

Unsupervised Learning: Clustering and Matrix Factorization

EE219: Large Scale Data Mining

Professor Roychowdhury

K-means clustering algorithm

Algorithm

0. Randomly initialize K cluster centers (centroids)
1. Iterate until convergence
 - 1.1 For each data point, find closest cluster center (partitioning step)
 - 1.2 Replace each centroid by average of data points in its partition

Objective function

Write $x_i = (x_{i1}, \dots, x_{ip})$:

Let the centroids be denoted by m_1, m_2, \dots, m_k and the clusters by c_1, c_2, \dots, c_k , then the objective function of K-means is to minimize Euclidean distance of the points with the centroids of corresponding clusters (within cluster sum of squares):

$$\sum_{k=1}^K \sum_{i \in c_k} \|x_i - m_k\|^2$$

$$m_k = \frac{1}{N_k} \sum_{i \in c_k} x_i$$

- ▶ Consider the assignment function $C(i)$:

$$C : 1, 2, \dots, N \rightarrow (1, 2, \dots, K)$$

- ▶ K -means minimizes $W(C)$

$$\begin{aligned} W(C) &= \frac{1}{2} \sum_{k=1}^K \sum_{C(i)=k} \sum_{C(j)=k} \|x_i - x_j\|^2 \\ &= \sum_{k=1}^K N_k \sum_{C(i)=k} \|x_i - \bar{x}_k\|^2 \end{aligned}$$

A proof for this equivalence is given in the following slide

- ▶ K -means solves the following problem to find assignment function C :

$$\min_{C, m_1, \dots, m_K} \sum_k \sum_{C(i)=k} \|x_i - m_k\|^2$$

The outer summation ($k = 1$ through K) is over different clusters. The summands for each k are data within the cluster k . So we just prove the equivalence for each cluster k . In other words, we show that if the number of data points in cluster k is N_k , then:

$$\sum_{C(i)=k} \sum_{C(j)=k} \|x_i - x_j\|^2 = \sum_{i=1}^{N_k} \sum_{j=1}^{N_k} \|x_i - x_j\| = N_k \sum_{i=1}^{N_k} \|x_i - \bar{x}\|^2$$

$$\begin{aligned} & \sum_{i=1}^{N_k} \left(\sum_{j=1}^{N_k} \|x_i - x_j\|^2 \right) \\ &= \sum_{i=1}^{N_k} \left(\sum_{j=1}^{N_k} \|(x_i - \bar{x}) - (x_j - \bar{x})\|^2 \right) \\ &= \sum_{i=1}^{N_k} \left(\sum_{j=1}^{N_k} [(x_i - \bar{x}) - (x_j - \bar{x})]^T [(x_i - \bar{x}) - (x_j - \bar{x})] \right) \\ &= \sum_{i=1}^{N_k} \left(\sum_{j=1}^{N_k} \left(\|x_i - \bar{x}\|^2 + \|x_j - \bar{x}\|^2 - 2(x_i - \bar{x})^T (x_j - \bar{x}) \right) \right) \end{aligned}$$

$$\begin{aligned}
&= \sum_{i=1}^{N_k} \left(\sum_{j=1}^{N_k} \left(\|x_i - \bar{x}\|^2 + \|x_j - \bar{x}\|^2 - 2(x_i - \bar{x})^T (x_j - \bar{x}) \right) \right) \\
&= \sum_{i=1}^{N_k} \left(N_k \|x_i - \bar{x}\|^2 + \sum_{j=1}^n \|x_j - \bar{x}\|^2 - 2 \sum_{j=1}^n (x_i - \bar{x})^T (x_j - \bar{x}) \right) \\
&= \sum_{i=1}^{N_k} \left(2N_k \|x_i - \bar{x}\|^2 - 2(x_i - \bar{x})^T \sum_{j=1}^n (x_j - \bar{x}) \right) \\
&= 2N_k \sum_{i=1}^{N_k} \left(\|x_i - \bar{x}\|^2 \right)
\end{aligned}$$

Thus:

$$\begin{aligned}
\sum_{i=1}^{N_k} \left(\sum_{j=1}^{N_k} \|x_i - x_j\|^2 \right) &= 2N_k \sum_{i=1}^{N_k} \|x_i - \bar{x}\|^2 \Rightarrow \\
\sum_{\substack{i,j=1 \\ i < j}}^{N_k} \|x_i - x_j\|^2 &= N_k \sum_{i=1}^{N_k} \|x_i - \bar{x}\|^2
\end{aligned}$$

Initial centroid problem

$$D \approx [u_1 u_2 \dots u_k] \Sigma \begin{bmatrix} \square & \square & \square & \square & \square & \square \end{bmatrix}$$

($t \times d$) \uparrow Eigen documents $\uparrow \uparrow \uparrow \uparrow$ d

- K-means converges to a local optimum whose quality largely depends on the initial choice of centroids
- Solution: multiple (e.g. 100) runs with random initial cluster centroids, then choosing the ones with the minimal final cost function

SVD:

$$D \approx [u_1 u_2 \dots u_k] \Sigma \begin{bmatrix} v_1^T \\ v_2^T \\ \vdots \\ v_k^T \end{bmatrix}$$

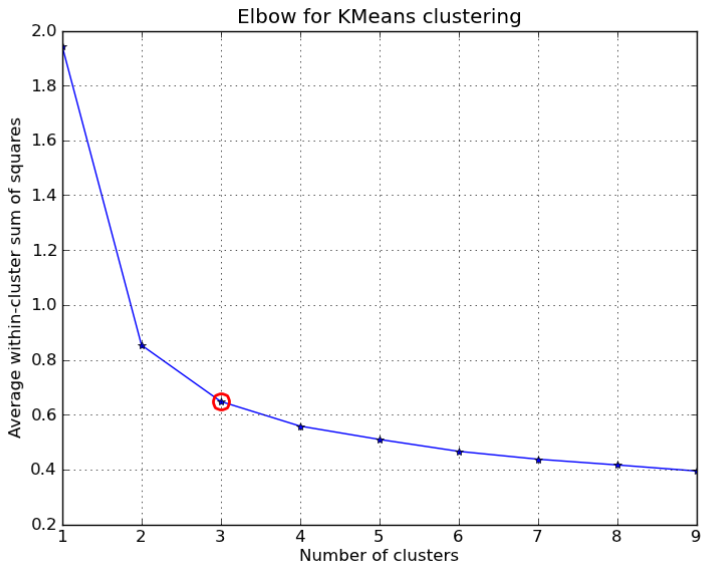
($t \times d$) $\tilde{d}_i = \mathcal{I}$

$d_i \rightarrow \begin{bmatrix} u_1^T \\ u_2^T \\ \vdots \\ u_k^T \end{bmatrix} d_i$

$\begin{bmatrix} \sigma_1 v_1^T \\ \sigma_2 v_2^T \\ \vdots \end{bmatrix}$

How to choose K?

- Elbow method



$$X = \begin{bmatrix} \uparrow & \uparrow & & \uparrow \\ x_1 & x_2 & \cdots & x_n \\ \downarrow & \downarrow & & \downarrow \end{bmatrix}$$

$$X X^T = \sum_{i=1}^N x_i x_i^T$$

(n x n)

$$(X X^T) z_i = \lambda_i z_i$$

$$\begin{bmatrix} z_1 & \cdots & z_n \end{bmatrix} \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix} \begin{bmatrix} z_1^T \\ \vdots \\ z_n^T \end{bmatrix}$$

($\lambda_1, \lambda_2, \lambda_3, \dots, \lambda_n$)

$$\begin{aligned}
 (X X^T) &\approx \underbrace{\left[\begin{array}{c} [z_1 z_2 \dots z_k] \\ \vdots \end{array} \right]}_K \underbrace{\left[\begin{array}{c|c} \begin{bmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_k \end{bmatrix} & \begin{bmatrix} z_1^T \\ \vdots \\ z_k^T \end{bmatrix} \end{array} \right]}_{K \times n} \\
 &\approx \sum_{i=1}^k \lambda_i z_i z_i^T
 \end{aligned}$$

$$\begin{array}{ccc}
 X_i \rightarrow & & \begin{bmatrix} z_1^T x_i \\ \vdots \\ z_k^T x_i \end{bmatrix} \\
 m \times 1 & K \times 1 &
 \end{array}$$

$$\begin{bmatrix} z_1^T \\ \vdots \\ z_k^T \end{bmatrix} X X^T \approx \begin{bmatrix} z_1^T \\ z_2^T \\ \vdots \\ z_k^T \end{bmatrix} \begin{bmatrix} z_1 & z_2 & \dots & z_k \\ \vdots & & & \end{bmatrix} \dots$$

$$\begin{bmatrix} z_1^T \\ \vdots \\ z_k^T \end{bmatrix} X X^T = \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_k \end{bmatrix} \begin{bmatrix} z_1^T \\ \vdots \\ z_k^T \end{bmatrix}$$

Projected
data vectors

$$\underbrace{\begin{bmatrix} z_1^T \\ \vdots \\ z_k^T \end{bmatrix}}_{k \times N} X \underbrace{X^T \begin{bmatrix} z_1 & z_2 & \dots & z_k \end{bmatrix}}_{N \times k} = \begin{bmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_k \end{bmatrix}$$

Non-negative matrix factorization (NMF)

Motivation

- ▶ PCA involves some basis adding some images and then subtracting others Eigen images (e.g. faces) are not intuitive

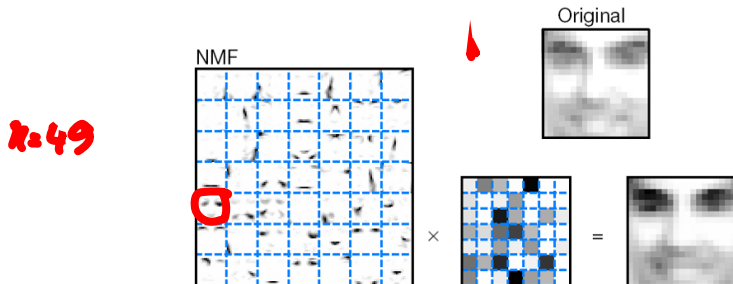
$K=15$
 $Z_1 = \begin{bmatrix} \vdots \\ 106 \end{bmatrix}$



$m=10$

- ▶ Similar arguments in the context of document classification: subtracting doesn't make intuitive sense

- ▶ NMF is like PCA, except the coefficients in the linear combination cannot be negative.
- ▶ In other words, we only allowing adding of basis images to make intuitive sense
- ▶ In the context of images, non-negative basis images represent *"parts"*



- ▶ We don't produce optimal (in norm-2 sense) basis images like PCA, however enforcing the non-negativity constraint gives better properties for classification tasks

Objective function

$$A = U \Sigma V^T$$
$$B^* = U_r \Sigma_r V_r^T$$

B is of rank $\leq r$

- Norm distance

$$\|A - B\|_F^2$$

where $\|\cdot\|_F$ denotes the Frobenius norm:

$$B = \sum_{i=1}^r \lambda_i v_i v_i^T$$

$$\|A\|_F = \left(\sum_{i,j=1}^n |a_{ij}|^2 \right)^{1/2}$$

- KL-divergence

$$D(A||B) = \sum_{ij} \left(A_{ij} \log \frac{A_{ij}}{B_{ij}} - A_{ij} + B_{ij} \right)$$

Objective function

$$A_{m \times n} \quad B = WH_{r \times n} = \sum_{i=1}^r w_i h_i^T$$

► Matrix V is the original $m \times n$ matrix, and matrices $W_{m \times r}$ and $H_{r \times n}$ give a lower rank approximation for V

► We can choose either of the following problems to solve

1. minimize $\|V - WH\|_F^2$ s.t. $W, H \geq 0$

2. minimize $D(V \| WH)$ s.t. $W, H \geq 0$

► Using gradient descent to find a local minimum, the gradient descent update rule for the norm distance cost function becomes [Lee, Seung '01]

$$H_{a\mu} \leftarrow H_{a\mu} + \eta_{a\mu} \left[(W^T V)_{a\mu} - (W^T WH)_{a\mu} \right]$$

$$W = [w_1 \ w_2 \ \dots \ w_r] \quad H = \begin{bmatrix} h_1^T \\ \vdots \\ h_r^T \end{bmatrix} \quad \begin{matrix} w_i \geq 0 \\ h_i \geq 0 \end{matrix}$$

$$V \approx WH$$

$$V_{ij} \approx \sum_k W_{ik} H_{kj}$$

V
cost

A nice property of the update rule

- ▶ The update rule is

$$\boxed{H_{a\mu}} \leftarrow H_{a\mu} + \eta_{a\mu} \left[(W^T V)_{a\mu} - (W^T W H)_{a\mu} \right]$$

$$H \leftarrow H + \eta (W^T V - W^T W H)$$

- ▶ We set $\eta_{a\mu} = \frac{H_{a\mu}}{(W^T W H)_{a\mu}}$
- ▶ Then the update rule becomes

$$\boxed{H_{a\mu} \leftarrow H_{a\mu} \frac{(W^T V)_{a\mu}}{(W^T W H)_{a\mu}}}$$

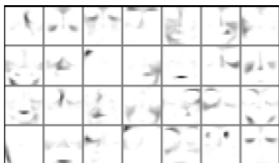
- ▶ Similarly for W we get

$$W_{ia} \leftarrow W_{ia} \frac{(V H^T)_{ia}}{(W H H^T)_{ia}}$$

- ▶ This is a multiplicative update instead of an additive update. If the initial values of W and H are both non-negative, then W and H will always remain non-negative.
- ▶ This guarantees a non-negative factorization.

Sparsity

- ▶ NMF doesn't always give parts based result [Hoyer '04]
- ▶ In figure below, on the left, a dataset of $m = 2429$ facial images, each consisting of $n = 19 \times 19$ pixels, is factorized over $r = 49$ basis images; and the basis images are displayed. On the right, another dataset of $m = 400$ images having $n = 112 \times 92$ pixels is factorized over $r = 25$ basis images.
- ▶ The factors on the right are different in that they are qualitatively global rather than local.
- ▶ The difference in results is attributed to the fact that in the dataset used for the experiment on the left, facial images are well aligned, whereas in the experiment on the right the faces have significant variations in face angles, illumination angle, facial expressions, etc.



Sparsity

- ▶ We put a constraint on a sparsity measure defined as follows:

$$sparseness(x) = \frac{\sqrt{n} - (\sum |x_i|) / \sqrt{\sum x_i^2}}{\sqrt{n} - 1}$$

- ▶ This function evaluates to unity if and only if x contains only a single non-zero component, and takes a value of zero if and only if all components are equal (up to signs), interpolating smoothly between the two extremes.

For $X \in \mathbb{R}^{n \times n}$, by definition, $\text{Tr}(X) = \sum_{i=1}^n X_{ii}$. Let us define the inner product of two matrices $A, B \in \mathbb{R}^{n \times n}$ as:

$$\langle A, B \rangle = \text{Tr}(AB^T)$$

$$\begin{aligned} \|A\|_F^2 &= \text{Tr}(AA^T) \\ &= \text{Tr}(A^T A) \end{aligned}$$

i.e. it has the following properties:

1. Symmetry:

$$\langle A, B \rangle = \langle B, A \rangle$$

2. Linearity in the first argument:

$$\begin{aligned} \langle aA, B \rangle &= a \langle A, B \rangle \\ \langle A, B + C \rangle &= \langle A, B \rangle + \langle A, C \rangle \end{aligned}$$

3. Positive-definiteness:

$$\begin{aligned} \langle A, A \rangle &\geq 0 \\ \langle A, A \rangle &= 0 \Leftrightarrow A = 0 \end{aligned}$$

Given the aforementioned definitions we have:

$$\begin{aligned}
 \|V - WH\|_F^2 &= \langle V - WH, V - WH \rangle \\
 &= \langle V, V \rangle + \langle WH, WH \rangle - 2 \langle V, WH \rangle \\
 &= \text{Tr}(V^T V) + \text{Tr}(H^T W^T WH) - 2 \text{Tr}(V^T WH)
 \end{aligned}$$

Taking the gradient of the objective function we get the following:

$$\begin{aligned}
 \nabla_H \|V - WH\|_F^2 &= \nabla_H \left[\text{Tr}(H^T W^T WH) - 2 \text{Tr}(V^T WH) \right] \\
 &= 2(W^T WH - V^T W)
 \end{aligned}$$

$H_{n \times m}$

Where we used:

$$\nabla_X \text{Tr}(AX) = A$$

$$\nabla_X \text{Tr}(X^T AX) = (A + A^T)X$$

$$\frac{\partial \mathcal{L}}{\partial H} (i,j) = \frac{\partial \mathcal{L}}{\partial H_{ij}}$$

$$\boxed{\nabla_H \mathcal{L}}_{n \times m}$$

- The proof of $Tr(.)$ properties are straightforward from definition. Below we prove the second property as an example:

$$\begin{aligned}
 y &= Tr(X^T A X) \\
 &= \sum_i \left(\sum_j \sum_k (X^T)_{ij} A_{jk} X_{ki} \right) \\
 &= \sum_{i,j,k} X_{ji} X_{ki} A_{jk}
 \end{aligned}$$

$$\frac{\partial y}{\partial x} = Ax + A^T x$$

$$\begin{aligned}
 \frac{\partial y}{\partial X_{mn}} &= \sum_{i,j,k} (X_{ki} A_{jk} |_{j=m, i=n} + X_{ji} A_{jk} |_{k=m, i=n}) \\
 &= \sum_{i,j,k} (X_{kn} A_{mk}) + \sum_{i,j,k} (X_{jn} A_{jm}) \\
 &= (AX)_{mn} + (A^T X)_{mn}
 \end{aligned}$$

$$\begin{aligned}
 H_{ij}(t+1) &= H_{ij}(t) \\
 &\quad - \eta \frac{\partial \text{Cost}}{\partial H_{ij}} \\
 H_{ij} & \\
 \eta > 0
 \end{aligned}$$

