

[9] Nonnegative Matrix Factorization (NMF)

Given $X \in \mathbb{R}_+^{d \times n}$ and $0 < r \leq \min\{d, n\}$, the NMF problem is to find $W \in \mathbb{R}_+^{d \times r}$ and $H \in \mathbb{R}_+^{r \times n}$ such that

$$X = WH$$

Main keyword: Interpretability

or that minimize $\min_{W \in \mathbb{R}_+^{d \times r}, H \in \mathbb{R}_+^{r \times n}} \|X - WH\|_F^2$

Preprocessing: remove zero columns from X , $0 = w \cdot 0$. Now everything can be normalized with respect to $\|\cdot\|_1$ vector norm:

$$X = WH \Leftrightarrow XD = W \cdot (HD) \Leftrightarrow XD = (\tilde{X}\tilde{D})(\tilde{D}^{-1}HD)$$

For any D, \tilde{D} diagonal invertible matrices. We can take

$$D = \text{diag}(\|x_{:,i}\|_1), \tilde{D} = \text{diag}(\|w_{:,i}\|_1) \text{ to obtain } \tilde{X} = \tilde{W}\tilde{H}, e^T \tilde{X} = e^T, e^T \tilde{W} = e^T, e^T \tilde{H} = e^T \tilde{W}\tilde{H} = e^T$$

→ We obtain an equivalent NMF $X = WH$ where $X, W, H \geq 0$ and H is column-stochastic

$$\Rightarrow x_{:,i} = h_{1,i} w_{:,1} + h_{2,i} w_{:,2} + \dots + h_{r,i} w_{:,r}$$

⇒ each data $x_{:,i}$ is now a convex comb. of the columns of W with weights given by $H_{:,i}$

We can thus interpret each $x_{:,i}$ as a superposition of "features" $w_{:,j}$ where $h_{j,i}$ says what % of $w_{:,j}$ is inside $x_{:,i}$

→ in case of clustering, we can take the biggest $h_{j,i}$ to choose which $w_{:,j}$ represents $x_{:,i}$ better

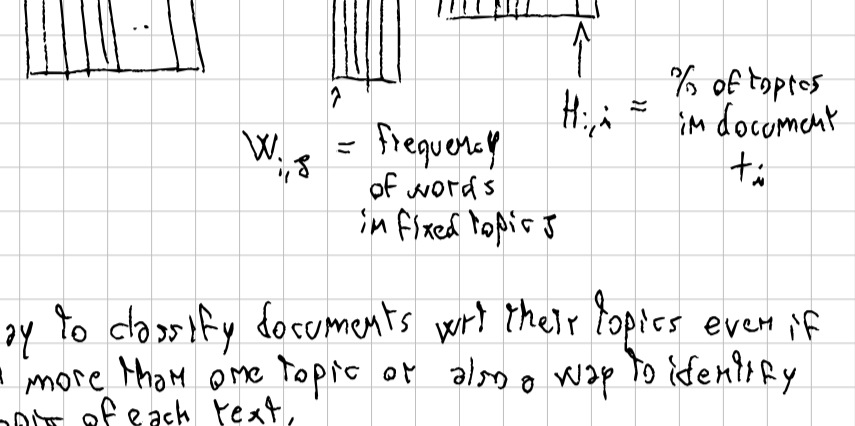
Example with Text Mining: Suppose we have n documents written over with a dictionary of d words. $x_{j,i}$ will be the frequency of word d_j in document t_i

$$x_{j,i} = \frac{\#\{\text{occurrence of word } d_j \text{ in } t_i\}}{\#\{\text{words in } t_i\}}$$

A document may cover different topics. The idea is that each topic is identified by a set of words that may intersect, and has some unique ones.

sports → score, play, goal; religion → church, priest, mass; politics → minister, government, leadership

document covering more topics will have a mix of those defining words inside



This is a way to classify documents wrt their topics even if they contain more than one topic or also a way to identify the main topics of each text.

In the case of Exact NMF $X = WH$, we call

$$\text{Nonnegative Rank} : \text{rank}_+(X) = \min \{r \mid W \in \mathbb{R}_+^{d \times r}, H \in \mathbb{R}_+^{r \times n}\}$$

where obviously $\text{rank}_+(X) \geq \text{rank}(X)$. There is a big literature about the NMF, but here we give only some defining properties

- If $\text{rank}(X) \leq 2$, then $\text{rank}_+(X) = \text{rank}(X)$
For $\text{rank}(X) = 3$, there are examples for $\text{rank}_+(X) = 2$ (Vd)
- $\text{rank}_+(X) \leq \min\{d, n\}$, since $X = I_d X = X I_n$
- Up to a translation of data $X \mapsto X + ve^T, v \geq 0$ we have $\text{rank}_+(X + ve^T) = \text{rank}_+(X) + 1$
- It is NP-Hard to check if $\text{rank}_+(X) = \text{rank}(X)$
- If W, H are randomly generated with continuous distr., then with prob. 1 $\text{rank}_+(X) = \text{rank}(X) = r$

→ In applications, it is usually supposed that $X \sim WH$ where $r = \text{rank}(W) = \text{rank}(H)$, and r is a numerical rank of X .

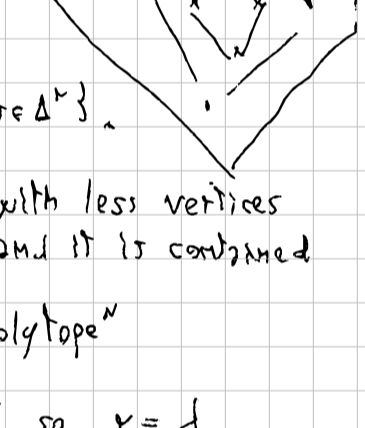
Geometric Interpretation $X = WH, \geq 0, \ell_1$ -norm.

Since $x_{:,i} = W h_{:,i}$ and $h_{:,i}$ are convex combinations, then

$$\text{Conv}(X) \subseteq \text{Conv}(W) \subseteq \Delta^d$$

where

$$\Delta^d = \{x \in \mathbb{R}_+^d \mid e^T x = 1\}, \text{Conv}(A) = \{Ax \mid x \in \Delta^d\}$$



We thus want the convex polytope $\text{Conv}(W)$ with less vertices possible ($\text{rank}_+(X)$) that contains $\text{Conv}(X)$ and it is contained in Δ^d .

→ This is called "Nested Polytope"

Notice: $\text{Conv}(W) = \Delta^d$ coincides with $W = I$ so $r = d$

$\text{Conv}(W) = \text{Conv}(X)$ " " $W = X$ so $r = n$

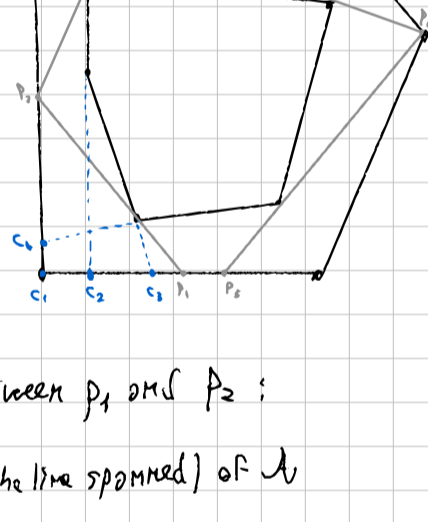
→ we are looking for $r \leq \min\{d, n\}$

NPP: Nested Polytope Problem

Given $A \subseteq B \subseteq \mathbb{R}^m$ polytopes full-dimensional with B bounded, and given $r > 0$ find a polytope E with r vertices and $A \subseteq E \subseteq B$

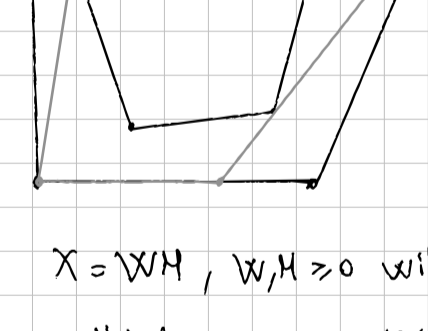
Theorem 4.1 IF $m=2$, NPP is in P [1e]

- Start from a point $p_1 \in \partial B$
- Compute $p_2 = F(p_1)$, i.e. From p_1 draw the tangent to A that intersects B in p_2 and $p_3 = F(p_2)$.
- If $\overline{p_1 p_3}$ does not intersect A , close the polytope, otherwise compute $p_{k+1} = F(p_k)$ until $p_k p_{k+1}$ does that.
- Identify the "contact charge points" between p_1 and p_2 :
 - The vertices of B between p_1, p_2
 - intersection between ∂B and an edge (the line spanned) of A call them c_i
- The optimal E is one of the polytopes generated with F starting from one of the c_i



Theorem 4.2 NPP with r vertices

and NMF with $r, X = WH$ and $\text{rank}(X) = \text{rank}(W) = r$ are equivalent: there exist polynomial reductions between the two algorithms. In particular, they have the same complexity.



Theorem 4.3

IF $X \in \mathbb{R}_+^{d \times n}$ admits $X = WH, W, H \geq 0$ with $\text{rank}(W) = \text{rank}(X) \leq 3$, then it is possible to recover W, H in poly time

Proof ($r=3$) ℓ_1 -normalize X and remove zero columns. Let U be a subset of r indep. columns of X so that $X = UV, U \geq 0$ and ℓ_1 -norm. This means $e^T = e^T X = e^T UV \Rightarrow V$ is ℓ_1 -norm. but may be negative. Call \bar{V} the matrix V without the last row and $A = \text{conv}(\bar{V})$. Since $r=3, A \subseteq \mathbb{R}^2$. Moreover if \bar{U} is U without last column,

$$B := \{x \in \mathbb{R}^2 \mid \bar{U}x + M_3(1-x_1-x_2) \geq 0\}$$

Notice that $A \subseteq B$ because if $v = \begin{bmatrix} v_1 \\ v_2 \\ v_3 \end{bmatrix}$ is a column of V , then $x_i = Uv \geq 0, Uv = \bar{U} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} + M_3(1-v_1-v_2)$.

Since V has $\text{rank} 3, A$ has dimension 2, and so has B . Notice also

$$\bar{U}x + M_3(1-x_1-x_2) = (M_1 - M_3)x_1 + (M_2 - M_3)x_2 + M_3 = Fx + M_3$$

where F is full rank since U was, and $e^T F = 0$. If B is unbounded then $\exists x, q$ s.t. for every $\lambda > 0, F(x + \lambda q) + M_3 \geq 0 \Rightarrow Fq \geq 0$, where $q \neq 0, Fq \geq 0$ since F is full rank and $e^T Fq = 0$, that is impossible since $Fq \geq 0$.

Now we have $X = UV = WH$, with $\text{rank}(U) = \text{rank}(W) = 3$, so they span the same column space $\Rightarrow W = U \cdot C \geq 0, UV = UCH \Rightarrow V = CH \Rightarrow \bar{V} = \bar{C}H \Rightarrow A = \text{conv}(\bar{V}) \subseteq \text{conv}(\bar{C})$ and $e^T = e^T W = e^T UC \Rightarrow e^T = e^T C, UC \geq 0 \Rightarrow$ every \bar{c}_i column of \bar{C} belongs to $B \Rightarrow A \subseteq \text{conv}(\bar{C}) \subseteq B \rightarrow$ it is a solution to NPP(3).

Vice versa, suppose $A \subseteq E \subseteq B \subseteq \mathbb{R}^2$ is a solution of NPP(3) with vertices $\{\bar{c}_1, \bar{c}_2, \bar{c}_3\} = \bar{C}$. Complete to C such that $e^T C = e^T$. Since $E \subseteq B$, we have $UC \geq 0$ and we take it as $W = UC$. Moreover, $A \subseteq E \Rightarrow V = \bar{C}H$ for some $H \geq 0, e^T H = e^T$, and $V = CH$ since $e^T V = e^T CH = e^T H = e^T \Rightarrow X = UV = UCH = WH$. Moreover, V is full rank, so C is full rank and as a consequence $\text{rank}(W) = \text{rank}(U) = 3$

- For $\text{rank}(X) \geq 4$ finding $\min \{r \mid X = WH, W, H \geq 0, \text{rk}(W) = \text{rk}(H) = r\}$ is NP-Hard
- If r is fixed, solving $X = WH, \text{rk}(W) = \text{rk}(H) = r$ takes $O(dn)^{O(r)}$ for some constant c , that is way too large.
- As of now, there is no algorithm to solve NMF in $(mn)^{o(r)}$

→ we solve instead $\min_{W, H \geq 0} \|X - WH\|_F^2$ for fixed r , that takes usually $O(dnr)$ per iteration with iterative solvers, but may converge to local minima.

• Another problem is that even if X has rational entries, its NMF W, H may have irrational entries, not approx. by a machine

Notice that NMF with $r=1$ is solved easily: take $X \sim m \cdot v^T$ best approx. by SVD. Then, since $X \geq 0, \|X - uv^T\|_F^2 \geq \|X - m \cdot v^T\|_F^2 \geq \|X - |uv^T|\|_F^2 \Rightarrow X \sim |uv^T|$ is the sol. to 1-NMF (akin to Perron-Frob.) but if $X \neq 0$ then even 1-NMF is NP-Hard.

Identifiability

The NMF $X = WH$ is very not unique. In fact, for every invertible diagonal D with positive diag. entries, and any permut. P

$$X = WH = (WDP)(P^T D^{-1} H)$$

but in application we do not care about the order of the features in W and we can always renormalize them to, for example, unit ℓ_1 norm. As a consequence we define

Uniqueness: We say that the NMF $X = WH$ is unique when

$$X = \tilde{W} \tilde{H} \text{ differ by scaling } D \text{ and permutation } P$$

Scaling and Permutation, in fact, are not the only problem. Even when $\text{rank}(W) = r$, there may be an invertible Q s.t.

$$X = WH = (WQ)(Q^{-1}H), \quad WQ \geq 0, \quad Q^{-1}H \geq 0$$

For example, $X = \begin{bmatrix} 1 & 2 \\ 1 & 1 \end{bmatrix} = X \cdot I = (X \cdot X^{-1})(X \cdot I) = IX$.

Uniqueness is important for the application and for stability since

Theorem 12.1 IF $X = WH$ NMF is unique, then for \tilde{X} in a neigh. of X , $\tilde{X} = \tilde{W} \tilde{H}$ NMF is unique and $(\tilde{W}, \tilde{H}), (\hat{W}, \hat{H})$ are close up to scaling and permutations.

For $\text{rank}(X) \leq 2$, it is easy:

Theorem 12.2 For $\text{rank}(X) \leq 2$, then $X = WH$ of size 2 always exists. It is unique iff X has a pd diagonal 2×2 submatrix.

To identify the 2×2 NMF it is enough to

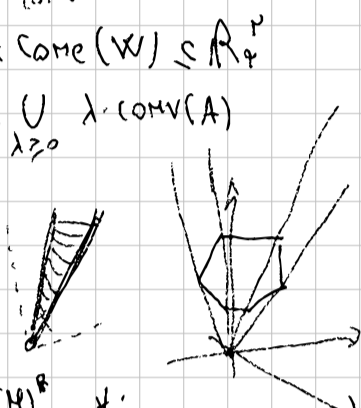
- Normalize X with ℓ_1 norm (after removing 0-columns)
- Take $y_1 = x_i$ with greatest ℓ_1 norm
- Take $y_2 = x_i$ furthest from y_1 , i.e. $\max \|x_i - y_1\|$
- identify j s.t. $y_{1,j} \neq y_{2,j}$
- for every i , solve $\begin{bmatrix} y_{1,i} & y_{2,i} \\ 1 & 1 \end{bmatrix} h_i = \begin{bmatrix} x_{i,1} \\ x_{i,2} \\ 1 \end{bmatrix}$
- $W = [y_1, y_2], H = [h_i]$

Let's go back to $X = WH \Leftrightarrow \text{Cone}(X) \subseteq \text{Cone}(W) \subseteq \Delta^r$

$$\text{where } \text{Cone}(A) = \{Ax \mid x \in \mathbb{R}_+^r\} = \bigcup_{\lambda \geq 0} \lambda \cdot \text{conv}(A)$$

Dual Cone: Given a cone C , then

$$C^* = \{x \mid y^T x \geq 0 \forall y \in C\}$$



Notice that $X = WH \geq 0 \Rightarrow w_i^T \in \text{Cone}(H)^*$ $\forall i$

$$\Rightarrow \text{Cone}(W^T) \subseteq \text{Cone}(H)^*$$

Separability $A \in \mathbb{R}^{r \times m}$ is separable if equivalently

- $\text{Cone}(A) = \mathbb{R}_+^r$
- $\text{Cone}(A)^* = \mathbb{R}_+^m$
- \exists subm. of A that is $D \cdot P$

Separable: We say that a NMF $X = WH$ is separable if H is sep.

Notice that H separable $\Leftrightarrow \Delta^m$ separable, so we can work with H ℓ_1 -norm. In this case, H separable $\Leftrightarrow I_r$ is a subm. of H up to permutation, or also said, if W is composed by col. of X .

In applications this has different names, for example in hyperspectral imaging, $X = WH$ is separable if for every material W there exists a pixel made only of that material (pure-pixel assumption). In text mining, sep. of W it means that each topic has a word that is exclusive of that topic (anchor word). In source sep. it means that each note must be played alone in a short span of time.

Theorem 12.3 IF $X = WH$, where $X \in \mathbb{R}^{d \times m}$, $W \in \mathbb{R}^{d \times r}$, $H \in \mathbb{R}_+^{r \times m}$ separable, then this is the unique separable MF.

Proof IF $X = WH = \tilde{W} \tilde{H}$ with H, \tilde{H} separable, then W and \tilde{W} are made of columns of X , but then $\text{Cone}(X) \supseteq \text{Cone}(W) \supseteq \text{Cone}(\tilde{W})$ and the same with \tilde{W} . Since W, \tilde{W} are full rank, their columns are the vertices of $\text{Cone}(X)$ up to scaling and permutation, since W is full rank, then H is also unique \square

\leadsto This has the advantage that retrieving a separable MF is doable in poly time or text mining.

\leadsto In some application, like face features extraction, is W to be sep. that has the same properties. The drawback is that separability is very strong and does not hold in many applications, so we need something weaker

Sufficiently Scattered Condition $A \in \mathbb{R}_+^{r \times m}$ is SSC if

$$\text{SSC1} \cdot C := \{x \in \mathbb{R}_+^r \mid e^T x \geq \sqrt{r-1} \|x\|\} \subseteq \text{Cone}(A)$$

$$\text{SSC2} \cdot \text{IF } Q \text{ orth. and } \text{Cone}(A) \subseteq \text{Cone}(Q), \text{ then } Q \text{ is a perm.}$$

Geometrically C is the cone of C_1 , the circle in Δ^r tangent to \mathbb{R}_+^r . In particular, $C \subseteq \mathbb{R}_+^r$

The idea is that for separability, $\text{Cone}(A) = \mathbb{R}_+^r$ that is, it must contain $\{e_1, e_2, \dots, e_r\}$ among its columns. With SSC, it must contain a smaller C .

SSC2 is saying that there is no other "rotated" Δ^r containing $\text{Cone}(A)$ except for Δ^r itself. This is satisfied for example when C is contained in the interior part of $\text{Cone}(A)$.

SSC is a sufficient condition for identifiability for certain problem with minimal, but to prove it we need some properties of dual cones.

• given two convex cones, $A \subseteq B \Leftrightarrow B^* \subseteq A^*$

• $C^* = \{y \in \mathbb{R}^r \mid e^T y \geq \|y\|\}$ is the circular cone circumscribed to Δ^r , and it contains also $y \notin \mathbb{R}_+^r$.

\leadsto bad news: checking SSC is NP-hard

\leadsto good news: IF H is randomly generated and each row has at least $r-1$ zeros, then it is SSC with high prob.

\leadsto ? news: Uniqueness of NMF do not require SSC of W, H necessarily

Minimum Volume Recall that when $X = WH$ ℓ_1 -norm, then

$$\text{Conv}(X) \subseteq \text{Conv}(W) \subseteq \Delta^r$$

but when the decomposition is not unique, one may try to find the W with the minimum volume that contain $\text{Conv}(X)$.

The volume of $\text{Conv}(W \cup \{0\})$ when $W \in \mathbb{R}^{d \times r}$, $d \times r$ is

$$\text{Vol}(W) = \frac{1}{r!} \sqrt{\det(W^T W)}$$

Theorem 12.4 Let $X = WH^*$ where $e^T H^* = e^T$, $H^* \in \mathbb{R}_+^{r \times m}$ SSC and $X \in \mathbb{R}^{d \times m}$, $W \in \mathbb{R}^{d \times r}$. Then (W, H^*) is the unique solution to

$$\min_{W \in \mathbb{R}^{d \times r}} \text{Vol}(W) : X = WH^*, e^T H^* = e^T, H^* \in \mathbb{R}_+^{r \times m}, W \in \mathbb{R}^{d \times r}$$

Proof Let $X = W^* H^* = WH^*$ where (W, H^*) is also feasible. Since $r = \text{rank}(X)$, then $W = W^* Q^{-1}$, $H^* = Q H^*$. Since H^* is full rank and $e^T H^* = e^T$, then $e^T = e^T H^* = e^T Q H^* \Rightarrow e^T Q = e^T$. Moreover, $H^* = Q H^* \geq 0 \Rightarrow \text{Cone}(Q^T) \subseteq \text{Cone}(H^*)^* \subseteq C^* (\text{SSC2}) = \{y \mid y^T e \geq \|y\|\}$

\Rightarrow $q_j^T e \geq \|q_j\| \forall j$. Then spm

$$|\det Q|^2 = \pi \|q_1\|^2 \dots \|q_r\|^2 = \pi \|q_1\|^2 \dots \|q_r\|^2 \leq \pi \|q_1\|^2 \dots \|q_r\|^2$$

$$|\det Q| \leq \pi \|q_1\| \dots \|q_r\| \leq \pi \left(\frac{e^T q_1}{r}\right)^r = \left(\frac{e^T q_1}{r}\right)^r = 1$$

But since W^* has the least volume, $\det(W^* W^*) \leq \det(W^T W) =$

$$= \det(Q^T W^* W^* Q^{-1}) = \det(W^* W^*) \cdot \det(Q)^{-2} \Rightarrow \det(Q) \leq 1$$

so $|\det Q| = 1$ and by AM-GM above, $q_j^T e = 1 = \|q_j\| \forall j$ and moreover

$$Q = I, Q^T Q = I \Rightarrow Q^T \text{ orth.}, \text{ but } \text{Cone}(Q^T)^* = \{y \mid Qy \geq 0\} = \{Q^T z \mid z \geq 0\} = \text{Cone}(Q^T), \text{ so}$$

$$\text{Cone}(Q^T) \subseteq \text{Cone}(H^*)^* \Rightarrow \text{Cone}(H^*) \subseteq \text{Cone}(Q^T) \xrightarrow{\text{SSC2}} Q^T \text{ permutation } \square$$

\leadsto Open question: $\text{minvol}(W)$ is poly? How does it fare with noise?

Possible solution: • maximum ellipsoid inscribed in W (exp)

• max vol of dual polar

Problem: $e^T H^* = e^T$ might not be a good condition when there are low norm x_i with loose connection with the features W . A better condition that is also more robust w.r.t noise is $e^T W = e^T$. Theorem 12.4 holds even for the modified min vol problem

$$\min_{W, H} \det(W^T W) : X = WH, e^T W = e^T, W \in \mathbb{R}^{d \times r}, H \in \mathbb{R}_+^{r \times m}$$

Recall that k-means can be recast as

$$\min_{C, H} \|X - CH\|_F^2 : C \in \mathbb{R}^{d \times k}, H \in \{0, 1\}^{k \times n}$$

$$e^T H = e^T$$

Notice that if $X \geq 0$, then $C \geq 0$ from k-means, since they are averaged points in X . Since $H \geq 0$, this is nothing else than an NMF. We can always bring X to positive by translation and if we normalize them in $\|\cdot\|_2$ norm, we obtain the classic

$$X = WH, \quad X, W, H \geq 0, \quad e^T X = e^T H = e^T, \quad e^T W = e^T$$

since in NMF H are not boolean, general NMF is much more flexible to different kind of representations.

A variant of k-means, called Spherical k-means, looks for directions W_i and clusters the data x_i whose angle with W_i are the least. In formulae, it is

$$\arg \max_{W \in \mathbb{R}^{d \times k}} \sum_i \max_{1 \leq l \leq k} \frac{W_{i,l}^T x_i}{\|W_{i,l}\| \|x_i\|} = SK(X)$$

so if X and w are normalized in $\|\cdot\|_2$, this is

$$\sum_i x_i^T W e_{z(i)} = \text{Tr}(X^T W H) = \frac{1}{2} [\|X - WH\|_F^2 - \|X\|_F^2 - \|WH\|_F^2]$$

$$\leadsto SK(X) = \arg \min_{H, W \in \mathbb{R}^{d \times k}} \|X - WH\|_F^2 : H \in \{0, 1\}^{k \times n}, e^T H = e^T$$

$$\|W_{i,l}\| = 1 \quad \forall i, l$$

Relax $X \geq 0$

$$\arg \min_{H, W \in \mathbb{R}_+^{d \times k}} \|X - WH\|_F^2 : H \in \mathbb{R}_+^{k \times n}, H H^T = I$$

notice: each col. of H has still only 1 entry > 0

ONMF: Assures that $\text{supp}(H_{i,:})$ are distinct and

$$\arg \min_{W \geq 0} \|X - WH\|_F^2 = X H^T$$

so it is equivalent to

$$\arg \min_{\substack{H \in \mathbb{R}_+^{k \times n} \\ H H^T = I}} \|X - X H^T H\|_F^2 = \arg \max_{\substack{H \in \mathbb{R}_+^{k \times n} \\ H H^T = I}} \|X H^T\|_F^2$$

The property $H H^T$ diag. assures that $\text{supp}(H_{i,:})$ are disjoint, so this is in itself a clustering algorithm. Moreover it tells us that $H^T H$ is a projector to some subspace when $H^T H = I$.

If we now reconsider SBM, we see that $A_n \sim E[A_n]$ where

$$E[A_n] = \begin{bmatrix} p_{11} & \dots & p_{1n} \\ \vdots & \ddots & \vdots \\ p_{n1} & \dots & p_{nn} \end{bmatrix} = \begin{bmatrix} \mathbb{1} \\ \vdots \\ \mathbb{1} \end{bmatrix} [P_{i,j}]_{i,j} \begin{bmatrix} \mathbb{1} & \dots & \mathbb{1} \end{bmatrix} = Z P Z^T$$

and $Z, P \geq 0 \leadsto$ This is a Tri-sym NMF

$$X \sim \begin{matrix} \text{mat} & \text{mat} & \text{mat} \\ \text{row} & \text{row} & \text{row} \end{matrix} W S W^T, \quad W, S \geq 0 \quad (S \text{ sym if } X \text{ sym})$$

Recall that for Spectral clustering we had $X \sim U \Lambda U^T$ where we needed to cluster the rows of U . Here instead $W \geq 0$ is more interpretable. In general clustering, W is the "belonging to cluster" matrix and S is instead the relations between different clusters.

Let's see an example of Topic Modeling:

In Text mining, we had A matrix of frequency of words in document, and $A = WH$ where W are the words in topics and H is the topic-to-documents matrix. The problem here is that if for example the documents are short (see Twitter or Facebook or X) they even if a post is discussing a topic, it won't use most of the words associated with that. For this reason, we form the co-occurrence matrix $X = A A^T$ that says how common is it that two words are in the same document. Now

$$X \sim W S W^T$$

has the usual word-to-topic W and also a topic-to-topic correlation S

KKT Conditions for NMF

$$\min_{\substack{W \in \mathbb{R}_+^{d \times r} \\ H \in \mathbb{R}_+^{r \times n}} D(X, WH) \quad \text{where } X \in \mathbb{R}_+^{d \times n}$$

in unconstrained optimization, the 1st order conditions are

$$\nabla_W D(X, WH) = 0, \quad \nabla_H D(X, WH) = 0$$

but with inequalities constraints we need something different. The idea is that if $W_{i,j} = 0$, then we don't need that $\partial/\partial W_{i,j} D(X, WH) = 0$; we just need to be sure that $D(X, WH)$ does not decrease if we increase $W_{i,j}$, i.e. we just need $\partial/\partial W_{i,j} D(X, WH) \geq 0$.

If instead $W_{i,j} > 0$, we actually need the gradient to be 0. As a consequence, this is written as

$$\begin{cases} W_{i,j} > 0, \quad \nabla_W D(X, WH) \geq 0, \quad W_{i,j} \cdot \nabla_W D(X, WH) = 0 \\ H_{i,j} > 0, \quad \nabla_H D(X, WH) \geq 0, \quad H_{i,j} \cdot \nabla_H D(X, WH) = 0 \end{cases}$$

↑ Karush-Kuhn-Tucker conditions

For $D(X, WH) = \|X - WH\|_F^2$, we get for example

$$\nabla_W \|X - WH\|_F^2 = -2XH^T + 2WHH^T \rightsquigarrow (WH - X)H^T \geq 0, \quad W_{i,j} \cdot (WH - X)_{i,j} = 0$$

and similar with H . Notice that $W=0, H=0$ is always a solution.

Another one is the Kullback-Leibler (KL) divergence, i.e.

$$D(A \| B) = \sum_{i,j} A_{i,j} \log\left(\frac{A_{i,j}}{B_{i,j}}\right) - A_{i,j} + B_{i,j}, \quad \begin{matrix} D(x \| 0) = \infty \quad (x > 0) \\ D(0 \| x) = x \\ D(0 \| 0) = 0 \end{matrix}$$

$$\frac{\partial D(x \| y)}{\partial y} = 1 - \frac{x}{y} \quad \text{so we can compute the grad. of the error:}$$

$$\nabla_W D(X \| WH) = \nabla_W \sum_{i,j} D(x_{i,j} \| (WH)_{i,j}) = \left[\sum_{i,j} \left(1 - \frac{x_{i,j}}{(WH)_{i,j}}\right) \frac{\partial (WH)_{i,j}}{\partial W_{k,h}} \right]_{k,h}$$

$$= \left[\sum_j \left(1 - \frac{x_{k,j}}{(WH)_{k,j}}\right) H_{j,h} \right]_{k,h} = e e^T H^T - (X ./ WH) H^T = \left(\frac{WH - X}{WH}\right) H^T$$

$$\nabla_H D(X \| WH) = \left[\sum_i \left(1 - \frac{x_{i,h}}{(WH)_{i,h}}\right) W_{i,k} \right]_{k,h} = W^T e e^T - W^T (X ./ WH)$$

The KKT conditions are thus

$$\begin{cases} \frac{WH - X}{WH} H^T \geq 0, \quad W \odot \left(\frac{WH - X}{WH} \cdot H^T\right) = 0 \\ W^T \frac{WH - X}{WH} \geq 0, \quad H \odot \left(W^T \cdot \frac{WH - X}{WH}\right) = 0 \end{cases}$$

Again, $(0,0)$ is a solution.

→ In general any rank-deficient stationary (W, H) are saddle points so they are not stable for most algorithms

Depending on the model of our noise/perturbation N s.t. $X = WH + N$, it is opportune to choose different $D(\cdot, \cdot)$.

Noise	$N(0, \Sigma)$	Uniform	Poisson
Distance	$\sum \frac{1}{2} (x_{i,j} - (WH)_{i,j})^2$	$\max_{i,j} x_{i,j} - (WH)_{i,j} $	$D(X \ WH)$

(we always suppose the entries of N are indep.)

Gaussian or Uniform noise may not be useful when X is sparse, since the noise may generate many negative entries. In these cases, it is better to use the divergence KL. In particular, images tend to be dense, so audio/text \rightsquigarrow KL, Images \rightsquigarrow Prob. or ℓ^p , Notice that

Theorem 13.1 Given $X \in \mathbb{R}_+^{d \times n}$ and $r > 0$, any stationary point (\bar{W}, \bar{H}) of $\min_{W, H} D(X \| WH)$ satisfies

$$Xe = WH e, \quad e^T X = e^T WH$$

Proof a stationary point satisfies $\bar{W} \odot \nabla_W D(X \| \bar{W}\bar{H}) = 0$ since if $\bar{W}_{i,j} \neq 0$, then the derivative must be zero (KKT condition) so

$$\left(\bar{W} \odot \frac{WH - X}{WH} H^T\right) e = 0 \Leftrightarrow \left(\bar{W} \odot e e^T H^T\right) e = \left(\bar{W} \odot \frac{X}{WH}\right) H^T e$$

$$\Leftrightarrow \text{diag}(\bar{W}) (e e^T H^T)^T = \text{diag}(\bar{W} \bar{H} \left(\frac{X}{WH}\right)^T)$$

$$\Leftrightarrow \bar{W} \bar{H} e = \left(\bar{W} \bar{H} \frac{X}{WH}\right) e = Xe, \quad \text{and the same with } H \quad \square$$

Lemma 13.2 Given $X \in \mathbb{R}_+^{d \times n}$, the unique solution to $r=1$ of $\min D(X \| WH)$ is $W = Xe, H^T = e^T X / e^T X e$ up to scaling.

Proof By th. 13.1, a solution (w, h) satisfies $Xe = w(h^T e)$ and $e^T X = (e^T w) h^T$ so we can take $w = Xe, h^T = e^T X / e^T X e$, and it is the unique solution up to scaling. \square

Problem: For every $D(\cdot, \cdot)$ "famous" i.e. KL, $\|\cdot\|_F, \|\cdot\|_2, \|\cdot\|_\infty$, etc. NMF is NP-hard, so we need good algorithm for approx. in $r \geq 2$. ℓ_1 -norm is also NP-hard for $r \geq 2$.

Other results: For stationary points (W, H) in $\|\cdot\|_F$ it holds that

$$\|X - WH\|_F^2 = \|X\|_F^2 - \|WH\|_F^2$$

$$\text{so } \|X\|_F^2 \geq \|WH\|_F^2$$

Lesson 9

SPA : Successive Projection Algorithm

Suppose $X = WH$, where H is separable and substochastic, i.e. $H \geq 0$, $e^T H \leq e^T$, $J_r = H(:, j)$ for $1 \leq j \leq r$. Suppose moreover that W is full rank. Notice that $W = X(:, J)$ and thus $\text{conv}(X) \subseteq \text{conv}(W) \subseteq \text{conv}(X) \Rightarrow \text{conv}(W) = \text{conv}(X)$ and the columns of W are the vertices of $\text{conv}(X)$.

SPA recovers the vertices of $\text{conv}(X)$ sequentially. In fact notice that the max ℓ_2 -norm column of X must be one of such vertices since the norm is convex, and when we consider it on $\text{conv}(X)$, it takes maximum only on the vertices. In formulae,

$$\|X(:, j)\| = \|WH(:, j)\| \leq \sum_i H(i, j) \cdot \|W(:, i)\| \leq \left(\max_i \|W(:, i)\| \right) \left(\sum_i H(i, j) \right) \leq \max_i \|W(:, i)\| \leq \max_j \|X(:, j)\|.$$

For a max- ℓ_2 norm $X(:, j)$ all of these are equalities, but the first is tr. ineq. that is = only when the $W(:, i)$ for which $H(i, j) \neq 0$ are collinear and W is full rank, so that's possible only when there's at most one coeff. $H(i, j) \neq 0$. The 2nd ineq. is = only when $H(i, j) \leq 1$, meaning that $H(i, j) = 1$ and so $X(:, j) = W(:, i)$.

After recovering k vertices $\{w(:, 1), \dots, w(:, k)\}$, SPA projects over the orthogonal space to all of these vertices, called V_k . Then it repeats the same procedure, finding the $k+1$ -index since $P_{V_k} X = (P_{V_k} W) \cdot H$ and we can remove the 0 columns from W and the corresponding rows from H to obtain a separable H and a rank $(r-k)$ $P_{V_k} W$.

SPA

```

Given  $X \in \mathbb{R}^{d \times m}$  and  $r \geq 0$ , let  $J = \emptyset$ 
while  $|J| < r$ 
    Find the max  $\|X(:, k)\| = \max_i \|X(:, i)\| = v$ 
     $J = J \cup \{k\}$ 
     $X = (I - \frac{v v^T}{\|v\|^2}) X$ 
end
    
```

Notice that X, W does not need to be ≥ 0 , but $X \geq 0 \Rightarrow W \geq 0$, so it is an NMF. The complexity is $O(dnr)$; and if X is sparse it is $O(r \cdot \text{nnz}(X))$.

This is also robust to noise in the sense that given $X = WH + N$ where H is separable and substochastic, W is full rank and

$$\epsilon = \max_j \|N(:, j)\| = O\left(\frac{\omega_{\min}(W)}{\sqrt{r} \kappa^2(W)}\right)$$

Then SPA can identify W with an ℓ_2 error up to $O(\epsilon \kappa^2(W))$. With further improvements (Fast Anchor Words / preconditioning) one gets

$$\epsilon = O\left(\frac{\omega_{\min}(W)}{\sqrt{r} \kappa(W)}\right), \text{ error} = O(\epsilon \kappa(W))$$

\rightarrow we get anyway an error proportional to $\frac{\omega_{\min}(W)}{\sqrt{r}}$.

Notice that if W_1 is very close to the span of the rest W_2, \dots, W_r then after the projection $\{w_1, \dots, w_r\}$, the projected w_1 will have very small norm, thus a random perturbation can easily prevent us to find w_1 . Notice that

$$\min_k \min_{\alpha} \|w_k - \sum_{i \neq k} \alpha w_i\| = \min_{k, \alpha} \|W \begin{bmatrix} \alpha \\ \vdots \\ 1 \\ \vdots \\ \alpha \end{bmatrix}\| \leq \min_{\|x\|_1=1} \|Wx\| = \omega_r(W)$$

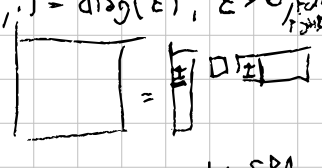
This is the reason why ϵ must be bounded by $\omega_{\min}(W)$.

\rightarrow It is also very sensitive to outliers. Smoothed/random version exist for this reason.

\rightarrow It can be proved that SPA is equivalent to $\max_{|J|=r} \text{Vol}(X(:, J))$ [33]

Tri-NMF

Suppose now we want the tri-sym NMF $A = WSW^T$ in the case W is separable, i.e. $W(:, j) = \text{diag}(z)$, $z \geq 0$ and everything is nonnegative, and A, S are sym. In this case there is an exact poly algorithm.



In fact, $A = (WS)W^T$ is a separable NMF for which we can apply SPA to find S : $A(:, j) = (WS) \cdot \text{diag}(z)$. Moreover,

$$A(:, j) = \text{diag}(z) S \text{diag}(z) \Rightarrow S = \text{diag}(z)^{-1} A(:, j) \text{diag}(z)^{-1} \quad (1)$$

and since we can always suppose $e^T W = e^T$, then

$$A(:, j) e = \text{diag}(z) \cdot S \cdot W^T e = A(:, j) \text{diag}(z)^{-1} e = A(:, j) \begin{bmatrix} 1/z_1 \\ \vdots \\ 1/z_j \end{bmatrix}$$

\Rightarrow we can compute z from $\begin{bmatrix} 1/z_i \\ \vdots \\ 1/z_i \end{bmatrix} = A(:, j)^{-1} A(:, j) e$ and then compute S by (1) and $A(:, j) = \text{diag}(z) S W^T \Rightarrow W^T = S^{-1} \text{diag}(z)^{-1} A(:, j)$.

Tri-sym NMF

```

Given  $A \in \mathbb{R}_+^{m \times m}$  symmetric and  $r \geq 0$ 
Compute the index set  $J$  from SPA( $A, r$ )
Compute  $y = A(:, j)^{-1} A(:, j) e$ 
Compute  $S = \text{diag}(y) \cdot A(:, j) \cdot \text{diag}(y)$ 
Compute  $W^T = S^{-1} \text{diag}(y) A(:, j)$ 
    
```

\rightarrow one could compute once $B = A(:, j)^{-1} A(:, j)$, so that $y = B e$ and $W^T = \text{diag}(y)^{-1} B$. If there's noise, instead of inverting $A(:, j)$ one solves

$$y = \arg \min \|A(:, j) e - A(:, j) x\|, \quad W^T = \arg \min \|A(:, j) \text{diag}(y)^{-1} A(:, j) - A(:, j)\|$$

or with ℓ_1 - ℓ_1 instead of $\|\cdot\|$.

\rightarrow In case one has SSC on W , one need to solve a min-vol Tri-NMF

In general $\min_{\substack{W \geq 0 \\ H \geq 0}} D(X, WH)$ is solved with alternating method i.e. optimize over W with H fixed and viceversa. So from now on we focus into solving

$$\min_{H \geq 0} D(X, WH) : W, X \geq 0$$

Active Set

Both for KL divergence and $\| \cdot \|_F^2$, the $D(\cdot, \cdot)$ can be decomposed as

$$D(X, WH) = \sum_i D(x_i, Wh_i)$$

so we can analyze separately $\min_{h \geq 0} D(x, Wh)$. The KKT conditions

$$h \geq 0, \nabla_h D(x, Wh) \geq 0, h \odot \nabla_h D(x, Wh) = 0$$

in this case are Necessary and Sufficient for global minimality. We can thus call $I(h) = \{i : h_i > 0\}$ and notice that for the solution h^* we need

$$\nabla_h D(x, Wh^*) \Big|_{I(h^*)} = 0$$

- In case of $\| \cdot \|_F^2$, we get $[W^T(WH - x)]_I = 0$, i.e.

$$W(i, I)^T W(i, I) h(I) = W(i, I)^T x$$

- In case of KL, we get $[W^T((WH - x) \odot Wh)]_I = 0$ i.e.

$$W(i, I)^T e = W(i, I)^T \text{diag}(x) \cdot \frac{1}{W(i, I)h(I)}$$

\leadsto In both cases, $h(I)$ and thus h^* is solvable as long as we know $I(h^*)$. So we can try to find I , that are finite but in theory exponential in r , so they tend to be slow even if accurate. The usual methods tend to update I to find the optimal one while looking for indices to add/remove in order to lower the error as much as possible.

MU

One of the, if not the, first algorithm proposed for NMF. Recall the KKT conditions for $\| \cdot \|_F^2$

$$\nabla_h \|X - WH\|_F^2 = 2W^T(WH - X) \geq 0, H \odot (W^T(WH - X)) = 0$$

so if $(\nabla_h)_{i,j} > 0$ and $H_{i,j} > 0$, then we have to decrease $H_{i,j}$

and we know that $(\nabla_h)_{i,j} > 0 \Rightarrow (W^T WH)_{i,j} > W^T x$

$\Rightarrow (W^T x)_{i,j} / (W^T WH)_{i,j} \leq 1$. Viceversa, $(\nabla_h)_{i,j} < 0$ then we have to increase $H_{i,j}$

and $(\nabla_h)_{i,j} < 0 \Rightarrow (W^T x)_{i,j} / (W^T WH)_{i,j} > 1$, so

$$H \leftarrow H \odot \frac{W^T x}{W^T WH} \quad \text{and} \quad W \leftarrow W \odot \frac{XH^T}{WHH^T}$$

MU

Given $X \in \mathbb{R}_+^{d \times n}$ and $\epsilon > 0$ initialize $W, H > 0$

Repeat until convergence

$$H \leftarrow H \odot \frac{W^T x}{W^T WH + \epsilon \cdot E}, \quad W \leftarrow W \odot \frac{XH^T}{WHH^T + \epsilon \cdot E}$$

Here $\epsilon > 0$ is a small enough constant to avoid dividing by zero, and E is the all-ones matrix.

There is also the version for the KL-div.

MU-kl

Given $X \in \mathbb{R}_+^{d \times n}$ and $\epsilon > 0$ initialize $W, H > 0$

Repeat until convergence

$$H \leftarrow H \odot \frac{W^T X / WH}{W^T E + \epsilon \cdot E}, \quad W \leftarrow W \odot \frac{X / WH^T}{E H^T + \epsilon \cdot E}$$

Theorem 14.1 MU makes $D(X, WH)$ decrease at each step both for $\| \cdot \|_F^2$ and $D(\cdot, \cdot)$

\leadsto Notice that if W or H have zero entries at some point, then they will continue having zeros forever on those entries. This says that MU generates sparse solutions, but at the same time, it tends to get stuck in local minima because it cannot get out of the zero entries, so it has convergence issues. $O(mnr)$ comp. cost for dense, or $O(r \cdot \text{nnz}(X))$ for sparse

Coordinate Gradient Descent

Let's split the problem even more.

$$\min_{H \geq 0} \|X - WH\|_F^2 = \min_{H \geq 0} \sum_i \|x_i - Wh_i\|^2$$

$$= \min_{H \geq 0} \sum_i \left\| \underbrace{x_i - \bar{w}_i}_{v} - \underbrace{h_{j,i}}_a \underbrace{w_j}_{z} \right\|^2$$

where \bar{w} is the submatrix of W where we removed column j and in \bar{h}_i we removed entry j . If we now fix everything except for $h_{j,i}$, we see that it becomes an easy problem of the kind

$$\arg \min_{\alpha \geq 0} \|v - \alpha z\|^2 = \arg \min_{\alpha \geq 0} \alpha^2 \|z\|^2 - 2\alpha \langle v, z \rangle$$

$$= \max \left\{ 0, \frac{\langle v, z \rangle}{\|z\|^2} \right\}$$

$$\text{so } h_{j,i}^* = \max \left\{ 0, \frac{w_j^T x_i - w_j^T \bar{w}_i}{\|w_j\|^2} \right\}$$

$$\leadsto H(j, i) \leftarrow \max \left\{ 0, \frac{1}{\|w_j\|^2} (w_j^T X - w_j^T \bar{W} \bar{H}) \right\}$$

HALS

Given $X \in \mathbb{R}_+^{d \times n}$ $r > 0$, initialize $W, H \geq 0$

Repeat until convergence

Compute XH^T and HH^T

$$\text{for } k=1, \dots, r$$

$$W(k, :) = \max \left\{ 0, \frac{(XH^T)(k, :) - \sum_{\ell \neq k} W(k, \ell) \cdot (HH^T)(\ell, k)}{(HH^T)_{k, k}} \right\}$$

Compute $W^T X$ and $W^T W$

$$\text{for } k=1, \dots, r$$

$$H(k, i) = \max \left\{ 0, \frac{(W^T X)(k, i) - \sum_{\ell \neq k} H(\ell, i) (W^T W)(k, \ell)}{(W^T W)_{k, k}} \right\}$$

One can accelerate the method by repeating the internal loop of update of W, H multiple time, since the computations of $XH^T, HH^T, W^T X, W^T W$ are the actual bottleneck. This is also highly parallelizable. $O(mnr)$ / $O(r \cdot \text{nnz}(X))$

\leadsto This is an effective Alternating Method, that solves perfectly each internal convex problem, so it is convergent by theorem...

It is moreover one of the most performing algorithms known nowadays and it guarantees the error to decrease at each step.

There are many other: Fast Pro. Gradient Method, or ASMM on

$$\min_{\substack{H \in \mathbb{R}^r \\ \varphi \in \mathbb{R}_+^r}} \|X - WH\|^2 : \varphi = h$$

$$\text{or ASMM directly on } \min_{\substack{W, H \\ U, V}} \|X - WH\|, \quad W=U \geq 0, H=V \geq 0$$

but all of these do not guarantee the error to decrease at each step