

Lecture Notes on Data Science: k -Means Clustering Is Matrix Factorization

Christian Bauckhage

B-IT, University of Bonn

In this note, we show that k -means clustering can be understood as a constrained matrix factorization problem. This insight will later allow us to recognize that k -means clustering is but a specific latent factor model and closely related to techniques such as non-negative matrix factorization or archetypal analysis.

Introduction

Previously, we discussed¹ that hard k -means clustering of a data set $X = \{x_1, x_2, \dots, x_n\} \subset \mathbb{R}^m$ into k clusters C_1, \dots, C_k boils down to the problem of finding appropriate cluster centroids μ_1, \dots, μ_k and that these will minimize the following objective function

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

where

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

OUR PURPOSE IN THIS NOTE is to show that there is yet another way of how to formalize the k -means objective in (1).

To this end, we note that we may understand the binary indicator variables z_{ij} in (2) as the elements of an *indicator matrix* $Z \in \mathbb{R}^{k \times n}$.

We also observe that we may think of the given data points x_j as the columns of a *data matrix*

$$X = \begin{bmatrix} x_1 & x_2 & \dots & x_n \end{bmatrix} \in \mathbb{R}^{m \times n} \quad (3)$$

and that we may furthermore introduce a *centroid matrix*

$$M = \begin{bmatrix} \mu_1 & \mu_2 & \dots & \mu_k \end{bmatrix} \in \mathbb{R}^{m \times k} \quad (4)$$

whose columns correspond to the cluster centroids that are to be determined.

Given the matrices defined in (2), (3), and (4), we will show that the k -means objective function in (1) can indeed be written as

$$\sum_{i=1}^k \sum_{j=1}^n z_{ij} \|x_j - \mu_i\|^2 = \|X - MZ\|_F^2 \quad (5)$$

where $\|\cdot\|_F$ denotes the matrix *Frobenius norm*.

IN OTHER WORDS, we will show that k -means clustering is a matrix factorization problem! If there were two appropriate matrices M and Z that would minimize the right hand side of (5), the data matrix X could be approximated as $X \approx MZ$.

¹ C. Bauckhage. Lecture Notes on Data Science: k -Means Clustering, 2015b. DOI: 10.13140/RG.2.1.2829.4886; and C. Bauckhage. Lecture Notes on Data Science: k -Means Clustering Is Gaussian Mixture Modeling, 2015a. DOI: 10.13140/RG.2.1.3033.2646



data matrix



centroid matrix



alternative form of the k -means objective function

Exercise: convince yourself that MZ is a $m \times n$ matrix.

Proving Equation (5)

In this section, we will prove that our claim in (5) does indeed hold. The basic idea is to expand both sides of the equation into several, more elementary terms and to show that the expressions we obtain for the left- and right hand side are indeed equivalent.

Yet, before we set out to do so, we will remind ourselves of general properties of the Frobenius norm and point out some of the peculiar features of the binary indicator matrix \mathbf{Z} .

General Properties of the Squared Frobenius Norm of a Matrix

Let $\mathbf{A} \in \mathbb{R}^{m \times n}$ be any real valued matrix of m rows and n columns. To denote individual elements of such a matrix, we either write a_{ij} or $(\mathbf{A})_{ij}$ and to refer to the j -th column vector of \mathbf{A} , we write \mathbf{a}_j .

The **squared Frobenius norm** of \mathbf{A} is defined as

$$\|\mathbf{A}\|_F^2 = \sum_{i=1}^m \sum_{j=1}^n a_{ij}^2 \quad (6)$$

and we recall the following properties

$$\|\mathbf{A}\|_F^2 = \sum_{j=1}^n \|\mathbf{a}_j\|^2 = \sum_{j=1}^n \mathbf{a}_j^T \mathbf{a}_j = \sum_{j=1}^n (\mathbf{A}^T \mathbf{A})_{jj} = \text{tr}[\mathbf{A}^T \mathbf{A}]. \quad (7)$$

Since our derivation below will frequently allude to the identities in (7), readers are encouraged to verify (7) for themselves.

Exercise: convince yourself that all the equalities in (7) do hold.

Peculiar Properties of the Indicator Matrix \mathbf{Z}

If the clusters C_1, \dots, C_k have distinct cluster centroids μ_1, \dots, μ_k , each of the n columns of \mathbf{Z} will contain a single element that is 1 and $k-1$ elements that are 0. Accordingly, each column j of \mathbf{Z} will sum to one

$$\sum_{i=1}^k z_{ij} = 1 \quad (8)$$

and the k different row sums will indicate the number of elements per cluster, that is, for each row i of \mathbf{Z} , we have

$$\sum_{j=1}^n z_{ij} = |C_i| = n_i. \quad (9)$$

Moreover, since $z_{ij} \in \{0, 1\}$ and each column of \mathbf{Z} only contains a single 1, the rows of \mathbf{Z} are pairwise perpendicular because

$$z_{ij} z_{i'j} = \begin{cases} 1, & \text{if } i = i' \\ 0, & \text{otherwise} \end{cases} \quad (10)$$

which is then to say that the matrix $\mathbf{Z}\mathbf{Z}^T$ is a diagonal matrix where

$$(\mathbf{Z}\mathbf{Z}^T)_{ii'} = \sum_j (\mathbf{Z})_{ij} (\mathbf{Z}^T)_{ji'} = \sum_j z_{ij} z_{i'j} = \begin{cases} n_i, & \text{if } i = i' \\ 0, & \text{otherwise.} \end{cases} \quad (11)$$

HAVING FAMILIARIZED OURSELVES with these properties of the indicator matrix, we are now positioned to establish the equalities in (5) which we will do in a step by step manner.

Step 1: Expanding the expression on the left of (5)

We begin by expanding the traditional k -means objective on the left of (5). For this expression, we have

$$\begin{aligned} \sum_{i,j} z_{ij} \|\mathbf{x}_j - \boldsymbol{\mu}_i\|^2 &= \sum_{i,j} z_{ij} (\mathbf{x}_j^T \mathbf{x}_j - 2\mathbf{x}_j^T \boldsymbol{\mu}_i + \boldsymbol{\mu}_i^T \boldsymbol{\mu}_i) \\ &= \underbrace{\sum_{i,j} z_{ij} \mathbf{x}_j^T \mathbf{x}_j}_{T_1} - 2 \underbrace{\sum_{i,j} z_{ij} \mathbf{x}_j^T \boldsymbol{\mu}_i}_{T_2} + \underbrace{\sum_{i,j} z_{ij} \boldsymbol{\mu}_i^T \boldsymbol{\mu}_i}_{T_3}. \end{aligned} \quad (12)$$

This expansion leads to further insights, if we examine the three terms T_1 , T_2 , and T_3 one by one.

FIRST OF ALL, we find

$$T_1 = \sum_{i,j} z_{ij} \mathbf{x}_j^T \mathbf{x}_j = \sum_{i,j} z_{ij} \|\mathbf{x}_j\|^2 \quad (13)$$

$$= \sum_j \|\mathbf{x}_j\|^2 \quad (14)$$

$$= \text{tr}[\mathbf{X}^T \mathbf{X}] \quad (15)$$

where we made use of (8) and (7).

SECOND OF ALL, we observe

$$T_2 = \sum_{i,j} z_{ij} \mathbf{x}_j^T \boldsymbol{\mu}_i = \sum_{i,j} z_{ij} \sum_l x_{lj} \mu_{li} \quad (16)$$

$$= \sum_{j,l} x_{lj} \sum_i \mu_{li} z_{ij} \quad (17)$$

$$= \sum_{j,l} x_{lj} (\mathbf{MZ})_{lj} \quad (18)$$

$$= \sum_j \sum_l (X^T)_{jl} (\mathbf{MZ})_{lj} \quad (19)$$

$$= \sum_j (X^T \mathbf{MZ})_{jj} \quad (20)$$

$$= \text{tr}[\mathbf{X}^T \mathbf{MZ}]. \quad (21)$$

THIRD OF ALL, we note that

$$T_3 = \sum_{i,j} z_{ij} \boldsymbol{\mu}_i^T \boldsymbol{\mu}_i = \sum_{i,j} z_{ij} \|\boldsymbol{\mu}_i\|^2 \quad (22)$$

$$= \sum_i \|\boldsymbol{\mu}_i\|^2 n_i \quad (23)$$

where we applied (9).

Step 2: Expanding the expression on the right of (5)

Next, we look at the expression on the right hand side of (5). As a squared Frobenius norm of a matrix difference, it can be written as

$$\begin{aligned} \|X - MZ\|^2 &= \text{tr}[(X - MZ)^T(X - MZ)] \\ &= \underbrace{\text{tr}[X^T X]}_{T_4} - 2 \underbrace{\text{tr}[X^T MZ]}_{T_5} + \underbrace{\text{tr}[Z^T M^T MZ]}_{T_6} \end{aligned} \quad (24)$$

GIVEN OUR RESULTS in (15) and (21), we immediately recognize that $T_1 = T_4$ and $T_2 = T_5$. Thus, to establish that (12) and (24) are indeed equivalent, it remains to verify whether $T_3 = T_6$?

REGARDING TERM T_6 , we note that, due to the cyclic permutation invariance of the [trace operator](#), we have

$$\text{tr}[Z^T M^T MZ] = \text{tr}[M^T MZZ^T]. \quad (25)$$

We also note that

$$\text{tr}[M^T MZZ^T] = \sum_i (M^T MZZ^T)_{ii} \quad (26)$$

$$= \sum_i \sum_l (M^T M)_{il} (ZZ^T)_{li} \quad (27)$$

$$= \sum_i (M^T M)_{ii} (ZZ^T)_{ii} \quad (28)$$

$$= \sum_i \|\mu_i\|^2 n_i \quad (29)$$

where we used the fact that ZZ^T is diagonal. This result, however, shows that $T_3 = T_6$ and, consequently, that (12) and (24) really are equivalent.

Summary and Outlook

Using rather tedious yet straightforward algebra, we have shown that the problem of hard k -means clustering can be understood as the following constrained matrix factorization problem

$$\begin{aligned} \underset{M, Z}{\text{argmin}} \quad & \|X - MZ\|^2 \\ \text{s.t.} \quad & z_{ij} \in \{0, 1\} \\ & \sum_i z_{ij} = 1 \end{aligned} \quad (30)$$

where

$$X \in \mathbb{R}^{m \times n} \text{ is a matrix of data vectors } x_j \in \mathbb{R}^m \quad (31)$$

$$M \in \mathbb{R}^{m \times k} \text{ is a matrix of cluster centroids } \mu_i \in \mathbb{R}^m \quad (32)$$

$Z \in \mathbb{R}^{k \times n}$ is a matrix of binary indicator variables such that

$$z_{ij} = \begin{cases} 1, & \text{if } x_j \in C_i \\ 0, & \text{otherwise.} \end{cases} \quad (33)$$

AT THIS POINT, readers who are not accustomed with the idea of matrix factorization for data analysis might be wondering what we could possibly gain from this insight.

Admittedly, the formulation of the k -means clustering problem in (30) appears to be more complicated and less intuitive than those found in the textbooks. However, in later notes, we will see that the expression in (30) allows for seamless insights into several important properties of the k -means clustering problem that are otherwise more difficult to uncover ².

²C. Bauckhage. k -Means Clustering via the Frank-Wolfe Algorithm. In *Proc. KDML-LWDA*, 2016

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.
- C. Bauckhage. k -Means Clustering via the Frank-Wolfe Algorithm. In *Proc. KDML-LWDA*, 2016.