

Lecture Notes on k -Means Clustering (I)

Christian Bauckhage

B-IT, University of Bonn

This is the first in a series of lecture notes on k -means clustering, its variants, and applications. In this note, we study basic ideas behind k -means clustering and identify common pitfalls in its use.

k -Means Clustering

Although the k -means clustering algorithm is frequently applied in practice, it seems that many users are not familiar with the theory behind it. This is unfortunate, because the algorithm operates on certain implicit assumptions which, if ignored, may lead to seemingly unreasonable results. We analyze these principles, discuss when and how to apply k -means, and point out common pitfalls in its use.

FOR THE TIME BEING, we focus on the problem of clustering data that can be embedded in Euclidean vector spaces; extensions to more abstract settings will be studied later on.

Accordingly, let us assume a data set $X = \{x_1, x_2, \dots, x_n\} \subset \mathbb{R}^m$ whose n elements x_j are m dimensional real-valued vectors that are to be clustered (see Fig. 1 for an example).

WHENEVER WE SET OUT TO CLUSTER DATA, we are interested in detecting latent structures or regularities within a given sample. In practice, we must therefore decide how to characterize structures within a set of data and how to search for them?

The famous k -means algorithm arguably provides the most popular answers to these questions. It represents structures in terms of *subsets* of X and attempts to divide the data into k different *clusters* $C_i \subset X$ which meet the following criteria:

1. clusters should be pairwise disjoint, i.e. $C_i \cap C_j = \emptyset$ for $i \neq j$
2. clusters should cover the data such that $C_1 \cup C_2 \cup \dots \cup C_k = X$
3. most importantly, data assigned to a cluster C_i should be *similar*.

While the first two criteria are well defined, the third one hinges on the rather intuitive notion of similarity. Indeed, there are many possible definitions of similarity and characteristic differences between clustering algorithms usually trace back to different notions of what it means to be similar.

The idea employed in k -means clustering is to represent each cluster C_i by means of a *centroid* $\mu_i \in \mathbb{R}^m$ and to measure similarities in terms of Euclidean distances between centroids and data points. Hence, two data points are considered similar if their distances to a common centroid are small and a point x_j will be assigned to cluster C_i if its distance to μ_i is less than that to any other centroid.

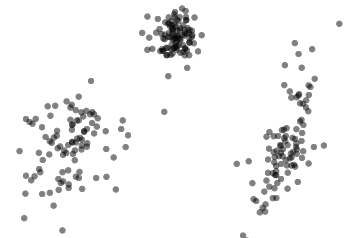


Figure 1: Didactic example of $n = 150$ data points $x_j \in \mathbb{R}^2$ sampled from three bivariate Gaussian distributions.



cluster centroids

THIS NOTION OF SIMILARITY reduces the problem of clustering to the problem of finding appropriate centroids. This, in turn, can be cast as the task of minimizing the following objective function

$$E(k) = \sum_{i=1}^k \sum_{x_j \in C_i} \|x_j - \mu_i\|^2 \quad (1)$$

where the minimization has to be carried out w.r.t. k centroids μ_i . Thus, minimizing (1) is to determine suitable centroids μ_i such that, if the data are partitioned into corresponding clusters C_i , distances between data points and their closest cluster centroid become as small as possible.

Already at this point, we note that the objective in (1) can also be expressed in terms of *indicator variables* z_{ij} which register for any x_j whether it belongs to cluster C_i . That is, by defining

$$z_{ij} = \begin{cases} 1, & \text{if } x_j \in C_i \\ 0, & \text{otherwise} \end{cases} \quad (2)$$

we find the above objective function to be equivalent to

$$E(k) = \sum_{i=1}^k \sum_{j=1}^n z_{ij} \|x_j - \mu_i\|^2. \quad (3)$$

ALTHOUGH BOTH OBJECTIVES LOOK RATHER INNOCENT and formalize an intuitive problem, finding an optimal solution, i.e. finding global minimizers μ_1, \dots, μ_k , is actually NP-hard¹. Even for 2D data and settings where $k = 2$, it cannot be guaranteed that the best possible solution will be found in reasonable time. This is important to know! It tells us that **the k -means algorithm is merely a heuristic for dealing with a surprisingly hard problem.**

The Classical k -means Algorithm

So, how does the k -means algorithm minimize the above objectives? In a nutshell, it realizes a greedy iterative update scheme. When started, say at iteration $t = 0$, it randomly initializes the parameters $\mu_1^{(t)}, \mu_2^{(t)}, \dots, \mu_k^{(t)}$. Given this initial guess, the data are then clustered accordingly. That is, the algorithm determines k clusters where

$$C_i^{(t)} = \left\{ x_j \in X \mid \|x_j - \mu_i^{(t)}\|^2 \leq \|x_j - \mu_l^{(t)}\|^2 \forall l \neq i \right\} \quad (4)$$

Once clusters have been determined, the algorithm updates the current estimates of the cluster centroids by computing

$$\mu_i^{(t+1)} = \frac{1}{n_i} \sum_{x_j \in C_i^{(t)}} x_j \quad (5)$$

where $n_i = |C_i^{(t)}|$ denotes the number of elements in cluster C_i . As these updates may lead to a new clustering, the iteration counter is increased to $t = t + 1$ and steps (4) and (5) are repeated until the



indicator variables

¹ D. Aloise, A. Deshpande, P. Hansen, and P. Popat. NP-Hardness of Euclidean Sum-of-Squares Clustering. *Machine Learning*, 75(2), 2009

```

t ← 0
done ← False
randomly initialize  $\mu_1^t, \mu_2^t, \dots, \mu_k^t$ 
while not done do
    // for each mean, update its cluster
    for i = 1 to k do
         $C_i^t = \{x \in X \mid \|x - \mu_i^t\| \leq \|x - \mu_j^t\|\}$ 
    if  $C_i^t \cap C_i^{t-1} = C_i^t \forall i$  then
        done ← True

    // for each cluster, update its mean
    for i = 1 to k do
         $n_i = |C_i^t|$ 
         $\mu_i^{t+1} = \frac{1}{n_i} \sum_{x \in C_i^t} x$ 
    if  $\|\mu_i^{t+1} - \mu_i^t\| \leq \epsilon \forall i$  then
        done ← True

    t ← t + 1
    if t >  $t_{\max}$  then
        done ← True

end while

```

Figure 2: Pseudo code for k -means.

assignment of data points to clusters does not change anymore or t exceeds a predefined number of iterations (see Fig. 3 for an example).

IT IS EASY TO SHOW that updating each centroid to the *sample mean* of its cluster will never increase the value of objectives in (1) or (3)². Each iteration therefore improves on the previous result and k -means clustering usually converges quickly.

However, as the algorithm starts from a random initialization, it cannot be guaranteed to find the global minimum of its objective function. As it typically converges to a local minimum, **it is pivotal to run k -means several times so as to empirically determine good clusters**. Yet, even then, k -means may produce questionable results when seeded with inappropriate centroids. Intelligent initializations are the topic of ongoing research^{3,4} and **arbitrary initialization are considered harmful; initial centroids should at least be selected among the available data points**.

Common Pitfalls in Practical Applications

Before we continue to study further theoretical aspects of k -means clustering, it seems instructive to discuss when and how to apply the k -means algorithm in practice.

It can be shown that the k -means algorithm is tailored to locally Gaussian data⁵. Hence, **if k -means clustering is applied to data sets that are not composed of compact convex subsets, it is overly optimistic to expect the algorithm to identify reasonable clusters**.

The following example illustrates what this may mean in practice. Figure 4(a) displays a set of 2D data which apparently consists of two distinct, linear clusters. Yet, if we set $k = 2$ and run k -means, the result will likely look like in Fig. 4(b). In particular, cluster centroids (\square) will be rather far from actual data points but form the centers of two imaginary circles which each contain about half the data. This practical result therefore illustrates the tendency of the k -means algorithm to produce compact, blob-like clusters. In other words, for data such as in this example, the algorithm will easily determine clusters but its results will hardly agree with what is salient to humans.

AS A REMEDY, the literature frequently suggests to normalize the

² We will prove this in a later note.

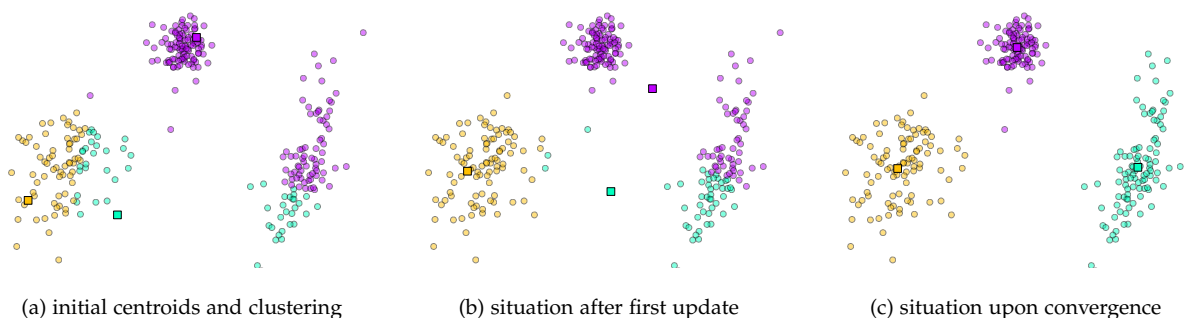
³ P.S. Bradley and U.M. Fayyad. Refining Initial Points for K-Means Clustering. In *Proc. ICML*, 1998

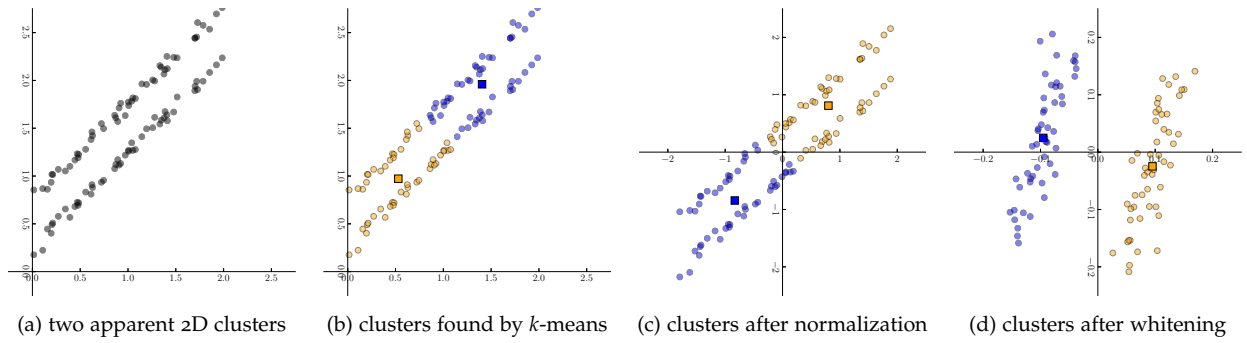
⁴ D. Arthur and S. Vassilvitskii. k -means++: The Advantages of Careful Seeding. In *Proc. SODA*, 2007

⁵ We will do this in a later note.

Figure 3: k -means clustering applied to determine $k = 3$ clusters in the data (\circ) from Fig. 1; (a) initial centroids (\square) and corresponding clustering, (b) situation after the first update, and (c) result upon convergence.

As k -means is tailored to Gaussian mixtures and since the initial choice of centroids was favorable, it took only four iterations to converge to this reasonable solution. However, there is no guarantee for k -means to always behave this way. In practice, one should therefore run it several times and base any further analysis on the result with the lowest value of the objective in Eq. (1).





data using $y_{dj} = (x_{dj} - \mu_d) / \sigma_d$ where x_{dj} is the d th component of vector x_j and μ_d and σ_d^2 denote the data's mean and variance along dimension d . Afterwards, the data will be of zero mean and unit variance and are commonly believed to be more amenable to k -means. But Fig. 4(c) shows that this is not necessarily the case. For our example, the transformation is too naïve to have a noticeable effect on the result; the data are simply not Gaussian enough for k -means to identify their “true” latent substructures.

THEREFORE, if we insist on applying the k -means algorithm to non-Gaussian data, we need better preprocessing. A reasonable idea is to *whiten* the data.

Computing the sample mean $\mu = \frac{1}{n} \sum_j x_j$ and transforming the data to zero-mean $y_j = x_j - \mu$ yields the covariance matrix $\Sigma = YY^T$. Using principal component analysis, Σ can be decorrelated so that $\Lambda = U^T \Sigma U$. Transformed vectors $z_j = U^T y_j$ thus have a diagonal covariance matrix Λ . Still, variances along different dimensions will usually differ; but since $\Lambda^{-1/2} \Lambda \Lambda^{-1/2} = I$, the transformation

$$w_j = \Lambda^{-1/2} U^T (x_j - \mu) \quad (6)$$

maps the data to vectors w_j for which $\frac{1}{n} \sum_j w_j w_j^T = I$. That is, after whitening, variances in the directions of the principal axes of the data are identically 1. Geometrically, we can picture this transformation as mapping the data (close) to the surface of a sphere where possible clusters will be more separated than in the original representation.

Figure 4(d) shows the whitened data and illustrates the corresponding effect on k -means clustering. Based on this example, it appears that **proper preprocessing enhances the chances for the k -means algorithm to identify structures that are not necessarily Gaussian**. Nevertheless, there is still no guarantee for it to do so.

THERE IS YET ANOTHER PATHOLOGICAL SITUATION, namely the case of very high dimensional data where we are given n data points $x_j \in \mathbb{R}^m$ but $n \ll m$.

For *high dimension low sample size data* like these, the notion of similarity invoked by k -means may be useless, because as dimensions grow so do distances and it becomes very likely that all the given

Figure 4: Example of the limitations of k -means clustering. As it implicitly assumes locally Gaussian data, the algorithm clusters data into compact convex subsets.

The obvious linear substructures in this example can not be uncovered by k -means. Generally, data needs to be whitened so that the algorithm stands a better chance of producing results close to human intuition. However, there is no guarantee for k -means to succeed in this regard.



data whitening



HDLSS data

data are equally far apart^{6,7,8}. Hence, centroids determined by k -means clustering will be far from any data point and the decision of whether to assign a data point to a cluster will basically depend on minuscule random variations rather than on structural properties. Therefore, **for very high dimensional data, k -means clustering should be applied only after dimensionality reduction.**

Notes and References

The origins of the idea of k -means clustering can be traced back to work by Steinhaus⁹ in the 1950s. The standard algorithm, that is the iterative algorithm discussed in this note, was first discussed by Lloyd in a Bell Labs technical report also in the 1950s but openly published only much later¹⁰. Finally, the now so commonly used term “ k -means” was coined by MacQueen¹¹ in the 1960s.

For readers interested in deeper expositions than the ones in this and our subsequent notes on the topic, we recommend the rigorous and thorough accounts by Bishop¹² or MacKay¹³.

⁶ D.L. Donoho and J. Tanner. Neighborliness of Randomly Projected Simplices in High Dimensions. *PNAS*, 102(27), 2005

⁷ P. Hall, J.S. Marron, and A. Neeman. Geometric Representations of High Dimension, Low Sample Size Data. *J. Royal Statistical Society B*, 67(3), 2005

⁸ F. Murtagh. The Remarkable Simplicity of Very High Dimensional Data: Applications of Model-based Clustering. *J. of Classification*, 26(3), 2009

⁹ H. Steinhaus. Sur la Division des Corps Matériels en Parties. *Bulletin Polish Acad. of Sciences*, 4(12):801–804, 1957

¹⁰ S.P. Lloyd. Least Squares Quantization in PCM. *IEEE Trans. on Information Theory*, 28(2), 1982

¹¹ J.B. MacQueen. Some Methods for Classification and Analysis of Multivariate Observations. In *Proc. Berkeley Symp. on Mathematical Statistics and Probability*, 1967

¹² C.M. Bishop. *Pattern Recognition and Machine Learning*. Springer, 2006

¹³ D.J.C. MacKay. *Information Theory, Inference, and Learning Algorithms*. Cambridge University Press, 2003