

NMF: Non-negative Matrix Factorization



Application \Rightarrow Dimensionality Reduction

1. Dimension Reduction

Formulation: n data $\{v_1, \dots, v_n\}$ $v_i \in \mathbb{R}^m$



find r basis vector $\{w_1, w_2, \dots, w_r\}$

Then
$$v_j \approx \sum_{i=1}^r w_i \underline{H_{ij}} \rightarrow \{ \underline{H_{ij}}, \dots, \underline{H_{rj}} \}$$

$$= [w_1, \dots, w_r] \begin{bmatrix} H_{1j} \\ \vdots \\ H_{rj} \end{bmatrix} \xrightarrow{\text{representation vector of } v_j}$$

matrix form $V = [v_1, \dots, v_n]$

coordinates

$$= [w_1, \dots, w_r] \cdot \begin{bmatrix} H_{11} & \dots & H_{1n} \\ \vdots & & \vdots \\ H_{r1} & \dots & H_{rn} \end{bmatrix}$$

$$:= W \cdot H$$

Here $\begin{cases} V \in \mathbb{R}^{m \times n} \\ W \in \mathbb{R}^{m \times r} \quad H \in \mathbb{R}^{r \times n} \end{cases}$

\rightarrow Error Measure $\begin{cases} 1. \|V - WH\|_F^2 \\ 2. \|V - WH\|_1 \end{cases}$ etc...

\rightarrow Constraints $\begin{cases} 1. \text{Non-negative } W \geq 0 \quad H \geq 0 \\ 2. \text{orthogonal } HH^T = I_r \\ 3. \text{symmetry } \underline{H = W^T} \quad (\underline{m=n!}) \end{cases}$ NMF

2. NMF Task

→ given $V \in \mathbb{R}_+^{m \times n}$ and a rank r (of V)

→ output $W \in \mathbb{R}_+^{m \times r}$ $H \in \mathbb{R}_+^{r \times n}$

This is achieved by:
$$\begin{cases} \min_{W, H} & \frac{1}{2} \|V - WH\|_F^2 \\ \text{s.t.} & W \geq 0, H \geq 0 \end{cases}$$

$W \geq 0$ is labeled basis and $H \geq 0$ is labeled coordinates.

Rmk: 1. object $f(W, H)$ is non-convex
but bi-convex (fix one)

interpretation

2. each column of V : represents one image!
(application)

flatten

3. Data



in each image dataset, images are well-aligned



noses are roughly in the same location

4. Application

① { Face Data V
manually set r → how many items can re-construct one face!

⇒ NMF ⇒ $W \in \mathbb{R}^{m \times r}$ $V \in \mathbb{R}^{r \times n}$ → Low-rank approximation

Interpretation :

(BASIS)

1. each column of $W \Rightarrow$ one type of feature in one image
(e.g. mouse, nose e.t.c.)

2. each column of $H \Rightarrow$ the ingredient of each feature
in matrix W (intensity)

② Text mining

a) term-document matrix

several sentences \rightarrow

word

Document Index

V

$$V_{ij} := (\# \text{ of } \underline{\text{word } i} \text{ appear in } \underline{\text{doc } j})$$

$$b) V \approx W \cdot H$$

interpretation :

{ column of W : topic

{ column of H : intensity of each topic in one doc

5. Exact NMF

→ given $V \in \mathbb{R}_+^{n \times n}$, if there exists $\begin{cases} W \in \mathbb{R}_+^{n \times r} \\ H \in \mathbb{R}_+^{r \times n} \end{cases}$

s.t. $V = W \cdot H$ (perfect recovery)

then $[WH \text{ is an exact NMF of } V]$

- Rmk:
1. exact NMF may not exist for some r
 2. exact NMF may not be unique
 3. [non-negative rank] $\text{rank}_+(V)$

smallest r s.t. $V = W_r \cdot H_r$
(SVD)

4. $\text{rank}(V) \leq \text{rank}_+(V) \leq \min(m, n)$



(trivial)

non-negative constraint

5. $V \geq 0$, $\text{rank}(V) \leq 2$.

then $\text{rank}(V) = \text{rank}_+(V)$

6. [counter-example] of 5.)

when $\text{rank}(V) = 3$, $\text{rank}_+(V)$ can be larger than 3

$V = \begin{pmatrix} 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 1 \end{pmatrix} \Rightarrow \text{rank}_+(V) = 4$

6. Variants of NMF \Rightarrow adding prior knowledge

- ↓
- ① Orthogonal NMF \rightarrow for H ? why
 - ② Symmetric NMF
 - ③ Sparse NMF ($\| \cdot \|_1$ penalty)

① Orthogonal NMF

Lemma: $\underline{H \geq 0} + \underline{HH^T = I_r}$



each column of H has at most 1 positive member

PF: $H = \begin{bmatrix} H_{1r} \\ \vdots \\ H_{nr} \end{bmatrix} \Rightarrow (HH^T)_{ij} = H_{i \cdot} H_{j \cdot}^T = 0 \quad \underline{i \neq j}$

\Rightarrow each row has disjoint support

\Rightarrow each column of H at most has 1 positive entry

Rmk: Lemma says that, Orthogonal NMF \Rightarrow clustering ^{FA}

Then, column of W \Rightarrow cluster centroid

column of H \Rightarrow cluster membership

membership matrix $H \in \mathbb{R}^{r \times n}$

To solve orthogonal NMF, we add penalty term

Orthogonal NMF Formulation:

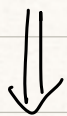
$$\begin{cases} \min_{W, H} & \frac{1}{2} \|V - WH\|_F^2 + \frac{\lambda}{2} \|HH^T - I_r\|_F^2 \\ \text{s.t.} & W \geq 0, H \geq 0 \end{cases}$$

② Sparse NMF \rightarrow Difficult to optimize

$$\begin{cases} \min_{W, H} & \frac{1}{2} \|V - WH\|_F^2 + \lambda_W \|W\|_1 + \lambda_H \|H\|_1 \\ \text{s.t.} & W \geq 0, H \geq 0 \end{cases}$$

③ Adjacency Matrix V (undirected graph G)

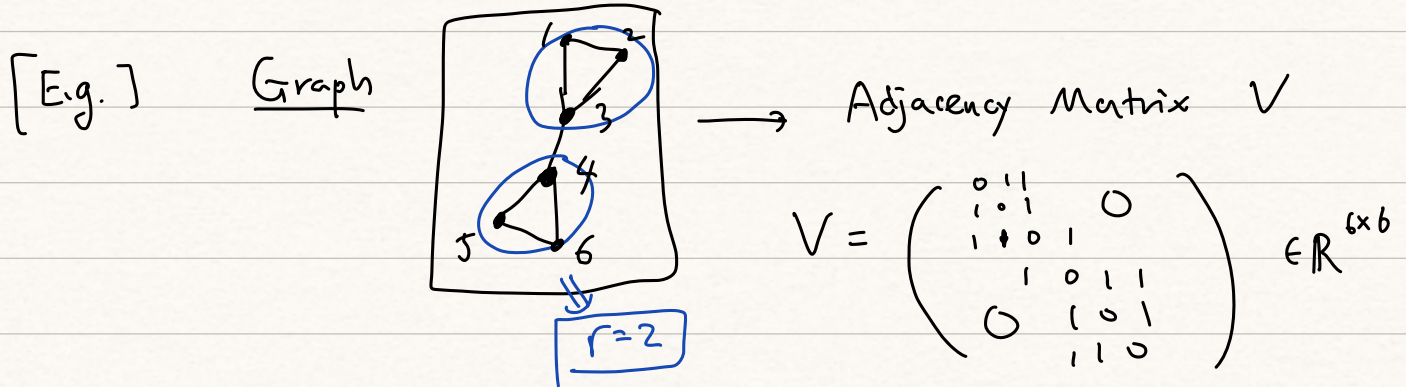
$$V_{ij} = \begin{cases} 0, & i \not\leftrightarrow j \\ 1, & i \leftrightarrow j \end{cases} \rightarrow \text{symmetry}$$



Symmetric NMF \rightarrow when V is symmetric

$$\begin{cases} \min_W & \frac{1}{2} \|V - WW^T\|_F^2 \\ \text{s.t.} & W \geq 0 \end{cases}$$

Application \rightarrow { Adjacency Matrix of Undirected Graph
Community detection
 \rightarrow dense connected nodes
 $r = (\# \text{ of communities})$



obviously, we have 2 communities

$W \in \mathbb{R}^{6 \times 2}$

node $\begin{matrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \\ 6 \end{matrix}$ community $\begin{matrix} 1 & 2 \\ a & 0 \\ b & 0 \\ c & d \\ d & c \\ 0 & b \\ 0 & a \end{matrix}$

$d \approx 0.24$
 $a \approx b \approx c \approx 0.8$

7. Algorithms for NMF

- ① PG & APG
 - ② MM
 - ③ ALS
 - ④ HALS
- \rightarrow block coordinate descent

$f(W, H)$

① PG

$$\min_{W, H} \frac{1}{2} \|V - WH\|_F^2 + \delta_{\mathbb{R}_+^{m \times r}}(W) + \delta_{\mathbb{R}_+^{r \times n}}(H)$$

for simplicity, define $\Pi_+(\cdot) := \Pi_{\mathbb{R}_+^{m \times r}}(\cdot) := \Pi_{\mathbb{R}_+^{r \times n}}(\cdot)$

$$\begin{cases} W^{(k+1)} = \Pi_+ (W^{(k)} - \alpha_k \nabla_W f(W^{(k)}, H^{(k)})) \\ H^{(k+1)} = \Pi_+ (H^{(k)} - \alpha_k \nabla_H f(W^{(k)}, H^{(k)})) \end{cases}$$

② MU Algorithm (Variant of PG)

↓
w.r.t step length

$$\begin{cases} W^{(k+1)} = \Pi_+ (W^{(k)} - S_{W^{(k)}} \odot \nabla_W f(W^{(k)}, H^{(k)})) \\ H^{(k+1)} = \Pi_+ (H^{(k)} - S_{H^{(k)}} \odot \nabla_H f(W^{(k)}, H^{(k)})) \end{cases}$$

element-wise multiplication

Here :

$$\begin{cases} S_{W^{(k)}} := W^{(k)} \odot [W^{(k)} H^{(k)} H^{(k)T}] \\ S_{H^{(k)}} := H^{(k)} \odot [W^{(k)T} W^{(k)} H^{(k)}] \end{cases}$$

the reason to choose this step length is:

can simplify update equation

⇓

Reason : $W - S_W \odot \nabla_W f$

$$= W + S_W \odot R H^T$$

$$= W + \frac{W}{W H H^T} \odot R H^T$$

$$= W + \frac{W}{WHH^T} \odot (VH^T - WHH^T)$$

$$= \underbrace{W \odot VH^T \odot WHH^T}_{\geq 0}$$

③ ALS (Alternating Least Square) (交替下降法)
alternative minimization (EM-like)
efficient solver

$$\left\{ \begin{array}{l} \underline{H \leftarrow \underset{H}{\operatorname{argmin}} \|V - WH\|_F^2} \quad (\text{given } W) \\ H \leftarrow \Pi_+(H) \\ W \leftarrow \underset{W}{\operatorname{argmin}} \|V - WH\|_F^2 \quad (\text{given } H) \\ W \leftarrow \Pi_+(W) \end{array} \right.$$

8. Choice of r

