ Remove iteratively through the disjoint trees.

Implementation

We can use the following for implementation

- The disjoint-set data structure with the following operations:
 - ▶ MakeSet
 - ▶ Find
 - ▶ Union
 - ▶ Remove iteratively through the disjoint trees.

Of course it is necessary to have extra fields for the necessary distances.

Thus, each node-element for the sets must have a field $dist(x_j)$ such that

$$dist(x_j) = L(x_j, \text{Find}(x_j)) \quad (12)$$

Implementation

We can use the following for implementation

- The disjoint-set data structure with the following operations:
 - ▶ MakeSet
 - ▶ Find
 - ▶ Union
 - ▶ Remove iteratively through the disjoint trees.

Of course it is necessary to have extra fields for the necessary distances:

Thus, each node-element for the sets must have a field $dist(x_j)$ such that

$$dist(x_j) = L(x_j, \text{Find}(x_j)) \quad (12)$$

Implementation

We can use the following for implementation

- The disjoint-set data structure with the following operations:
 - ▶ MakeSet
 - ▶ Find
 - ▶ Union
- ⇒ Remove iteratively through the disjoint trees.

Of course it is necessary to have extra fields for the necessary distances:

Thus, each node-element for the sets must have a field $dist(x_j)$ such that

$$dist(x_j) = L(x_j, \text{Find}(x_j)) \quad (12)$$

Implementation

We can use the following for implementation

- The disjoint-set data structure with the following operations:
 - ▶ MakeSet
 - ▶ Find
 - ▶ Union
 - ▶ Remove Iteratively through the disjoint trees.

Of course it is necessary to have extra fields for the necessary distances.

Thus, each node-element for the sets must have a field $dist(x_j)$ such that

$$dist(x_j) = L(x_j, \text{Find}(x_j)) \quad (12)$$

Implementation

We can use the following for implementation

- The disjoint-set data structure with the following operations:
 - ▶ MakeSet
 - ▶ Find
 - ▶ Union
 - ▶ Remove Iteratively through the disjoint trees.

Of course it is necessary to have extra fields for the necessary distances

Thus, each node-element for the sets must have a field $dist(\mathbf{x}_j)$ such that

$$dist(\mathbf{x}_j) = L(\mathbf{x}_j, Find(\mathbf{x}_j)) \quad (12)$$

Although

Other things need to be taken in consideration

I will allow to you to think about them

Outline

1

K-Means Clustering

- The NP-Hard Problem
- *K*-Means Clustering Heuristic
- Convergence Criterion
- The Distance Function
- Example
- Properties of *K*-Means

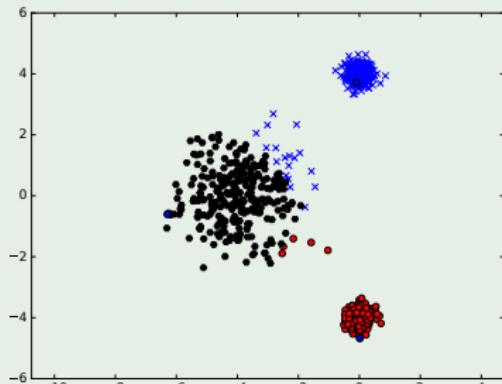
2

The *K*-Center Criterion Clustering

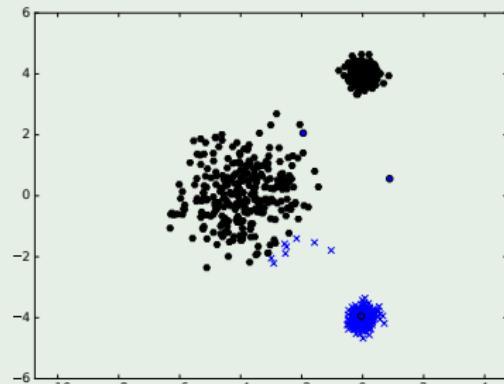
- Introduction
- Re-Statting the *K*-center as a Clustering Problem
- Comparison with *K*-means
- The Greedy *K*-Center Algorithm
- Pseudo-Code
- The *K*-Center Algorithm
- Notes in Implementation
- Examples**
- *K*-Center Algorithm Properties

Example

Running the k-center and k-means algorithms allows to see that for different densities k-center is more robust



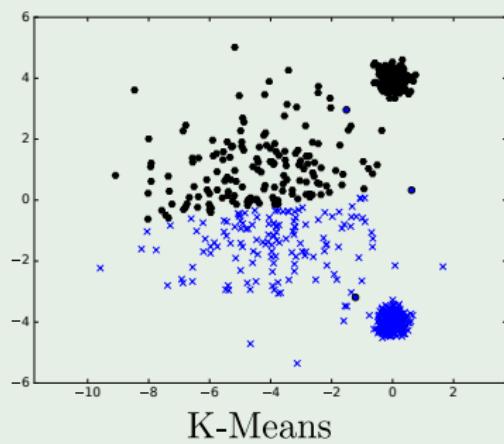
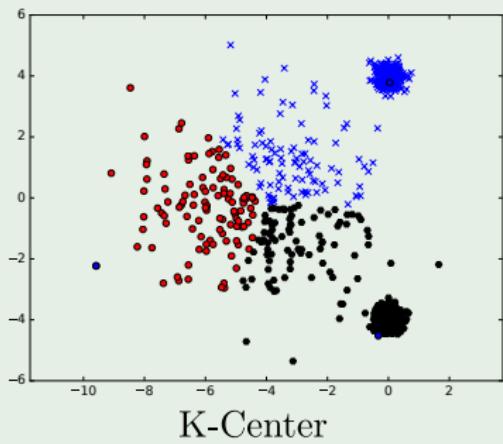
K-Center



K-Means

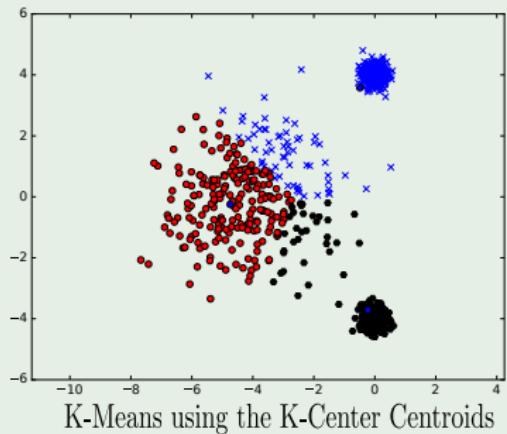
Example

Decreasing the density of one of the clusters, we see a degradation on the clusters

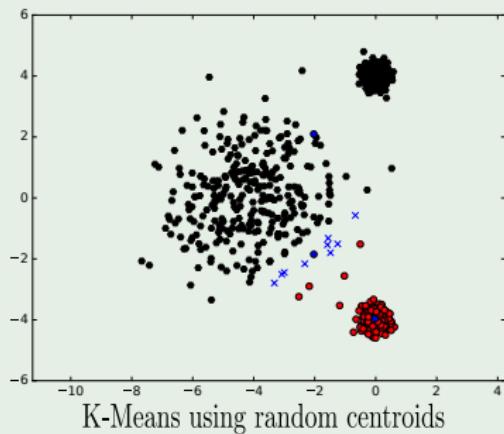


Using Centroids of the K-center to initialize K-mean

Thus, we can use the centroids of K-center to try to improve upon K-means to a certain degree



K-Means using the K-Center Centroids



K-Means using random centroids

Outline

1

K-Means Clustering

- The NP-Hard Problem
- *K*-Means Clustering Heuristic
- Convergence Criterion
- The Distance Function
- Example
- Properties of *K*-Means

2

The *K*-Center Criterion Clustering

- Introduction
- Re-Stating the *K*-center as a Clustering Problem
- Comparison with *K*-means
- The Greedy *K*-Center Algorithm
- Pseudo-Code
- The *K*-Center Algorithm
- Notes in Implementation
- Examples
- ***K*-Center Algorithm Properties**

The Running Time

We have that

The running time of the algorithm is $O(KN)$, where K is the number of clusters generated and N is the size of the data set.



Because K -center only requires pairwise distance between any point and the centroids.

The Running Time

We have that

The running time of the algorithm is $O(KN)$, where K is the number of clusters generated and N is the size of the data set.

Why?

Because K -center only requires pairwise distance between any point and the centroids.