# Lecture Notes on Data Science: Soft k-Means Clustering

1 author:

Christian Bauckhage
University of Bonn
**366** PUBLICATIONS   **5,640** CITATIONS

Some of the authors of this publication are also working on these related projects:

Project    lectures on pattern recognition View project

Project    P3ML - ML Engineering Knowledge View project

# Lecture Notes on Data Science: Soft k-Means Clustering

*Christian Bauckhage*

*B-IT, University of Bonn*

In this note, we study the idea of soft *k*-means clustering which yields soft assignments of data points to clusters. We discuss the basic theory, a simple algorithm, and examples for how this algorithm behaves.

## Introduction

So far, our discussion of *k*-means clustering and its properties has been restricted to the paradigm of *hard k-means clustering*.

In other words, we considered the problem of partitioning a set

$$X = \{x_1, x_2, \ldots, x_n\} \subset \mathbb{R}^m \tag{1}$$

of data points $x_j$ into $k$ clusters $C_1, \ldots, C_k$ for which we required that they are proper subsets $C_i \subset X$ such that $C_1 \cup C_2 \cup \ldots \cup C_k = X$ and

$$C_i \cap C_l = \varnothing \tag{2}$$

for $i \neq l$. Given these requirements, we then derived the classical *k*-means procedure[1] as a method for computing *hard cluster assignments*[2] where

$$\forall\, i \neq l : x_j \in C_i \Rightarrow x_j \notin C_l. \tag{3}$$

QUITE OFTEN, HOWEVER, hard cluster assignments may be overzealous. Consider, for example, Figure 1. While most everybody would agree that this data contains three distinct clusters, the two on the left appear to overlap in that there are a few data points for which it is not quite so obvious to which cluster they belong.

IN THIS NOTE, we will therefore relax the assumption in (2) and study the idea of *soft k-means clustering*.

There are two fundamental innovations with regard to this idea. First of all, soft *k*-means clustering allows every data point to belong to several clusters simultaneously. Second of all, this is not meant in the simple sense that clusters may overlap. Rather, soft *k*-means clustering determines to which *degree* each data point belongs to each cluster. This is to say, that the resulting clusters are fuzzy rather than crisp sets so that soft *k*-means clustering is sometimes also referred to as *fuzzy k-means clustering*[3].

## Soft k-Means Clustering

Recall that, when we studied the relationship between (hard) *k*-means clustering and Gaussian mixture modeling[4], we saw that clustering is essentially a minimization problem. In particular, we saw that
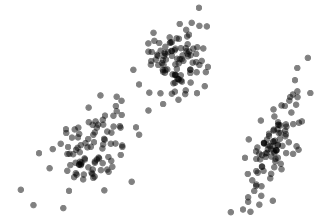


Figure 1: Didactic example of $n = 300$ data points $x_j \in \mathbb{R}^2$ sampled from three bivariate Gaussian distributions.

[1] C. Bauckhage. Lecture Notes on Data Science: *k*-Means Clustering, 2015b. DOI: 10.13140/RG.2.1.2829.4886

[2] While the logical expression in (3) may look daunting, it is just a fancy way of saying that classical *k* means clustering assigns each data point to exactly one cluster.

[3] In fact, there seems to be some controversy as to whether or not soft *k*-means and fuzzy *k*-means clustering really are the same. We do not partake in this silly debate but consider them equivalent and explain why in our discussion below.

[4] C. Bauckhage. Lecture Notes on Data Science: *k*-Means Clustering Is Gaussian Mixture Modeling, 2015a. DOI: 10.13140/RG.2.1.3033.2646

$k$-means clustering can be formalized as the problem of finding $k$ cluster centroids $\mu_1, \mu_2, \ldots, \mu_k$ which minimize the error function

$$E(\mu_1, \mu_2, \ldots, \mu_k) = \sum_{i=1}^{k} E(\mu_i) = \sum_{i=1}^{k} \sum_{j=1}^{n} z_{ij} \left\| x_j - \mu_i \right\|^2 \qquad (4)$$

where the $z_{ij}$ are ***indicator variables***

$$z_{ij} = \begin{cases} 1, & \text{if } x_j \in C_i \\ 0, & \text{otherwise.} \end{cases} \qquad (5)$$

which indicate whether or not $x_j$ belongs to cluster $C_i$.

Looking at (5), the indicator variables seem not to depend on the cluster centroids $\mu_i$. Yet, if we recall that the classical $k$-means procedure assigns data point $x_j$ to cluster $C_i$ if and only if $\mu_i$ is the closest centroid, we realize that (5) is equivalent to

$$z_{ij} = \begin{cases} 1, & \text{if } \left\| x_j - \mu_i \right\|^2 \leq \left\| x_j - \mu_l \right\|^2 \ \forall i \neq l \\ 0, & \text{otherwise.} \end{cases} \qquad (6)$$

This, however, is to say that the values of the $z_{ij}$ do indeed depend on distances between data points and centroids.

THE FACT THAT the $z_{ij}$ in (6) are binary, i.e. $z_{ij} \in \{0, 1\}$, once again reflects the crisp nature of the clusters produced by hard $k$-means clustering. But what if we would allow the $z_{ij}$ to vary continuously? That is, what if there were degrees of cluster membership $z_{ij} \in [0, 1]$? And how could we reasonably define these?

*Soft Cluster Assignments*

Given our discussion so far, it appears reasonable to decide cluster membership based on distances to cluster centroids. Hence, given a set of centroids $\mu_1, \mu_2, \ldots, \mu_k$, the problem of computing degrees of cluster membership can be understood as the problem of turning distances $\left\| x_j - \mu_i \right\|^2$ into numbers $z_{ij} \in [0, 1]$.

A POPULAR WAY OF ACHIEVING SUCH A TRANSFORMATION is to consider softmax functions[5].

A corresponding, rather common choice in the context of soft $k$-means clustering is

$$z_{ij} = \frac{e^{-\beta \left\| x_j - \mu_i \right\|^2}}{\sum_{l=1}^{k} e^{-\beta \left\| x_j - \mu_l \right\|^2}} \qquad (7)$$

where the parameter $\beta > 0$ is called the ***stiffness parameter***.

GIVEN THIS PARTICULAR DEFINITION of the cluster membership degrees $z_{ij}$, we observe that the numerator in (7) guarantees that $0 < z_{ij} \leq 1$ and the denominator in (7) guarantees that $\sum_i z_{ij} = 1$. But this is to say that the $z_{ij}$ can be interpreted as probabilities $p(\mu_i \mid x_j)$ which may be desirable in practice.

[5] Softmax functions are a staple of machine learning and frequently found in neural computing algorithms.

⚠️
**softmax function**

⚠️
**stiffness parameter**

⚠️
**probabilistic interpretation**

*Centroid Updates*

Having devised a way of computing soft rather than hard cluster membership indicators, we next need to address the question of how to adapt the updating of cluster centroids $\mu_i$ to the situation of soft $k$-means clustering?

In order to answer this question, let us single out a specific cluster, say, $C_l$ with corresponding centroid $\mu_l$. Looking at (4), we note that this cluster contributes a term of

$$E(\mu_l) = \sum_{j=1}^{n} z_{lj} \left\| x_j - \mu_l \right\|^2 \tag{8}$$

to the overall error we want to minimize. Hence, if we could find a value for $\mu_l$ such that this contribution was as small as possible, the overall error would decrease. In other words, we are dealing with the task of having to solve the following optimization problem

$$\mu_l = \operatorname*{argmin}_{\mu} \sum_{j=1}^{n} z_{lj} \left\| x_j - \mu \right\|^2. \tag{9}$$

Luckily, this is a simple problem because the sum on the right hand side of (9) is a convex function in $\mu$. It therefore has a unique minimizer which corresponds to the root of the corresponding gradient.

So let us proceed and find that root. To begin with, we note that

$$E(\mu) = \sum_{j=1}^{n} z_{lj} \left\| x_j - \mu \right\|^2 \tag{10}$$

$$= \sum_{j=1}^{n} z_{lj} \left( x_j - \mu \right)^T \left( x_j - \mu \right) \tag{11}$$

$$= \sum_{j=1}^{n} z_{lj} \left( x_j^T x_j - 2 x_j^T \mu + \mu^T \mu \right). \tag{12}$$

Hence, if we derive[6] $E(\mu)$ with respect to $\mu$, we find

$$\frac{\partial}{\partial \mu} E(\mu) = \sum_{j=1}^{n} z_{lj} \frac{\partial}{\partial \mu} \left( x_j^T x_j - 2 x_j^T \mu + \mu^T \mu \right) = \sum_{j=1}^{n} z_{lj} \left( 2\mu - 2x_j \right) \tag{13}$$

so that equating (13) to zero provides

$$\sum_{j=1}^{n} z_{lj} \, 2 \, \mu = \sum_{j=1}^{n} z_{lj} \, 2 \, x_j \tag{14}$$

$$\Leftrightarrow \qquad \mu = \frac{\sum_{j=1}^{n} z_{lj} \, x_j}{\sum_{j=1}^{n} z_{lj}} \tag{15}$$

as the solution to (9). In other words, setting $\mu_l$ to the weighted mean

$$\mu_l = \frac{\sum_{j=1}^{n} z_{lj} \, x_j}{\sum_{j=1}^{n} z_{lj}} \tag{16}$$

will minimize the contribution of cluster $C_l$ to the overall error in (4).

[6] Recall the following rules from vector calculus

$$\frac{\partial}{\partial v} u^T v = u$$

and

$$\frac{\partial}{\partial v} v^T v = 2v.$$
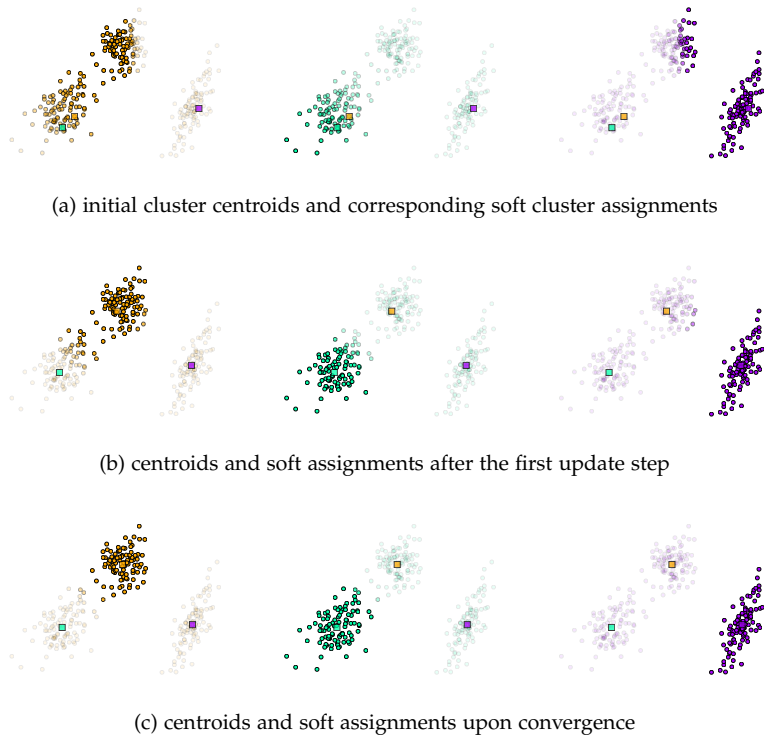
⚠ **weighted mean**

(a) initial cluster centroids and corresponding soft cluster assignments



(b) centroids and soft assignments after the first update step



(c) centroids and soft assignments upon convergence

Figure 2:   Soft $k$-means clustering for $k = 3$ of the data ($\circ$) in Fig. 1; (a) initial centroids ($\square$) and cluster assignments, (b) situation after the first update, and (c) result upon convergence.

In this example, the procedure in Fig. 3 converged within four iterations to a good solution. However, the caveats as to potentially questionable results we pointed out in our earlier discussions of the $k$-means procedure generally also apply to soft $k$-means clustering.

A video that visualizes the process shown in this figure can be found here: www.youtube.com/watch?v=Np9VuEg_aqo

## Iterative Algorithm and Examples

Now that we know how to compute cluster assignments and cluster centroids under a soft clustering paradigm, it is easy to guess how the classical procedure (or Lloyd's algorithm) for hard $k$-means can be adapted to our current setting. Hence, without further ado: Fig. 3 presents pseudo code for soft $k$-means clustering.

FIGURE 2 ATTEMPTS TO VISUALIZE the behavior of this algorithm. Given the data in Fig. 1, we set the stiffness parameter $\beta$ in (7) to 1.5, randomly initialized 3 cluster centroids, and ran soft $k$-means. For each cluster, Fig. 2 shows cluster centroids and cluster memberships during different iterations. Note that membership degrees are indicated in terms of shaded colors, the more opaque a data point is plotted, the more it belongs to the corresponding cluster[7].

In this simple example, it took the algorithm only four iterations to move the cluster centroids to positions which most human observers would agree make sense. Generally, soft $k$-means clustering will be more stable than hard $k$-means clustering since the idea of degrees of cluster memberships causes the algorithm to optimize more globally rather than just locally. Yet, the performance of soft $k$-means, too, will depend on the initialization of the centroids and there are still no guarantees for the algorithm to find a globally optimal solution. Moreover, result obtained from soft $k$-means clustering will depend on the choice of the stiffness paramtere $\beta$.

FIGURE 4 PROVIDES AN IMPRESSION for how the choice of $\beta$ may

---

$t \leftarrow 0$
randomly initialize centroids $\mu_1^t, \mu_2^t, \ldots, \mu_k^t$
**repeat**
    // compute all indicator variables
$$z_{ij} = \frac{\exp\left[-\beta \left\|x_j - \mu_i^t\right\|^2\right]}{\sum_l \exp\left[-\beta \left\|x_j - \mu_l^t\right\|^2\right]}$$

    // update all centroids
$$\mu_i^{t+1} = \frac{\sum_j z_{ij} x_j}{\sum_j z_{ij}}$$

    $t \leftarrow t + 1$
**until** centroids stabilize

Figure 3: Pseudo code for soft $k$-means.

[7] Also note that we cheated a bit when creating these plots; colors for low membership values are exaggerated for better visibility.

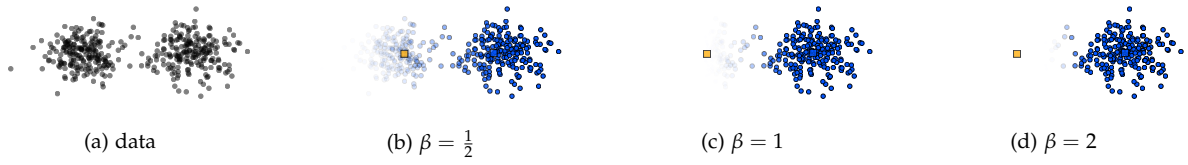(a) data      (b) $\beta = \frac{1}{2}$      (c) $\beta = 1$      (d) $\beta = 2$

Figure 4: Examples of how the choice of the stiffness parameter $\beta$ in (7) may impact soft $k$-means clustering.

influence the overall result of soft $k$-means clustering. In order to make sense of these results, we note that, w.r.t. each cluster $C_i$, the $z_{ij}$ in (7) essentially follow a Gaussian distribution so that $\beta$ can be seen as an inverse variance. Hence, larger values of $\beta$ will cause the Gaussian to be narrower so that degrees of membership of data points $x_j$ far from the centroid $\mu_i$ will become small[8].

[8] The obvious question then is how to choose $\beta$ in practice? When in doubt, $\beta = 1$ is always a pragmatic choice. A generally better idea would be to set it as a multiple of the average distance between the given data points.

## Notes and References

Before we conclude, we feel it is warranted to address the issue of "soft vs. fuzzy $k$-means" which seems to cause confusion in blogs and forums on the Web. To be specific, the debate revolves around the question of whether or not soft $k$-means clustering as discussed in this note and fuzzy $k$-means clustering are the same?

Note that the defining feature of fuzzy computing paradigms is that they are concerned with fuzzy sets. A fuzzy set $S$ is a set whose elements $s$ have degrees of membership often written as $m_S(s)$ where $m_S$ is a function $m_S : S \rightarrow [0, 1]$.

Looking at (7), we observe once again that the degrees of cluster membership defined there are values $z_{ij} \in [0, 1]$. In fact, each $z_{ij}$ is determined w.r.t. a cluster (a set) $C_i$ and an element $x_j$. Except for notational clutter, there is therefore no reason that would prevent us from writing $z_{C_i}(x_j)$ instead of $z_{ij}$.

In other words, the probabilistic cluster membership indicators in (7) constitute but a specific choice of a fuzzy membership function. We could have chosen other functions such as triangle functions or trapezoids that are common in fuzzy logic. We could have also chosen a different function for each cluster. But, in the end, (7) can be understood as fuzzy membership function. Hence, soft $k$-means clustering as discussed above is but a specific instance of fuzzy $k$-means clustering. It is neither "better" or "worse". Any differences between the two "paradigms" boil down to parametrizations.

To FINALLY CONCLUDE THIS NOTE, we would like to point readers interested in textbooks dealing with soft $k$-means clustering to the excellent texts by MacKay[9] or Duda, Hart, and Stork[10].

[9] D.J.C. MacKay. *Information Theory, Inference, & Learning Algorithms*. Cambridge University Press, 2003

[10] R.O. Duda, P.E. Hart, and D.G. Stork. *Pattern Classification*. Wiley-Interscience, 2nd edition, 2001

*References*

C. Bauckhage.    Lecture  Notes  on  Data  Science:    $k$-Means Clustering  Is  Gaussian  Mixture  Modeling,  2015a.    DOI: 10.13140/RG.2.1.3033.2646.

C. Bauckhage. Lecture Notes on Data Science: $k$-Means Clustering, 2015b. DOI: 10.13140/RG.2.1.2829.4886.

R.O. Duda, P.E. Hart, and D.G. Stork. *Pattern Classification*. Wiley-Interscience, 2nd edition, 2001.

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