Log-determinant Non-negative Matrix Factorization via Successive Trace Approximation Optimizing an optimization algorithm

Andersen Ang

Mathématique et recherche opérationnelle UMONS, Belgium

Email: manshun.ang@umons.ac.be Homepage: angms.science

File created : November 30, 2017 Last update : May 18, 2018

Joint work with my supervisor : Nicolas Gillis (UMONS, Belgium)

Overview

- Background and the research problem
 - Non-negative Matrix Factorization
 - Separable NMF
 - Why Separable NMF is not enough
 - The minimum volume criterion

2 Minimum volume log-determinant (logdet) NMF

- The logdet regularization
- A logdet inequality
- An upper bound of logdet

Solving logdet NMF - Successive Trace Approximation

- STA algorithm
- The NNQPs
- Optimizing STA
- Experiments

Discussions, extensions and directions

- Coordinate acceleration by randomization
- On robustness and robust formulations
 - ★ De-noising norm
 - ★ Iterative Reweighted Least Sqaures
- Theoretical convergences
- Comparisons with benchmark algorithms (skiped here)
- On selecting the regularization parameter λ
- On automatic detection of r
- Further acceleration with weighted column fitting norms

What is Non-negative Matrix Factorization (NMF) ?

Given $\mathbf{X} \in \mathbb{R}^{m \times n}_+$ and integer r, find matrices $\mathbf{W} \in \mathbb{R}^{m \times r}_+, \mathbf{H} \in \mathbb{R}^{r \times n}_+$ s.t.

$\mathbf{X}=\mathbf{W}\mathbf{H}.$

1 This is called : extact NMF and it is NP-hard (Vavasis2007).

What is Non-negative Matrix Factorization (NMF)?

Given $\mathbf{X} \in \mathbb{R}^{m \times n}_+$ and integer r, find matrices $\mathbf{W} \in \mathbb{R}^{m \times r}_+, \mathbf{H} \in \mathbb{R}^{r \times n}_+$ s.t.

$\mathbf{X}=\mathbf{W}\mathbf{H}.$

- **1** This is called : extact NMF and it is NP-hard (Vavasis2007).
- We consider
 - low rank/complexity model $1 \le r \le \min\{m, n\}$.
 - approximate NMF

$$[\mathbf{W},\mathbf{H}] = \operatorname*{arg\,min}_{\mathbf{W} \ge 0,\mathbf{H} \ge 0} \frac{1}{2} \|\mathbf{X} - \mathbf{W}\mathbf{H}\|_{F}^{2}.$$

What is Non-negative Matrix Factorization (NMF) ?

Given $\mathbf{X} \in \mathbb{R}^{m \times n}_+$ and integer r, find matrices $\mathbf{W} \in \mathbb{R}^{m \times r}_+, \mathbf{H} \in \mathbb{R}^{r \times n}_+$ s.t.

$\mathbf{X}=\mathbf{W}\mathbf{H}.$

- This is called : extact NMF and it is NP-hard (Vavasis2007).
- We consider
 - low rank/complexity model $1 \le r \le \min\{m, n\}$.
 - approximate NMF

$$[\mathbf{W},\mathbf{H}] = \underset{\mathbf{W} \ge 0,\mathbf{H} \ge 0}{\operatorname{arg\,min}} \frac{1}{2} \|\mathbf{X} - \mathbf{W}\mathbf{H}\|_{F}^{2}.$$

- Such minimization problem is
 - also NP-hard
 - a non-convex problem
 - an ill-posed problem

What is Non-negative Matrix Factorization (NMF)?

Given $\mathbf{X} \in \mathbb{R}^{m \times n}_+$ and integer r, find matrices $\mathbf{W} \in \mathbb{R}^{m \times r}_+, \mathbf{H} \in \mathbb{R}^{r \times n}_+$ s.t.

$\mathbf{X}=\mathbf{W}\mathbf{H}.$

- This is called : extact NMF and it is NP-hard (Vavasis2007).
- We consider
 - low rank/complexity model $1 \le r \le \min\{m, n\}$.
 - approximate NMF

$$[\mathbf{W},\mathbf{H}] = \operatorname*{arg\,min}_{\mathbf{W} \ge 0,\mathbf{H} \ge 0} \frac{1}{2} \|\mathbf{X} - \mathbf{W}\mathbf{H}\|_F^2.$$

- Such minimization problem is
 - also NP-hard
 - a non-convex problem
 - an ill-posed problem
- **3** Assumptions (1) **W**, **H** full rank, (2) r is known.
- **③** Notation note : we use $\mathbf{W}\mathbf{H}$ instead of $\mathbf{W}\mathbf{H}^ op$

Given $\mathbf{X} \in \mathbb{R}^{m \times n}_+$ and integer r, find matrices $\mathbf{W} \in \mathbb{R}^{m \times r}_+, \mathbf{H} \in \mathbb{R}^{r \times n}_+$ s.t.

$$\mathbf{X} = \mathbf{WH}$$
 and $\underbrace{\mathbf{W} = \mathbf{X}(:, \mathcal{K})}_{\text{separability}}$.

Q Separable NMF = NMF + additional condition $\mathbf{W} = \mathbf{X}(:, \mathcal{K})$

- Meaning : W is some columns of X.
- ▶ 𝒴 : column index set

$$\blacktriangleright |\mathcal{K}| = r$$

Given $\mathbf{X} \in \mathbb{R}^{m \times n}_+$ and integer r, find matrices $\mathbf{W} \in \mathbb{R}^{m \times r}_+, \mathbf{H} \in \mathbb{R}^{r \times n}_+$ s.t.

$$\mathbf{X} = \mathbf{W}\mathbf{H}$$
 and $\underbrace{\mathbf{W} = \mathbf{X}(:, \mathcal{K})}_{\text{separability}}$.

Q Separable NMF = NMF + additional condition $\mathbf{W} = \mathbf{X}(:, \mathcal{K})$

- Meaning : W is some columns of X.
- K : column index set

$$\mid \mathcal{K} \mid = r$$

- In the second second
- Separable NMF :
 - XRAY (Kumar,2013)
 - SPA, SNPA (Gillis, 2013)

Illustrative examples



Figure: NMF, m = 5, n = 10, r = 3.



Figure: Separable NMF, $\mathcal{K} = \{8, 1, 3\}$. H has some special columns : only one '1', and other elements are 0.

"Pure pixels" assumption : matrix $\mathbf{H} = [\mathbf{I}_r \ \mathbf{H}']\mathbf{\Pi}$. Related terms : self-expressive, self-dictionary

Motivation — why study NMF

- In hyperspectral imaging application, $\mathbf{X}=$ image.
- W = absorption behaviour of materials : non-negative spectrum.
- r = # fundamental materials (e.g. rock, vegetation, water).
- H = abundance of materials : non-negative and sum-to-1.



Figure: Hypersectral images decomposition. Figure copied shamelessly from N. Gillis. Interpretation in short :

 $\mathbf{X} = \mathsf{data}, \ \mathbf{W} = \mathsf{basis}, \ r = \#\mathsf{basis}, \ \mathbf{H} = \mathsf{membership}$ of data w.r.t. basis

NMF has many other signal processing applications.

NMF vs SVD : SVD has lower fitting error (in fact SVD achieve the optimal fitting), but basis of SVD are not interpretable. Separable NMF basis comes form data, interpretable $\frac{10}{95}$

H tells the membership of data points in X w.r.t basis W. Column form expression of X = WH is

 $\mathbf{X}(:,\,j\,)=\mathbf{W}\mathbf{H}(:,\,j\,).$



 ${\bf H}$ tells the membership of data points in ${\bf X}$ w.r.t basis ${\bf W}.$

Column form expression of $\mathbf{X}=\mathbf{W}\mathbf{H}$ is

$$\mathbf{X}(:, j) = \mathbf{WH}(:, j).$$



H tells the membership of data points in X w.r.t basis W. Column form expression of X = WH is

$$\mathbf{X}(:, j) = \mathbf{WH}(:, j).$$





 \mathbf{H} tells the membership of data points in \mathbf{X} w.r.t basis \mathbf{W} . Column form expression of $\mathbf{X} = \mathbf{W}\mathbf{H}$ is

 $\mathbf{X}(:, j) = \mathbf{WH}(:, j).$

• Nonnegativity : $\mathbf{H}_{ij} \ge 0$ so it is in fact conical combination !



 ${\bf H}$ tells the membership of data points in ${\bf X}$ w.r.t basis ${\bf W}.$

Column form expression of $\mathbf{X} = \mathbf{W}\mathbf{H}$ is

$$\mathbf{X}(:, j) = \mathbf{WH}(:, j).$$



 ${\bf H}$ tells the membership of data points in ${\bf X}$ w.r.t basis ${\bf W}.$

```
Column form expression of \mathbf{X}=\mathbf{W}\mathbf{H} is
```

```
\mathbf{X}(:, j) = \mathbf{WH}(:, j).
```

• If columns of **H** are normalized : it means **W** is forming a *convex hull* encapsulating the data columns of **X**.



 \mathbf{H} tells the membership of data points in \mathbf{X} w.r.t basis \mathbf{W} . Column form expression of $\mathbf{X} = \mathbf{W}\mathbf{H}$ is

$$\mathbf{X}(:, j) = \mathbf{WH}(:, j).$$

• Algebarically, column normalization of **H** removes the scaling ambiguity of factorization, prevents huge **H** and super small **W** as

$$\mathbf{W}_{1}\mathbf{H}_{1} = \underbrace{\mathbf{W}_{1}\Lambda\Pi}_{\mathbf{W}_{2}}\underbrace{\Pi^{-1}\Lambda^{-1}\mathbf{H}_{1}}_{\mathbf{H}_{2}}$$

No normalization : convex hull \rightarrow conical hull \implies scaling ambiguity!

What has been done in literature

The problem statement : given the data points that has pure pixel (i.e. data points are **distributed in the entire** data subspace).

Goal : find the vertices = (1) find r (#vertices), (2) locate them.



Figure: A 2D PCA projection of a high dimensional data, showing the data points (black dots) encapsulated inside convex hull spanned by the generator (vertices).

This problem \subset Blind source identification with even distributed data Existing methods : XRAY, SPA, SNPA ...

This work : what if pure pixel $(\mathbf{H}_{ij} = 1)$ is hidden in data?

The problem statement : given the data points that are at least $1 - \theta$ away from the vertices (i.e. the pure pixel are hidden, data points are **not distributed in the entire** data subspace).

Goal – find the vertices : (1) find r (#vertices), (2) locate them.



Figure: A 2D PCA projection of a high dimensional data, showing the **data points** encapsulated inside convex hull spanned by the **generator vertices**.

This work : what if pure pixel $(\mathbf{H}_{ij} = 1)$ is hidden in data?

The problem statement : given the data points that are at least $1 - \theta$ away from the vertices (i.e. the pure pixel are hidden, data points are **not distributed in the entire** data subspace).

Goal – find the vertices : (1) find r (#vertices), (2) locate them.



Figure: A 2D PCA projection of a high dimensional data, showing the **data points** encapsulated inside convex hull spanned by the **generator vertices**.

- θ : "purity" of the data, $\theta \in [0 \ 1]$.
- $\theta = 1$: Separable NMF.
- Problem (1) is not considered here : r = # vertices (red dots) is assume known.

Existing algorithms aimed solving the Separable NMF ($\theta = 1$) work poorly on the problems with $\theta < 1$.



Figure: Results (2d PCA projection) of SNPA with decreasing θ . The dimensions are (m, n, r) = (8, 1000, 3).

Why separable NMF is not enough ... (2/2)

Due to the nature of high dimensional geometry, when (m, r) increase, the data points are getting more and more concentrated around the annulus of the origin and thus they are **not distributed in entire data subspace**. Making approches that use L_2 norm of data points (such as SNPA) less workable.



Figure: Results (2d PCA proj.) of SNPA with increasing (m, r). Red dots : ground truth vertices. Blue dots : estimated vertices. The dimensions are $(n, \theta) = (1000, 0.999)$.

• An idea from 1994¹ : fit a low rank convex hull with minimum volume.

 $^{^1\}mbox{Craig},$ Minimum-volume transforms for remotely sensed data. IEEE Trans. Geosci. Remote Sensing

 $^{^{2}}$ Lin, et. al, Identifiability of the simplex volume minimization criterion for blind hyperspectral unmixing: The no-pure-pixel case. IEEE Trans Geosci. Remote Sensing 23 / 95

- An idea from 1994¹ : fit a low rank convex hull with minimum volume.
- Theoretical result in 2015^2 : such hull is identifiable if the data points are well spreaded (the underlying θ is not too small)
 - \implies **volume** regularized NMF.

 $^{^1\}mathrm{Craig},$ Minimum-volume transforms for remotely sensed data. IEEE Trans. Geosci. Remote Sensing

 $^{^{2}}$ Lin, et. al, Identifiability of the simplex volume minimization criterion for blind hyperspectral unmixing: The no-pure-pixel case. IEEE Trans Geosci. Remote Sensing 24 / 95

- An idea from 1994¹ : fit a low rank convex hull with minimum volume.
- Theoretical result in 2015² : such hull is identifiable if the data points are well spreaded (the underlying θ is not too small)
 ⇒ volume regularized NMF.
- ullet Volume is related to determinant $\implies \det \mathbf{W}$ regularization

25 / 95

 $^{^1\}mathrm{Craig},$ Minimum-volume transforms for remotely sensed data. IEEE Trans. Geosci. Remote Sensing

²Lin, et. al, Identifiability of the simplex volume minimization criterion for blind hyperspectral unmixing: The no-pure-pixel case. IEEE Trans Geosci. Remote Sensing

- An idea from 1994¹ : fit a low rank convex hull with minimum volume.
- Theoretical result in 2015² : such hull is identifiable if the data points are well spreaded (the underlying θ is not too small)
 ⇒ volume regularized NMF.
- $\bullet~$ Volume is related to determinant $\implies \det \mathbf{W}$ regularization
- det W only works for sqaure $W \implies \mathsf{consider} \det W^\top W$ or $\log \det W^\top W$

26 / 95

 $^{^1\}mathrm{Craig},$ Minimum-volume transforms for remotely sensed data. IEEE Trans. Geosci. Remote Sensing

²Lin, et. al, Identifiability of the simplex volume minimization criterion for blind hyperspectral unmixing: The no-pure-pixel case. IEEE Trans Geosci. Remote Sensing

- An idea from 1994¹ : fit a low rank convex hull with minimum volume.
- Theoretical result in 2015² : such hull is identifiable if the data points are well spreaded (the underlying θ is not too small)
 ⇒ volume regularized NMF.
- $\bullet~$ Volume is related to determinant $\implies \det \mathbf{W}$ regularization
- det W only works for sqaure $W \implies \mathsf{consider} \det W^\top W$ or $\log \det W^\top W$

$$\begin{split} \det \mathsf{NMF} &: \min_{\substack{\mathbf{W} \geq 0 \\ \mathbf{H} \geq 0 \\ \mathbf{1}_r^\top \mathbf{H} \leq \mathbf{1}_n}} \frac{\frac{1}{2} \|\mathbf{X} - \mathbf{WH}\|_F^2}{\mathsf{data fitting term } \mathcal{F}} + \frac{\lambda}{2} \underbrace{\det(\mathbf{W}^\top \mathbf{W})}_{\mathsf{volume regularizer } \mathcal{G}} \\ \mathsf{logdetNMF} &: \min_{\substack{\mathbf{W} \geq 0 \\ \mathbf{H} \geq 0 \\ \mathbf{1}_r^\top \mathbf{H} \leq \mathbf{1}_n}} \frac{\frac{1}{2} \|\mathbf{X} - \mathbf{WH}\|_F^2}{\mathsf{data fitting term } \mathcal{F}} + \frac{\lambda}{2} \underbrace{\log \det(\mathbf{W}^\top \mathbf{W} + \delta \mathbf{I})}_{\mathsf{volume regularizer } \mathcal{G}} . \end{split}$$

 $^1\mathrm{Craig},$ Minimum-volume transforms for remotely sensed data. IEEE Trans. Geosci. Remote Sensing

 2 Lin, et. al, Identifiability of the simplex volume minimization criterion for blind hyperspectral unmixing: The no-pure-pixel case. IEEE Trans Geosci. Remote Sensing 27 / 95

Log-determinant Non-negative Matrix Factorization

Given $\mathbf{X} \in \mathbb{R}^{m imes n}$, find matrices $\mathbf{W} \in \mathbb{R}^{m imes r}_+, \mathbf{H} \in \mathbb{R}^{r imes n}_+$ by solving

$$\min_{\substack{\mathbf{W} \ge 0\\ \mathbf{H} \ge 0\\ \mathbf{I}_r^\top \mathbf{H} \le \mathbf{1}_n}} \underbrace{\frac{1}{2} \|\mathbf{X} - \mathbf{W}\mathbf{H}\|_F^2}_{\text{data fitting term}} + \frac{\lambda}{2} \underbrace{\log \det(\mathbf{W}^\top \mathbf{W} + \delta \mathbf{I}_r)}_{\text{volume regularizer}}.$$

• $\lambda > 0$: tuning parameters (regularization parameter).

• ($r \in \mathbb{N}_+$ assumed known, \mathbf{W}, \mathbf{H} assumed full rank).

Log-determinant Non-negative Matrix Factorization

Given $\mathbf{X} \in \mathbb{R}^{m imes n}$, find matrices $\mathbf{W} \in \mathbb{R}^{m imes r}_+, \mathbf{H} \in \mathbb{R}^{r imes n}_+$ by solving

$$\min_{\substack{\mathbf{W} \ge 0\\ \mathbf{H} \ge 0\\ \mathbf{I}_r^\top \mathbf{H} \le \mathbf{1}_n}} \underbrace{\frac{1}{2} \|\mathbf{X} - \mathbf{W}\mathbf{H}\|_F^2}_{\text{data fitting term}} + \frac{\lambda}{2} \underbrace{\log \det(\mathbf{W}^\top \mathbf{W} + \delta \mathbf{I}_r)}_{\text{volume regularizer}}.$$

- $\lambda > 0$: tuning parameters (regularization parameter).
- δ : fix small positive constant (e.g. 1).
- Why $\delta \mathbf{I}_r$: to bound $\log \det$ (otherwise $\lim_{\|\mathbf{W}\| \to 0} \log \det \mathbf{W}^\top \mathbf{W} \to -\infty$).
- ($r \in \mathbb{N}_+$ assumed known, \mathbf{W}, \mathbf{H} assumed full rank).

Two Block Coordinate Descent solution framework

The optimization problem : given $\mathbf{X} \in \mathbb{R}^{m \times n}$ and $r \geq 1$

$$\min_{\substack{\mathbf{W} \ge 0\\ \mathbf{H} \ge 0\\ \mathbf{I}_r^\top \mathbf{H} \le \mathbf{1}_n}} \|\mathbf{X} - \mathbf{W}\mathbf{H}\|_F^2 + \lambda \log \det(\mathbf{W}^\top \mathbf{W} + \delta \mathbf{I}_r),$$

1: INPUT : $\mathbf{X} \in \mathbb{R}^{m \times n}$, $r \in \mathbb{N}_+$ and $\lambda \ge 0$

2:
$$\mathsf{OUTPUT}: \mathbf{W} \in \mathbb{R}^{m imes r}_+$$
 and $\mathbf{H} \in \mathbb{R}^{r imes n}_+$

- 3: INITIALIZATION : $\mathbf{W} \in \mathbb{R}^{m imes r}_+$ and $\mathbf{H} \in \mathbb{R}^{r imes n}_+$
- 4: for k = 1 to itermax do

5:
$$\mathbf{W} \leftarrow \underset{\mathbf{W} \ge 0}{\operatorname{arg\,min}} \|\mathbf{X} - \mathbf{W}\mathbf{H}\|_F^2 + \lambda \log \det(\mathbf{W}^\top \mathbf{W} + \delta \mathbf{I}_r).$$

- 6: $\mathbf{H} \leftarrow \operatorname*{arg\,min}_{\mathbf{H} \ge 0, \mathbf{1}_r^\top \mathbf{H} \le \mathbf{1}_n} \|\mathbf{X} \mathbf{W}\mathbf{H}\|_F^2 + \lambda \log \det(\mathbf{W}^\top \mathbf{W} + \delta \mathbf{I}_r).$
- 7: end for

From now on, line 1-3 will be skipped (for space) Subproblems lines 5-6 can be solve by projected gradient.

On solving ${\bf H}$ and ${\bf W}$

To solve

$$\mathbf{H} \leftarrow \underset{\mathbf{H} \geq 0}{\operatorname{arg\,min}} \|\mathbf{X} - \mathbf{W}\mathbf{H}\|_{F}^{2} + \underbrace{\lambda \log \det(\mathbf{W}^{\top}\mathbf{W} + \delta \mathbf{I}_{r})}_{\text{independent of }\mathbf{H}},$$

the FGM (Fast gradient method on constrainted least sqaure on unit simplex) from N. Gillis^{\dagger} will be used :

1: for k = 1 to itermax do

2:
$$\mathbf{W} \leftarrow \underset{\mathbf{W} \ge 0}{\operatorname{arg\,min}} \|\mathbf{X} - \mathbf{W}\mathbf{H}\|_{F}^{2} + \lambda \log \det(\mathbf{W}^{\top}\mathbf{W} + \delta \mathbf{I}_{r}).$$

- 3: Update \mathbf{H} using FGM[†] with { $\mathbf{X}, \mathbf{W}, \mathbf{H}$ }.
- 4: end for

† N. Gillis, "Successive Nonnegative Projection Algorithm for Robust Nonnegative Blind Source Separation", SIAM J. on Imaging Sciences 7 (2), pp. 1420-1450, 2014.

On solving ${\bf H}$ and ${\bf W}$

So the key problem is to solve for W.

- Data fitting part $\|\mathbf{X} \mathbf{W}\mathbf{H}\|_F^2$ is easy to handle.
- The regularizer $\log \det(\mathbf{W}^{\top}\mathbf{W} + \delta \mathbf{I}_r)$ is problematic : non-convex, column-coupled, non-proximable.
- *Don't forget we can always solve this problem just by vanilla projected gradient, but such approach does not utilize the structure of the problem not good !
- 1: for k = 1 to itermax do
- 2: $\mathbf{W} \leftarrow \underset{\mathbf{W} \ge 0}{\operatorname{arg\,min}} \|\mathbf{X} \mathbf{W}\mathbf{H}\|_{F}^{2} + \lambda \log \det(\mathbf{W}^{\top}\mathbf{W} + \delta \mathbf{I}_{r}).$
- 3: Update \mathbf{H} using FGM[†] with { $\mathbf{X}, \mathbf{W}, \mathbf{H}$ }.
- 4: end for

† N. Gillis, "Successive Nonnegative Projection Algorithm for Robust Nonnegative Blind Source Separation", SIAM J. on Imaging Sciences 7 (2), pp. 1420-1450, 2014.

Previous lines of attack on $\log \det$

Previous lines of attack (fails or no result) :

- Proximal operator on $+\log \det \mathbf{W}^{\top}\mathbf{W}$
- Hadamard's inequality : $|\det(\mathbf{A})| < \prod_i ||a_i||_2^2$ (Upper bound of volume spanned by $\mathbf{A} = [a_1 a_2 ...]$)
- exp-trace-log equation : $det(\mathbf{I}_r + \mathbf{A}) = exp tr log(\mathbf{I}_r + \mathbf{A})$
- Approximating the determinant 1. Diagonal Approximations
- Approximating the determinant 2. Eigenspectrum approximations

$$det(\mathbf{I}_r + \delta \mathbf{A}) = det(\mathbf{I}_r + \delta P^{-1}JP)$$

= $det(P^{-1}(\mathbf{I} + \delta J)P)$
= $det(P^{-1}P(\mathbf{I} + \delta J))$
= $\prod_i (1 + \delta J_{ii})$
= $1 + \delta \sum_i J_{ii} + \delta^2 \sum_{i,j,j \neq i} J_{ii}J_{jj} + .$

• A bound from telecom research : $\log \det(\mathbf{W}^{\top}\mathbf{W} + \delta \mathbf{I}_r) \leq \operatorname{tr}(D^{\mathsf{Taylor}}\mathbf{A}^{\top}\mathbf{A}) - \log \det(\mathbf{D}^{\mathsf{Taylor}}) - r$ where $\mathbf{D}^{\mathsf{Taylor}} = (\mathbf{W}_{-1}^{\top}\mathbf{W}_{-1} + \delta \mathbf{I}_r)^{-1}$, this is infact the first order Taylor convex upper bound of the function $\log \det(X)$. Reference includes : S.Christensen et al., "Weighted Sum-Rate Maximization using Weighted MMSE for MIMO-BC Beamforming Design", IEEE Trans. Wireless Com., pp. 4792-4799, vol. 7, issue 12, 2008

The key iequality : logdet-trace inequality

Given a positive definite matrix $\mathbf{A} \in \mathbb{R}^{r \times r}$, we have

$$\operatorname{tr}(\mathbf{I}_r - \mathbf{A}^{-1}) \le \log \det \mathbf{A} \le \operatorname{tr}(\mathbf{A} - \mathbf{I}_r)$$

Sowe now have an upper bound : put $\mathbf{A} = \mathbf{W}^{ op} \mathbf{W} + \delta \mathbf{I}_r$

 $\log \det(\mathbf{W}^t \mathbf{W} + \delta \mathbf{I}_r) \leq \operatorname{tr}(\mathbf{W}^\top \mathbf{W} + (\delta - 1)\mathbf{I}_r)$

 $\|\mathbf{X} - \mathbf{W}\mathbf{H}\|_F^2 + \lambda \log \det(\mathbf{W}^\top \mathbf{W} + \delta \mathbf{I}_r) \leq \|\mathbf{X} - \mathbf{W}\mathbf{H}\|_F^2 + \lambda \operatorname{tr}(\mathbf{W}^\top \mathbf{W} + (\delta - 1)\mathbf{I}_r)$

• $\log \det(\mathbf{W}^{\top}\mathbf{W} + \delta \mathbf{I}_r)$ is not convex w.r.t. \mathbf{W} but the trace is.

The key iequality : logdet-trace inequality

Given a positive definite matrix $\mathbf{A} \in \mathbb{R}^{r \times r}$, we have

$$\operatorname{tr}(\mathbf{I}_r - \mathbf{A}^{-1}) \le \log \det \mathbf{A} \le \operatorname{tr}(\mathbf{A} - \mathbf{I}_r)$$

Sowe now have an upper bound : put $\mathbf{A} = \mathbf{W}^{ op} \mathbf{W} + \delta \mathbf{I}_r$

$$\log \det(\mathbf{W}^t \mathbf{W} + \delta \mathbf{I}_r) \leq \operatorname{tr}(\mathbf{W}^\top \mathbf{W} + (\delta - 1)\mathbf{I}_r)$$

 $\|\mathbf{X} - \mathbf{W}\mathbf{H}\|_F^2 + \lambda \log \det(\mathbf{W}^\top \mathbf{W} + \delta \mathbf{I}_r) \leq \|\mathbf{X} - \mathbf{W}\mathbf{H}\|_F^2 + \lambda \operatorname{tr}(\mathbf{W}^\top \mathbf{W} + (\delta - 1)\mathbf{I}_r)$

- $\log \det(\mathbf{W}^{\top}\mathbf{W} + \delta \mathbf{I}_r)$ is not convex w.r.t. \mathbf{W} but the trace is.
- Algorithm that minimizes this upper bound :
- 1: for k = 1 to itermax do
- 2: $\mathbf{W} \leftarrow \underset{\mathbf{W} \ge 0}{\operatorname{arg\,min}} \|\mathbf{X} \mathbf{W}\mathbf{H}\|_F^2 + \lambda \operatorname{tr}(\mathbf{W}^\top \mathbf{W} + (\delta 1)\mathbf{I}_r).$
- 3: Update \mathbf{H} using FGM with $\mathbf{X}, \mathbf{W}, \mathbf{H}$.
- 4: end for

The key iequality : logdet-trace inequality

Given a positive definite matrix $\mathbf{A} \in \mathbb{R}^{r \times r}$, we have

$$\operatorname{tr}(\mathbf{I}_r - \mathbf{A}^{-1}) \le \log \det \mathbf{A} \le \operatorname{tr}(\mathbf{A} - \mathbf{I}_r)$$

Sowe now have an upper bound : put $\mathbf{A} = \mathbf{W}^{ op} \mathbf{W} + \delta \mathbf{I}_r$

$$\log \det(\mathbf{W}^t \mathbf{W} + \delta \mathbf{I}_r) \leq \operatorname{tr}(\mathbf{W}^\top \mathbf{W} + (\delta - 1)\mathbf{I}_r)$$

 $\|\mathbf{X} - \mathbf{W}\mathbf{H}\|_F^2 + \lambda \log \det(\mathbf{W}^\top \mathbf{W} + \delta \mathbf{I}_r) \leq \|\mathbf{X} - \mathbf{W}\mathbf{H}\|_F^2 + \lambda \operatorname{tr}(\mathbf{W}^\top \mathbf{W} + (\delta - 1)\mathbf{I}_r)$

• $\log \det(\mathbf{W}^{\top}\mathbf{W} + \delta \mathbf{I}_r)$ is not convex w.r.t. \mathbf{W} but the trace is.

• Algorithm that minimizes this upper bound :

- 1: for k = 1 to itermax do
- 2: $\mathbf{W} \leftarrow \underset{\mathbf{W} \ge 0}{\operatorname{arg\,min}} \|\mathbf{X} \mathbf{W}\mathbf{H}\|_F^2 + \lambda \operatorname{tr}(\mathbf{W}^\top \mathbf{W} + (\delta 1)\mathbf{I}_r).$
- 3: Update \mathbf{H} using FGM with $\mathbf{X}, \mathbf{W}, \mathbf{H}$.

4: end for

• Don't stop here, it can be better !!
• Given a positive definite matrix $\mathbf{A} \in \mathbb{R}^{r imes r}$, we have

 $\log \det \mathbf{A} \leq \operatorname{tr}(\mathbf{A} - \mathbf{I}_r).$

• Given a positive definite matrix $\mathbf{A} \in \mathbb{R}^{r imes r}$, we have

 $\log \det \mathbf{A} \leq \operatorname{tr}(\mathbf{A} - \mathbf{I}_r).$

• Let μ denotes eigenvalues, we have

$$\det \mathbf{A} = \prod_i \mu_i$$
 and $\operatorname{tr} \mathbf{A} = \sum_i \mu_i.$

$$\log \det \mathbf{A} \leq \operatorname{tr}(\mathbf{A} - \mathbf{I}_r) \iff \sum_i \log \mu_i \leq \sum_i (\mu_i - 1).$$

• Given a positive definite matrix $\mathbf{A} \in \mathbb{R}^{r imes r}$, we have

 $\log \det \mathbf{A} \leq \operatorname{tr}(\mathbf{A} - \mathbf{I}_r).$

• Let μ denotes eigenvalues, we have

det
$$\mathbf{A} = \prod_{i} \mu_{i}$$
 and $\operatorname{tr} \mathbf{A} = \sum_{i} \mu_{i}$.

$$\log \det \mathbf{A} \leq \operatorname{tr}(\mathbf{A} - \mathbf{I}_r) \iff \sum_i \log \mu_i \leq \sum_i (\mu_i - 1).$$

• $\log \mu_i$ means matrix **A** has to be positive definite $(\mu_i > 0 \forall i)$, which is satisfied for $\mathbf{A} = \mathbf{W}^\top \mathbf{W} + \delta \mathbf{I}_r$.

• Given a positive definite matrix $\mathbf{A} \in \mathbb{R}^{r imes r}$, we have

 $\log \det \mathbf{A} \leq \operatorname{tr}(\mathbf{A} - \mathbf{I}_r).$

• Let μ denotes eigenvalues, we have

$$\det \mathbf{A} = \prod_i \mu_i$$
 and $\operatorname{tr} \mathbf{A} = \sum_i \mu_i$.

$$\log \det \mathbf{A} \leq \operatorname{tr}(\mathbf{A} - \mathbf{I}_r) \iff \sum_i \log \mu_i \leq \sum_i (\mu_i - 1).$$

- $\log \mu_i$ means matrix **A** has to be positive definite $(\mu_i > 0 \forall i)$, which is satisfied for $\mathbf{A} = \mathbf{W}^\top \mathbf{W} + \delta \mathbf{I}_r$.
- $\sum_i \log \mu_i \leq \sum_i (\mu_i 1) \Leftrightarrow \log \mu_i \leq \mu_i 1 \,\forall i$. We can focuse on the inequality $\log \mu_i \leq \mu_i 1$ with $\mu_i \geq 0$

On $\log x \le x - 1$, $x \ge 0$

- $\log x$ is concave.
- x 1 is the first order Taylor approximation of $\log x$ at x = 1.
- x-1 is the only **convex-tight** upper bound of $\log x$.[†]
- Tight : x 1 touch $\log x$ at the point x = 1.

On $\log x \le x - 1$, $x \ge 0$

- $\log x$ is concave.
- x 1 is the first order Taylor approximation of $\log x$ at x = 1.
- x-1 is the only **convex-tight** upper bound of $\log x$.[†]
- Tight : x 1 touch $\log x$ at the point x = 1.
- Generalize to point x_0 : $\log x \le g(x|x_0) = a_1(x_0)x + a_0(x_0)$ is

$$\log x \le \frac{1}{x_0}x + \log x_0 - 1.$$



 \dagger Higher order Taylor approximation of $\log x$ is tight, more accurate but not convex.

A parametric trace upper bound for $\log \det \mathbf{A}$

$$\log \det \mathbf{A} = \sum \log \mu_i$$

$$\leq \sum \frac{1}{\mu_i^-} \mu_i + \log \mu_i^- - 1$$

$$\leq \sum \frac{1}{\mu_{\min}^-} \mu_i + \log \mu_i^- - 1$$

$$= \operatorname{tr}(\mathbf{D}^1 \mathbf{A} + \mathbf{D}^0)$$

 $\mathbf{D}^1 = \frac{1}{\mu_{\min}^-} \mathbf{I}_r, \ \mathbf{D}^0 = \mathsf{Diag}(\log \mu_i^- - 1), \ \mu_i^- \text{ is } \mu_i \text{ of the previous step}$

Put $\mathbf{A} = \mathbf{W}^{\top} \mathbf{W} + \delta \mathbf{I}_r$, we have

$$\begin{split} \log \det(\mathbf{W}^{\top}\mathbf{W} + \delta \mathbf{I}_{r}) &\leq \operatorname{tr}(\mathbf{D}^{1}\mathbf{W}^{\top}\mathbf{W} + \delta \mathbf{D}^{1} + \mathbf{D}^{0}) \\ (\text{ignore constants}) &= \operatorname{tr}\mathbf{D}^{1}\mathbf{W}^{\top}\mathbf{W} \end{split}$$

Consider the matrix A has eigen-decomposition as $\mathbf{A} = \mathbf{V}\Lambda\mathbf{V}^T$. Let the weighting $\frac{1}{\mu_i^-}$ be a_i , then $\sum_i \frac{1}{\mu_i^-} \mu_i = \sum_i a_i \mu_i$ and



 $\neq \operatorname{tr} \mathbf{D}^{a} \mathbf{A}$

:

The original function $F(\mathbf{W}) = \log \det(\mathbf{W}^{\top}\mathbf{W} + \delta \mathbf{I}_r)$ is upper bounded by

Eigen bound (B_1) : tr $\mathbf{D}^{\mathsf{T}}\mathbf{W}^{\mathsf{T}}\mathbf{W}$ + constants. Taylor bound (B_2) : tr $\mathbf{D}^{\mathsf{Taylor}}\mathbf{W}^{\mathsf{T}}\mathbf{W}$ + constants.

- Constants \mathbf{D}^1 , \mathbf{D}^0 are defined as before, and constant $\mathbf{D}^{\mathsf{Taylor}} = (\mathbf{W}_{-1}^{\top}\mathbf{W}_{-1} + \delta \mathbf{I}_r)^{-1}.$
- Both bounds are trace functional with an relaxation gap :
 - (B_1) has eigen gap $\mu_i \ge \mu_{\min}$
 - (B_2) has convexification gap
 - ▶ D¹ is diagonal but D^{Taylor} is not (it is dense) ⇒ column-wise decomposition is possible

Algorithm 1 Successive Trace Approximation

- 1: INPUT: $\mathbf{X} \in \mathbb{R}^{m \times n}_+$, $r \in \mathbb{N}_+$, $\lambda > 0$, $\delta > 0$.
- 2: OUTPUT: $\mathbf{W} \in \mathbb{R}^{m \times r}_+$ and $\mathbf{H} \in \mathbb{R}^{r \times n}_+$.
- 3: INITIALIZATION : $\mathbf{W}\in\mathbb{R}_+^{m imes r}$, $\mathbf{H}\in\mathbb{R}_+^{r imes n}$, $\mathbf{D}^1=\mathbf{I}_r$
- 4: for K = 1 to itermax do
- 5: for k = 1 to itermax do
- 6: $\mathbf{W} \leftarrow \operatorname*{arg\,min}_{\mathbf{W} \ge 0} \|\mathbf{X} \mathbf{W}\mathbf{H}\|_F^2 + \lambda \operatorname{tr} \mathbf{D}^1 \mathbf{W}^\top \mathbf{W}.$
- 7: Update \mathbf{H} using FGM with $\mathbf{X}, \mathbf{W}, \mathbf{H}$.
- 8: end for
- 9: $\mu_i \leftarrow \mathsf{svd}(\mathbf{W}^\top \mathbf{W} + \delta \mathbf{I}_r), \ \mathbf{D}^1 = \mathsf{Diag}(\mu_{\min}^{-1})$
- 10: end for

Small summary : a model relaxation



In one sentence : Eigenval-wise convex relaxation of a non-convex problem using logdet – trace ineqaulity.

Consider one vector w_i while fixing all other things :

 $\|\mathbf{X} - \mathbf{W}\mathbf{H}\|_F^2 + \lambda \operatorname{tr}(\mathbf{D}^1\mathbf{W}^\top\mathbf{W} + \mathbf{D}^0)$

$$= \|\mathbf{X} - \sum_{i} w_{i}h_{i}\|_{F}^{2} + \lambda \sum_{i} \left(D_{ii}^{1} \|w_{i}\|_{2}^{2} + \mathbf{D}_{ii}^{0}\right)$$

Consider one vector w_i while fixing all other things :

 $\|\mathbf{X} - \mathbf{W}\mathbf{H}\|_F^2 + \lambda \operatorname{tr}(\mathbf{D}^1\mathbf{W}^\top\mathbf{W} + \mathbf{D}^0)$

$$= \|\mathbf{X} - \sum_{i} w_{i}h_{i}\|_{F}^{2} + \lambda \sum_{i} \left(D_{ii}^{1} \|w_{i}\|_{2}^{2} + \mathbf{D}_{ii}^{0}\right)$$

$$= \|\underbrace{\left(X - \sum_{j \neq i} w_{j}h_{j}\right) - w_{i}h_{i}\|_{F}^{2} + \lambda \left(\mathbf{D}_{ii}^{1} \|w_{i}\|_{2}^{2} + \sum_{j \neq i} \left(\mathbf{D}_{jj}^{1} \|w_{j}\|_{2}^{2} + \mathbf{D}_{ii}^{0}\right)\right)}_{c}$$

Consider one vector w_i while fixing all other things :

 $\|\mathbf{X} - \mathbf{W}\mathbf{H}\|_F^2 + \lambda \operatorname{tr}(\mathbf{D}^1\mathbf{W}^\top\mathbf{W} + \mathbf{D}^0)$

$$= \|\mathbf{X} - \sum_{i} w_{i}h_{i}\|_{F}^{2} + \lambda \sum_{i} \left(D_{ii}^{1} \|w_{i}\|_{2}^{2} + \mathbf{D}_{ii}^{0}\right)$$

$$= \|(X - \sum_{j \neq i} w_{j}h_{j}) - w_{i}h_{i}\|_{F}^{2} + \lambda \left(\mathbf{D}_{ii}^{1} \|w_{i}\|_{2}^{2} + \sum_{j \neq i} \left(\mathbf{D}_{jj}^{1} \|w_{j}\|_{2}^{2} + \mathbf{D}_{ii}^{0}\right)\right)$$

$$= \|\mathbf{X}_{i} - w_{i}h_{i}\|_{F}^{2} + \lambda \mathbf{D}_{ii}^{1} \|w_{i}\|_{2}^{2} + c$$

Consider one vector w_i while fixing all other things :

 $\|\mathbf{X} - \mathbf{W}\mathbf{H}\|_F^2 + \lambda \operatorname{tr}(\mathbf{D}^1\mathbf{W}^\top\mathbf{W} + \mathbf{D}^0)$

$$= \|\mathbf{X} - \sum_{i} w_{i}h_{i}\|_{F}^{2} + \lambda \sum_{i} \left(D_{ii}^{1} \|w_{i}\|_{2}^{2} + \mathbf{D}_{ii}^{0}\right)$$

$$= \|\underbrace{\left(X - \sum_{j \neq i} w_{j}h_{j}\right) - w_{i}h_{i}\|_{F}^{2} + \lambda \left(\mathbf{D}_{ii}^{1} \|w_{i}\|_{2}^{2} + \sum_{j \neq i} \left(\mathbf{D}_{jj}^{1} \|w_{j}\|_{2}^{2} + \mathbf{D}_{ii}^{0}\right)\right)}{\mathbf{x}_{i}}$$

$$= \|\mathbf{X}_{i} - w_{i}h_{i}\|_{F}^{2} + \lambda \mathbf{D}_{ii}^{1} \|w_{i}\|_{2}^{2} + c$$

$$\leq \|\mathbf{X}_{i} - w_{i}h_{i}\|_{F}^{2} + \lambda \mathbf{D}_{ii}^{1} \|w_{i}\|_{2}^{2} + \frac{\gamma}{2} \|w_{i} - w_{i}^{-}\|_{2}^{2} + c.$$

Ignoring constants, we have a constrainted regularized QP

$$\min_{w_i \ge 0} \|\mathbf{X}_i - w_i h_i\|_F^2 + \lambda \mathbf{D}_{ii}^1 \|w_i\|_2^2 + \frac{\gamma}{2} \|w_i - w_i^-\|_2^2.$$

where w_i^- is the previous iterate of w_i , $\gamma > 0$ is a (small) constant. The proximal term $||w_i - w_i^-||_2^2$ penalizes w for leaving w^- too far. 51/95

$$\min_{w_i \ge 0} \|\mathbf{X}_i - w_i h_i\|_F^2 + \lambda \mathbf{D}_{ii}^1 \|w_i\|_2^2 + \frac{\gamma}{2} \|w_i - w_i^-\|_2^2.$$

Introducing the proximal term :

- turns the QP problem strongly convex.
- gurantees $\|h\|_2^2 + \lambda \mathbf{D}_{ii}^1 + \gamma > 0$: Lemma Given $\mathbf{X}, h, z, \alpha, \beta$, the optimal solution of

$$\min_{w \ge 0} \|\mathbf{X} - wh\|_F^2 + \alpha \|w\|_2^2 + \beta \|w - z\|_2^2$$

is

$$w = \frac{\left[\mathbf{X}h^T + \beta z\right]_+}{\|h\|_2^2 + \alpha + \beta}.$$

Proof: skiped.

Note. This is also related to solving the problem using Newton iteration.

The STA algorithm with decomposed $f(w_i)$

- 1: INPUT: $\mathbf{X} \in \mathbb{R}^{m \times n}_+$, $r \in \mathbb{N}_+$, $\lambda > 0$ and $\delta > 0$
- 2: OUTPUT: $\mathbf{W} \in \mathbb{R}^{m imes r}_+$ and $\mathbf{H} \in \mathbb{R}^{r imes n}_+$
- 3: INITIALIZATION : $\mathbf{W} \in \mathbb{R}^{m imes r}_+$, $\mathbf{H} \in \mathbb{R}^{r imes n}_+$ and $\mathbf{D}^1 = I_r$, $\gamma = 10^{-6}$
- 4: for K = 1 to itermax do
- 5: for k = 1 to itermax do
- 6: for i = 1 to r do
- 7: $w_i = \operatorname*{arg\,min}_{w_i \ge 0} f(w_i) = \|\mathbf{X}_i w_i h_i\|_F^2 + \lambda \mathbf{D}_{ii}^1 \|w_i\|_2^2 + \frac{\gamma}{2} \|w_i w_i^-\|_2^2$
- 8: Update \mathbf{H} by FGM with $\mathbf{X}, \mathbf{W}, \mathbf{H}$
- 9: end for
- 10: end for
- 11: $\mu_i \leftarrow \mathsf{svd}(\mathbf{W}^\top \mathbf{W} + \delta I) \text{ and } \mathbf{D}^1 = \mathsf{Diag}(\mu_{\min}^{-1})$

12: end for

Apply close form solution to $f(w_i)$

1: INPUT: $\mathbf{X} \in \mathbb{R}^{m \times n}_+$, $r \in \mathbb{N}_+$, $\lambda > 0$ and $\delta > 0$ 2: OUTPUT: $\mathbf{W} \in \mathbb{R}^{m \times r}_+$ and $\mathbf{H} \in \mathbb{R}^{r \times n}_+$ 3: INITIALIZATION : $\mathbf{W} \in \mathbb{R}^{m imes r}_+$, $\mathbf{H} \in \mathbb{R}^{r imes n}_+$ and $\mathbf{D}^1 = \mathbf{I}_r$, $\gamma = 10^{-6}$ 4 for K = 1 to itermax do **for** k = 1 to itermax **do** 5: for i = 1 to r do 6: $w_i = \frac{\left[X_i h_i^T + \gamma w_i^-\right]_+}{\|h_i\|_2^2 + \lambda \mathbf{D}_{ii}^1 + \frac{\gamma}{2}} \quad \text{where} \quad \mathbf{X}_i = \mathbf{X} - \sum_{j \neq i} w_j h_j$ 7: Update H by FGM with X, W, H8: end for 9: end for 10: $\mu_i \leftarrow \mathsf{svd}(\mathbf{W}^\top \mathbf{W} + \delta \mathbf{I}_r) \text{ and } \mathbf{D}^1 = \mathsf{Diag}(\mu_{\min}^{-1})$ 11: 12: end for

Move the update of ${\bf H}$ outside the loop to reduce computation burden

- 1: INPUT: $\mathbf{X} \in \mathbb{R}^{m \times n}_+$, $r \in \mathbb{N}_+$, $\lambda > 0$ and $\delta > 0$
- 2: OUTPUT: $\mathbf{W} \in \mathbb{R}^{m \times r}_+$ and $\mathbf{H} \in \mathbb{R}^{r \times n}_+$
- 3: INITIALIZATION : $\mathbf{W} \in \mathbb{R}^{m imes r}_+$, $\mathbf{H} \in \mathbb{R}^{r imes n}_+$ and $\mathbf{D}^1 = \mathbf{I}_r$, $\gamma = 10^{-6}$
- 4: for K = 1 to itermax do
- 5: for k = 1 to itermax do

6: **for**
$$i = 1$$
 to r **do**
7: $w_i = \frac{\left[\mathbf{X}_i h_i^T + \gamma w_i^-\right]_+}{\|h_i\|_2^2 + \lambda \mathbf{D}_{ii}^1 + \frac{\gamma}{2}}$ where $\mathbf{X}_i = \mathbf{X} - \sum_{j \neq i} w_j h_j$

8: end for

- 9: Update ${\bf H}$ by FGM with ${\bf X}, {\bf W}, {\bf H}$
- 10: end for

11:
$$\mu_i \leftarrow \mathsf{svd}(\mathbf{W}^\top \mathbf{W} + \delta \mathbf{I}_r)$$
 and $\mathbf{D}^1 = \mathsf{Diag}(\mu_{\min}^{-1})$

12: end for

Move the update of D^1 inside the loop ($\mathbf{W}^{\top}\mathbf{W}$ is *r*-by-*r*, small !)

1: INPUT:
$$\mathbf{X} \in \mathbb{R}^{m \times n}_+$$
, $r \in \mathbb{N}_+$, $\lambda > 0$ and $\delta > 0$
2: OUTPUT: $\mathbf{W} \in \mathbb{R}^{m \times r}_+$ and $\mathbf{H} \in \mathbb{R}^{r \times n}_+$
3: INITIALIZATION : $\mathbf{W} \in \mathbb{R}^{m \times r}_+$, $\mathbf{H} \in \mathbb{R}^{r \times n}_+$ and $\mathbf{D}^1 = \mathbf{I}_r$, $\gamma = 10^{-6}$
4: for $K = 1$ to itermax do
5: for $k = 1$ to itermax do
6: for $i = 1$ to r do
7: $w_i = \frac{[\mathbf{X}_i h_i^T + \gamma w_i^-]_+}{\|h_i\|_2^2 + \lambda \mathbf{D}_{ii}^1 + \frac{\gamma}{2}}$ where $\mathbf{X}_i = \mathbf{X} - \sum_{j \neq i} w_j h_j$
8: $\mu_i \leftarrow \operatorname{svd}(\mathbf{W}^\top \mathbf{W} + \delta \mathbf{I}_r)$ and $\mathbf{D}^1 = \operatorname{Diag}(\mu_{\min}^{-1})$
9: end for
10: Update H by FGM with $\mathbf{X}, \mathbf{W}, \mathbf{H}$
11: end for
12: end for

Now consider line 7, it can be show that it can be further optimized³.

1: INPUT:
$$\mathbf{X} \in \mathbb{R}^{m \times n}_+$$
, $r \in \mathbb{N}_+$, $\lambda > 0$ and $\delta > 0$

- 2: OUTPUT: $\mathbf{W} \in \mathbb{R}^{m \times r}_+$ and $\mathbf{H} \in \mathbb{R}^{r \times n}_+$
- 3: INITIALIZATION : $\mathbf{W} \in \mathbb{R}^{m imes r}_+$, $\mathbf{H} \in \mathbb{R}^{r imes n}_+$ and $\mathbf{D}^1 = \mathbf{I}_r$, $\gamma = 10^{-6}$
- 4: for K = 1 to itermax do
- 5: for k = 1 to itermax do

6: for
$$i = 1$$
 to r do
7: $w_i = \frac{\left[\mathbf{X}_i h^T + \gamma w_i^-\right]_+}{\|h\|_2^2 + \lambda \mathbf{D}_{ii}^1 + \frac{\gamma}{2}}$ where $\mathbf{X}_i = \mathbf{X} - \sum_{j \neq i} w_j h_j$

8:
$$\mu_i \leftarrow \mathsf{svd}(\mathbf{W}^\top \mathbf{W} + \delta \mathbf{I}_r) \text{ and } \mathbf{D}^1 = \mathsf{Diag}(\mu_{\min}^{-1})$$

- 9: end for
- 10: Update \mathbf{H} by FGM with $\mathbf{X}, \mathbf{W}, \mathbf{H}$
- 11: end for
- 12: end for

 $^{^3}$ N. Gillis and F. Glineur, "Accelerated Multiplicative Updates and Hierarchical ALS Algorithms for Nonnegative Matrix Factorization", Neural Computation 24 (4), 2012 $57\/95$

Line 7

$$w_i = \frac{\left[\mathbf{X}_i h_i^T + \gamma w_i^-\right]_+}{\|h_i\|_2^2 + \lambda \mathbf{D}_{ii}^1 + \frac{\gamma}{2}} \quad \text{where} \quad \mathbf{X}_i = \mathbf{X} - \sum_{j \neq i} w_j h_j.$$



Line 7

$$w_i = \frac{\left[\mathbf{X}_i h_i^T + \gamma w_i^-\right]_+}{\|h_i\|_2^2 + \lambda \mathbf{D}_{ii}^1 + \frac{\gamma}{2}} \quad \text{where} \quad \mathbf{X}_i = \mathbf{X} - \sum_{j \neq i} w_j h_j.$$

• Consider w_2 .





Line 7

$$w_i = \frac{\left[\mathbf{X}_i h_i^T + \gamma w_i^-\right]_+}{\|h_i\|_2^2 + \lambda \mathbf{D}_{ii}^1 + \frac{\gamma}{2}} \quad \text{where} \quad \mathbf{X}_i = \mathbf{X} - \sum_{j \neq i} w_j h_j.$$



Line 7

$$w_i = \frac{\left[\mathbf{X}_i h_i^T + \gamma w_i^-\right]_+}{\|h_i\|_2^2 + \lambda \mathbf{D}_{ii}^1 + \frac{\gamma}{2}}$$

where
$$\mathbf{X}_i = \mathbf{X} - \sum_{j
eq i} w_j h_j.$$

• The term $\mathbf{X}_i \mathbf{h}_i^T$



Line 7

$$w_i = \frac{\left[\mathbf{X}_i h_i^T + \gamma w_i^-\right]_+}{\|h_i\|_2^2 + \lambda \mathbf{D}_{ii}^1 + \frac{\gamma}{2}}$$

where
$$\mathbf{X}_i = \mathbf{X} - \sum_{j \neq i} w_j h_j.$$



Line 7

$$w_i = \frac{\left[\mathbf{X}_i h_i^T + \gamma w_i^-\right]_+}{\|h_i\|_2^2 + \lambda \mathbf{D}_{ii}^1 + \frac{\gamma}{2}} \quad \text{where} \quad \mathbf{X}_i = \mathbf{X} - \sum_{j \neq i} w_j h_j.$$



Line 7

$$w_i = \frac{\left[\mathbf{X}_i h_i^T + \gamma w_i^-\right]_+}{\|h_i\|_2^2 + \lambda \mathbf{D}_{ii}^1 + \frac{\gamma}{2}} \quad \text{where} \quad \mathbf{X}_i = \mathbf{X} - \sum_{j \neq i} w_j h_j.$$

is equivalent to

$$w_{i} = \frac{\left[P_{i} - \sum_{j=1}^{i-1} w_{j}Q_{ji} - \sum_{j=i+1}^{r} w_{j}^{-}Q_{ji} + \gamma w_{i}^{-}\right]_{+}}{Q_{ii} + \lambda \mathbf{D}_{ii}^{1} + \frac{\gamma}{2}}$$

where $P = \mathbf{X}\mathbf{H}^{\top}$, $P_i = P(:, i)$, $Q = \mathbf{H}\mathbf{H}^{\top}$ and $Q_{ji} = Q(j, i)$. P and Q can be pre-computed. Algorithm 2 STA

1: INPUT:
$$\mathbf{X} \in \mathbb{R}^{m \times n}_+$$
, $r \in \mathbb{N}_+$, $\lambda > 0$ and $\delta > 0$
2: OUTPUT: $\mathbf{W} \in \mathbb{R}^{m \times r}_+$ and $\mathbf{H} \in \mathbb{R}^{r \times n}_+$
3: INITIALIZATION : $\mathbf{W} \in \mathbb{R}^{m \times r}_+$, $\mathbf{H} \in \mathbb{R}^{r \times n}_+$ and $\mathbf{D}^1 = \mathbf{I}_r$, $\gamma = 10^{-6}$
4: for $K = 1$ to itermax do
5: for $k = 1$ to itermax do
6: $P = \mathbf{X}\mathbf{H}^\top$ and $Q = \mathbf{H}\mathbf{H}^\top$.
7: for $i = 1$ to r do
8: $w_i = \frac{\left[P_i - \sum_{j=1}^{i-1} w_j Q_{ji} - \sum_{j=i+1}^r w_j^- Q_{ji} + \gamma w_i^-\right]_+}{2}$

8:
$$w_{i} = \frac{1}{Q_{ii} + \lambda \mathbf{D}_{ii}^{1} + \frac{\gamma}{2}}$$

$$Q_{ii} + \lambda \mathbf{D}_{ii}^{1} + \frac{\gamma}{2}$$

$$Q_{ii} + \lambda \mathbf{D}_{ii}^{1} + \frac{\gamma}{2}$$

9:
$$\mu_i \leftarrow \mathsf{svd}(\mathbf{W}^\top \mathbf{W} + \delta \mathbf{I}_r) \text{ and } D^1 = \mathsf{Diag}(\mu_{\min}^{-1})$$

- 10: end for
- 11: Update \mathbf{H} by FGM with $\mathbf{X}, \mathbf{W}, \mathbf{H}$
- 12: end for

13: end for

In fact, line 7 - 10 can be run multiple times for "better" convergence. 65 / 95

Small summary

$$\|\mathbf{X} - \mathbf{W}\mathbf{H}\|_{F}^{2} \rightarrow \|\mathbf{X} - \mathbf{W}\mathbf{H}\|_{F}^{2} + \lambda \log \det(\mathbf{W}^{\top}\mathbf{W} + \delta\mathbf{I}_{r})$$

$$\|\mathbf{X} - \mathbf{W}\mathbf{H}\|_{F}^{2} + \lambda \operatorname{tr}(\mathbf{W}^{\top}\mathbf{W} + (\delta - 1)\mathbf{I}_{r})$$

$$\|\mathbf{X} - \mathbf{W}\mathbf{H}\|_{F}^{2} + \lambda \operatorname{tr}(\mathbf{D}^{1}\mathbf{W}^{\top}\mathbf{W} + \mathbf{D}^{0})$$

$$\sum_{i} \|\mathbf{X}_{i} - w_{i}h_{i}\|_{F}^{2} + \lambda \mathbf{D}_{ii}^{1}\|w_{i}\|_{2}^{2} + \frac{\gamma}{2}\|w_{i} - w_{i}^{-}\|_{2}^{2}$$

Synthetic data

- Ground truth : $\mathbf{W}_0 \in \mathbb{R}^{m imes r}_+$ and $\mathbf{H}_0 \in \mathbb{R}^{r imes n}_+$ (with $\mathbf{H}^T \mathbf{1} = \alpha$)
- m, n are sizes, r is rank, $\alpha \in (0 \ 1]$ tells how "well spread" the data are ($\alpha = 1$ means pure pixel)
- Form \mathbf{X}_0 as $\mathbf{X}_0 = \mathbf{W}_0 \mathbf{H}_0 \in \mathbb{R}^{m imes n}_+$
- Add noise $N \in \mathbb{R}^{m \times n}$ and $N \sim \mathcal{N}(0, R)$ as $\mathbf{X} = \mathbf{X}_0 + N$.
- As $N \in \mathbb{R}$ not \mathbb{R}_+ , corrupted data points may lie outside the convex hull.

Experiments and fancy figures - setup (2/3)

Real data (An example : Copperas Cove Texas Walmart)



Figure: RGB image of Copperas Cove Texas Walmart



Figure: Three spectral images of Copperas Cove Texas Walmart $\mathbf{X} \in \mathbb{R}^{94249 \times 162}_+$, or $\mathbf{X} \in \mathbb{R}^{307 \times 307 \times 162}_+$, pick r = 5, 6, 7.

Performance measurements. Algorithm produces \hat{W} and \hat{H} and $\hat{X}=\hat{W}\hat{H}$

- For simulation with known \mathbf{W}_0 , \mathbf{H}_0 , \mathbf{X}_0 :
 - Data fitting error : $\frac{\|\mathbf{X}_0 - \hat{\mathbf{X}}\|_F}{\|\mathbf{X}_0\|_F}$ Endmember fitting error : $\frac{\|\mathbf{W}_0 - \hat{\mathbf{W}}\|_F}{\|\mathbf{W}_0\|_F}$
 - Computational time
- For real data without knowing \mathbf{W}_0 , \mathbf{H}_0 , \mathbf{X}_0 :
 - Data fitting errror : $\frac{\|\mathbf{X} \hat{\mathbf{X}}\|_F}{\|\mathbf{X}\|_F}$
 - Volume of convex hull : $\log \det(\hat{\mathbf{W}}^{\top}\hat{\mathbf{W}} + \delta \mathbf{I}_r)$
 - Computational time

- What iterations of the algorithm looks like : the gif file
- Effect of fix lambda : EXAMPLE Rotate
- Effect of very big lambda : EXAMPLE Big lambda

Comparing the two logdet inequality - error



Comparing the two logdet inequality - fitting


Theoretical stuff — convergence

We have problems :

(P₀) minimizes $\|\mathbf{X} - \mathbf{W}\mathbf{H}\|_F^2 + \lambda \log \det(\mathbf{W}^\top \mathbf{W} + \delta \mathbf{I}_r),$ (P₁) minimizes $\|\mathbf{X} - \mathbf{W}\mathbf{H}\|_F^2 + \lambda \operatorname{tr}(\mathbf{D}^1 \mathbf{W}^\top \mathbf{W} + \mathbf{D}^0),$

under the constaints : $\mathbf{W} > 0, \mathbf{H} > 0, \mathbf{H}^{\top} \mathbf{1} < \mathbf{1}.$

Convergence properties :

- **(**) STA algorithm produces a stationary point for problem P_1 .
- **2** The solution of P_1 obtained by STA converges to the solution of P_0 .

Sonvergence rate of STA algorithm.

Idea : consider W, H and D¹, (and D⁰) as variables and treat STA is as an *Inexact Block Coordinate Descent* (BCD) algorithm / *Block Successive Upper bound Minimization* (BSUM) : at each iteration on variable W, we are not considering the original problem but an upper bound function (with the proximal term) $\|\mathbf{X}_i - w_i h_i\|_F^2 + \lambda \mathbf{D}_{ii}^1 + \frac{\gamma}{2} \|w_i - w_i^-\|_2^2$, notice that the inexactness is also contributed by the eigen-gap introduced by $\mu_i \ge \mu_{\min}$.

Extension : accelerated coordinate descent ... (1/3)

The "simplified" STA algorithm (the updates of D and H are hidden here)

Algorithm	3	STA	with	cyclic	indexing
-----------	---	-----	------	--------	----------

- 1: for k = 1 to itermax do
- 2: for i = 1 to r do
- 3: Update w_i by doing something
- 4: end for
- 5: end for

STA is a Inexact BCD with cycling indexing. That is, w_i is selected according to i = 1, 2, ..., r, 1, 2, ..., r, ...

In fact, cyclic indexing is not optimal. Acceleration can be made on using random indexing and extrapolation. The next page will discuss accelerated STA that, **in expectation**, converges faster.

Extension : accelerated coordinate descent ... (2/3)

Algorithm 4 Accelerated STA with random indexing

1: Set
$$v_i^k = w_i$$
 for $i = 1, 2, ..., r$.

- 2: for k = 1 to itermax do
- 3: Pick $i = i_k$ by random (with uniform probability)

4:
$$y_i^k = a_k v_i^k + (1 - a_k) w_i^k$$

5:
$$w_i^{k+1} = y_i^k - \frac{1}{L_{i_k}} \nabla_{i_k} f(y_i^k)$$

6:
$$v_i^{k+1} = b_k v_i^k + (1 - b^k) y^k - \frac{c_k}{L_{i_k}} \nabla_{i_k} f(y_i^k)$$

7: end for

where σ is strong convexity parameter and L_{i_k} is smoothness parameter of function $f(w_i)$. And parameters a, b, c are obtained by solving the following non-linear equations

$$c_k^2 - \frac{c_k}{r} = \left(1 - \frac{c_k\sigma}{r}\right)c_{k-1}^2, \ a_k = \frac{r - c_k\sigma}{c_k(r^2 - \sigma)}, \ b_k = 1 - \frac{c_k\sigma}{r}$$

Idea is similar to the "extrapolation induced acceleration" of the Nesterov's accelerated (full-)graident. 75/95

The "simplified" acceelrated STA algorithm :

Algorithm 5 Accelerated STA with random indexing

- 1: for k = 1 to itermax do
- 2: Pick $i = i_k$ by random
- 3: update w_i by doing something, together with two additional series v^k_i and y^k_i
- 4: end for

The trade off of faster convergence is the "randomness" : picking up repeated index is possible $i=1,{\bf 3},{\bf 3},{\bf 3},4,2,5,\ldots$

Solution is to use random shuffle instead of totally random. For example, $r=3,\,\mathrm{and}$

 $i=[1,3,2],[3,2,1],[3,1,2],[1,2,3],[2,1,3],\ldots$

But the analysis of the convergence property becomes complicated.

Extension : noise, outlier and robustness (1/3)

Formulation (P_0) and (P_1) are sensitive to **outlier** and **noise**.

• Outlier : a single outlier can kill the algorithm.



Extension : noise, outlier and robustness (1/3)

Formulation (P_0) and (P_1) are sensitive to **outlier** and **noise**.

• Outlier : a single outlier can kill the algorithm.



• Solution : robust norm on data fitting:

 $\|\mathbf{X} - \mathbf{W}\mathbf{H}\|_F^2 \xrightarrow{\text{changed to}} \|\mathbf{X} - \mathbf{W}\mathbf{H}\|_\phi^2$

• Examples of ϕ

- $\|\mathbf{X} \mathbf{W}\mathbf{H}\|_{2-1}^{2} + \lambda \operatorname{tr}(\mathbf{D}^{1}\mathbf{W}^{\top}\mathbf{W} + \mathbf{D}^{0}) (L_{2-1} \text{ norm, column robust})$ $\|\mathbf{X} - \mathbf{W}\mathbf{H}\|_{1}^{2} + \lambda \operatorname{tr}(\mathbf{D}^{1}\mathbf{W}^{\top}\mathbf{W} + \mathbf{D}^{0}) (L_{1} \text{ norm, matrix robust})$ $\|\mathbf{X} - \mathbf{W}\mathbf{H}\|_{p}^{2} + \lambda \operatorname{tr}(\mathbf{D}^{1}\mathbf{W}^{\top}\mathbf{W} + \mathbf{D}^{0}) (L_{p} \text{ norm, tunable})$ $\|\mathbf{X} - \mathbf{W}\mathbf{H}\|_{p}^{2} = \sum_{ij} B_{ij}(\mathbf{X} - \mathbf{W}\mathbf{H})_{ij}^{2} = \|B \odot (\mathbf{X} - \mathbf{W}\mathbf{H})\|_{F}^{2}$
 - 78 / 95

Extension : noise, outlier and robustness (2/3)

Formulation (P_0) and (P_1) are sensitive to **outlier** and **noise**. • **Noise** : *F*-norm is for additive Gaussian noise.

Extension : noise, outlier and robustness (2/3)

Formulation (P_0) and (P_1) are sensitive to **outlier** and **noise**.

- **Noise** : *F*-norm is for additive Gaussian noise.
- Other noises and corresponding denoising norms :
 - Kullback Leibler divergence for Poisson noise
 - Itakura Saito divergence for Gamma Expoential noise
 - ► Laplacian / Double Expoential noise , Uniform noise, Lorentz noise



Figure: Noises. Source : https://gimper.net/resources/noise-generator.576/95

Extension : noise, outlier and robustness (3/3)

Formulations (P_0) and (P_1) are sensitive to **outlier** and **noise**. Solution - 1 : solve $\|\mathbf{X} - \mathbf{W}\mathbf{H}\|_p^2 + \lambda \operatorname{tr}(\mathbf{D}^1\mathbf{W}^\top\mathbf{W} + \mathbf{D}^0)$ by Iterative Reweighted Least Squares (IRLS) Idea :

$$||x||_p = \left(\sum_i |x_i|^p\right)^{1/p} = \left(\sum_i w_i^2 |x_i|^2\right)^{1/2}$$

where weight $w_i = x_i^{\frac{p-2}{2}}$. For application w_i can be set as x_i^- .

i.e. Approximate $\|\mathbf{X} - \mathbf{WH}\|_p^2$ as a weighted L_2 norm problem.

Extension : noise, outlier and robustness (3/3)

Formulations (P_0) and (P_1) are sensitive to **outlier** and **noise**. Solution - 1 : solve $\|\mathbf{X} - \mathbf{W}\mathbf{H}\|_p^2 + \lambda \operatorname{tr}(\mathbf{D}^1\mathbf{W}^\top\mathbf{W} + \mathbf{D}^0)$ by Iterative Reweighted Least Squares (IRLS) Idea :

$$||x||_p = \left(\sum_i |x_i|^p\right)^{1/p} = \left(\sum_i w_i^2 |x_i|^2\right)^{1/2}$$

where weight $w_i = x_i^{\frac{p-2}{2}}$. For application w_i can be set as x_i^- .

i.e. Approximate $\|\mathbf{X} - \mathbf{WH}\|_p^2$ as a weighted L_2 norm problem.

Or, approximate matrix L_1 norm by IRLS :

$$\|\mathbf{X} - \mathbf{W}\mathbf{H}\|_{B}^{2} = \sum_{ij} B_{ij} (\mathbf{X} - \mathbf{W}\mathbf{H})_{ij}^{2}$$
$$\approx \|\mathbf{X} - \mathbf{W}\mathbf{H}\|_{1}^{2}$$
1

for $B_{ij} = \frac{1}{\|\mathbf{X} - \mathbf{W}\mathbf{H}\|_{ij} + \epsilon}$

82 / 95

Formulations (P_0) and (P_1) are sensitive to **outlier** and **noise**. **Solution - 2** : use L_{2-1} norm, or more specific :

$$\sum_{j=1}^{n} \frac{1}{2} \left(\| \mathbf{X}(:,j) - \mathbf{W}\mathbf{H}(:,j) \|_{2}^{2} + \epsilon \right)^{\frac{p}{2}} + \lambda \operatorname{tr}(\mathbf{D}^{1}\mathbf{W}^{\top}\mathbf{W} + \mathbf{D}^{0})$$

where $\epsilon > 0$ is smoothness constant and $p \in (0 \ 2]$ is robustness parameter.

FGM cannot be used directly on this formulation, need some modifications.

Small summary

$$\|\mathbf{X} - \mathbf{W}\mathbf{H}\|_{F}^{2} \longrightarrow \|\mathbf{X} - \mathbf{W}\mathbf{H}\|_{F}^{2} + \lambda \log \det(\mathbf{W}^{\top}\mathbf{W} + \delta I)$$

$$\|\mathbf{X} - \mathbf{W}\mathbf{H}\|_{\phi}^{2} \|\mathbf{X} - \mathbf{W}\mathbf{H}\|_{F}^{2} + \lambda \operatorname{tr}(\mathbf{W}^{\top}\mathbf{W} + (\delta - 1)\mathbf{I}_{r})$$

$$\|\mathbf{X} - \mathbf{W}\mathbf{H}\|_{F}^{2} + \lambda \operatorname{tr}(\mathbf{D}^{1}\mathbf{W}^{\top}\mathbf{W} + \mathbf{D}^{0})$$

$$\sum_{i} \|\mathbf{X}_{i} - w_{i}h_{i}\|_{F}^{2} + \lambda \operatorname{tr}(\mathbf{D}^{1}\mathbf{W}^{\top}\mathbf{W} + \mathbf{D}^{0})$$

84 / 95

Standard open problems

- How to tune λ
 - small noise $\lambda \to 0$
 - large noise $\lambda \to \infty$
 - $\lambda(N) = ?$
- How to find r (in real world application you don't know the r !)

Other directions

- On solving ${f H}$
- On super-big data
- Parallelism: divide-and-conquer / decompose-and-recombine.
- Acceleration by weighted formulation

A very difficult problem. How to tune λ dynamically within the iterations ?? That is, find an expression in the form as :

$$\lambda = \lambda(\mathbf{X}, \mathbf{W}_k, \mathbf{H}_k, k)$$

where k is the current number of iteration.

Don't expect I can give a global solution ! A rough idea of approaches :

- Simulated Annealing
- Dynamic approach
- Hybrid approach
- Pros : easy to implement
- Cons : very hard to establish theoretical convergence gurantee

Simulated Annealing parameter tuning of λ

- Annealing (metallurgy) = heat treatment of metal
- At the beginning, start with very high temperature to make coarse adjustment of the metal (hammering)
- Temperature is gradually decrease in the process, and graudally moving from coarse adjustment to fine adjustment
- Finally the metal is cooled down

Simulated Annealing parameter tuning of λ

- Annealing (metallurgy) = heat treatment of metal
- At the beginning, start with very high temperature to make coarse adjustment of the metal (hammering)
- Temperature is gradually decrease in the process, and graudally moving from coarse adjustment to fine adjustment
- Finally the metal is cooled down
- On the problem $f(\mathbf{W},\mathbf{H}) + \lambda G(\mathbf{W})$,
 - high temperature = starting with very large λ
 - Coarse adjustment of the metal = rotation of the convex hull
 - Temperature is gradually decrease = gradually decrease the value of λ
 - Fine adjustment = growth of convex hull
 - metal is cooled down $= \lambda_k$ is very close to zero
 - (Reminder to myself : refer to the "rotate.gif")

- Very primitive, not robust, non-determinstic
- Only works with cases that \mathbf{W}_0 is known : is $\log \det(\mathbf{W}^\top \mathbf{W} + \delta \mathbf{I}_r)$ really a good estimator of $\|\mathbf{W}_0 - \hat{\mathbf{W}}\|_F$? How about $\det \mathbf{W}^\top \mathbf{W}$?

 \implies Compare the models $\log \det(\mathbf{W}^{\top}\mathbf{W} + \delta \mathbf{I}_r)$ with $\det \mathbf{W}^{\top}\mathbf{W}$

Equivalent problems:

$$\begin{aligned} & (P_0) \quad \min \quad \|\mathbf{X} - \mathbf{W}\mathbf{H}\|_F^2 + \lambda \log \det(\mathbf{W}^\top \mathbf{W} + \delta \mathbf{I}_r), \\ & (P_1) \quad \min. \quad \operatorname{tr}(\mathbf{D}^1 \mathbf{W}^\top \mathbf{W} + \mathbf{D}^0) \text{ s.t. } \|\mathbf{X} - \mathbf{W}\mathbf{H}\|_F^2 \leq \epsilon \end{aligned}$$

For every ϵ in P_1 , there exists a λ in P_0 that both of them share the same solution. Solving P_1 does not involve parameter tuning.

Further improving STA : on H(1/3)

Recall the STA algorithm (with first 3 lines remoeved)



This line is to solve

$$\mathbf{H} \leftarrow \underset{\mathbf{H} \ge 0, \mathbf{1}_{r}^{\top} \mathbf{H} \le \mathbf{1}_{n}}{\operatorname{arg\,min}} \|\mathbf{X} - \mathbf{W}\mathbf{H}\|_{F}^{2} + \underbrace{\lambda \log \det(\mathbf{W}^{\top}\mathbf{W} + \delta \mathbf{I}_{r})}_{\text{a constant for }\mathbf{H}} \mathbf{M} + \underbrace{\lambda \log \det(\mathbf{W}^{\top}\mathbf{W} + \delta \mathbf{I}_{r})}_{\text{a constant for }\mathbf{H}} \mathbf{M} + \underbrace{\lambda \log \det(\mathbf{W}^{\top}\mathbf{W} + \delta \mathbf{I}_{r})}_{\text{a constant for }\mathbf{H}} \mathbf{M} + \underbrace{\lambda \log \det(\mathbf{W}^{\top}\mathbf{W} + \delta \mathbf{I}_{r})}_{\text{a constant for }\mathbf{H}} \mathbf{M} + \underbrace{\lambda \log \det(\mathbf{W}^{\top}\mathbf{W} + \delta \mathbf{I}_{r})}_{\text{a constant for }\mathbf{H}} \mathbf{M} + \underbrace{\lambda \log \det(\mathbf{W}^{\top}\mathbf{W} + \delta \mathbf{I}_{r})}_{\text{a constant for }\mathbf{H}} \mathbf{M} + \underbrace{\lambda \log \det(\mathbf{W}^{\top}\mathbf{W} + \delta \mathbf{I}_{r})}_{\text{a constant for }\mathbf{H}} \mathbf{M} + \underbrace{\lambda \log \det(\mathbf{W}^{\top}\mathbf{W} + \delta \mathbf{I}_{r})}_{\text{a constant for }\mathbf{H}} \mathbf{M} + \underbrace{\lambda \log \det(\mathbf{W}^{\top}\mathbf{W} + \delta \mathbf{I}_{r})}_{\text{a constant for }\mathbf{H}} \mathbf{M} + \underbrace{\lambda \log \det(\mathbf{W}^{\top}\mathbf{W} + \delta \mathbf{I}_{r})}_{\text{a constant for }\mathbf{H}} \mathbf{M} + \underbrace{\lambda \log \det(\mathbf{W}^{\top}\mathbf{W} + \delta \mathbf{I}_{r})}_{\text{a constant for }\mathbf{H}} \mathbf{M} + \underbrace{\lambda \log \det(\mathbf{W}^{\top}\mathbf{W} + \delta \mathbf{I}_{r})}_{\text{a constant for }\mathbf{H}} \mathbf{M} + \underbrace{\lambda \log \det(\mathbf{W}^{\top}\mathbf{W} + \delta \mathbf{I}_{r})}_{\text{a constant for }\mathbf{H}} \mathbf{M} + \underbrace{\lambda \log \det(\mathbf{W}^{\top}\mathbf{W} + \delta \mathbf{I}_{r})}_{\text{a constant for }\mathbf{H}} \mathbf{M} + \underbrace{\lambda \log \det(\mathbf{W}^{\top}\mathbf{W} + \delta \mathbf{I}_{r})}_{\text{a constant for }\mathbf{H}} \mathbf{M} + \underbrace{\lambda \log \det(\mathbf{W}^{\top}\mathbf{W} + \delta \mathbf{H}_{r})}_{\text{a constant for }\mathbf{H}} \mathbf{M} + \underbrace{\lambda \log \det(\mathbf{W}^{\top}\mathbf{W} + \delta \mathbf{H}_{r})}_{\text{a constant for }\mathbf{H}} \mathbf{M} + \underbrace{\lambda \log \det(\mathbf{W}^{\top}\mathbf{W} + \delta \mathbf{H}_{r})}_{\text{a constant for }\mathbf{H}} \mathbf{M} + \underbrace{\lambda \log \det(\mathbf{W}^{\top}\mathbf{W} + \delta \mathbf{H}_{r})}_{\text{a constant for }\mathbf{H}} \mathbf{M} + \underbrace{\lambda \log \det(\mathbf{W}^{\top}\mathbf{W} + \delta \mathbf{H}_{r})}_{\text{a constant for }\mathbf{H}} \mathbf{M} + \underbrace{\lambda \log \det(\mathbf{W}^{\top}\mathbf{W} + \delta \mathbf{H}_{r})}_{\text{a constant for }\mathbf{H}} \mathbf{M} + \underbrace{\lambda \log \det(\mathbf{W}^{\top}\mathbf{W} + \delta \mathbf{H}_{r})}_{\text{a constant for }\mathbf{H}} \mathbf{M} + \underbrace{\lambda \log \det(\mathbf{W}^{\top}\mathbf{W} + \delta \mathbf{H}_{r})}_{\text{a constant for }\mathbf{H}} \mathbf{M} + \underbrace{\lambda \log \det(\mathbf{W}^{\top}\mathbf{W} + \delta \mathbf{H}_{r})}_{\text{a constant for }\mathbf{H}} \mathbf{M} + \underbrace{\lambda \log \det(\mathbf{W}^{\top}\mathbf{W} + \delta \mathbf{H}_{r})}_{\text{a constant for }\mathbf{H}} \mathbf{M} + \underbrace{\lambda \log \det(\mathbf{W}^{\top}\mathbf{W} + \delta \mathbf{H}_{r})}_{\text{a constant for }\mathbf{H}} \mathbf{M} + \underbrace{\lambda \log \det(\mathbf{W}^{\top}\mathbf{W} + \delta \mathbf{H}_{r})}_{\text{a constant for }\mathbf{H}} \mathbf{M} + \underbrace{\lambda \log \det(\mathbf{W}^{\top}\mathbf{W} + \delta \mathbf{H}_{r})}_{\text{a constant for }\mathbf{H}} \mathbf{M} + \underbrace{\lambda \log \det(\mathbf{W}^{\top}\mathbf{W} + \delta \mathbf{H}_{r})}_{\text{a constant f$$

Further improving STA : on H(2/3)

Is the FGM (Fast gradient method on constrainted least sqaure on unit simplex) really the best ?

1: for K = 1 to itermax do for k = 1 to itermax do 2. $P = \mathbf{X}\mathbf{H}^{\top}$ and $Q = \mathbf{H}\mathbf{H}^{\top}$. 3. for i = 1 to r do 4 $w_{i} = \frac{\left[P_{i} - \sum_{j=1}^{i-1} w_{j} Q_{ji} - \sum_{j=i+1}^{r} w_{j}^{-} Q_{ji} + \gamma w_{i}^{-}\right]_{+}}{Q_{ii} + \lambda \mathbf{D}_{ij}^{1} + \frac{\gamma}{2}}$ 5: $\mu_i \leftarrow \mathsf{svd}(\mathbf{W}^\top \mathbf{W} + \delta \mathbf{I}_r) \text{ and } \mathbf{D}^1 = \mathsf{Diag}(\mu_{\min}^{-1})$ 6: 7: end for Update \mathbf{H} by FGM with $\mathbf{X}, \mathbf{W}, \mathbf{H}$ 8: end for 9:

10: end for

Further improving STA : on H(3/3)

Currently, yes.



Figure: Four methods on H.

* primal algorithm vs primal-dual algorithm.

We assumed r is known, it is only true for synthetic experiment. In real application, no one know the true r.

Idea : if input r is larger than the true r, when the minimum volume is 'achieved', then there will be (almost-)colinear columns in $\mathbf{W} \implies$ a way to auto-detect r!!

Reminder to me : run the large_r.gif in chrome.

Further issues : acceleration by weighted sum formulation

Since

$$\|\mathbf{X} - \mathbf{W}\mathbf{H}\|_{F}^{2} = \sum_{j} \|\mathbf{X}(:, j) - \mathbf{W}\mathbf{H}(:, j)\|_{2}^{2}.$$

What if the data fitting terms becomes

$$\sum_{j} \alpha_i \| \mathbf{X}(:,j) - \mathbf{W}\mathbf{H}(:,j) \|_2^2,$$

where α_i are weights.

Idea : increase the weight on the points that $\notin \mathsf{conv}(\mathbf{W})$ to speed up the fitting of the vertices of the next iteration



94 / 95

Last page - summary

- Introduce / review NMF $\|\mathbf{X} \mathbf{W}\mathbf{H}\|_F^2$
- Minimum volume NMF $\|\mathbf{X} \mathbf{W}\mathbf{H}\|_F^2 + \lambda \log \det(\mathbf{W}^\top \mathbf{W} + \delta \mathbf{I}_r)$
- Successive Trace Approximation $\mathrm{tr}(\mathbf{D}^1\mathbf{W}^\top\mathbf{W}+\mathbf{D}^0)$
- The STA algorithm and refinments
- BCD acceleration by randomization (on index)
- Convergence of STA (just rough idea)
- Robust STA (just rough idea) : $\|\mathbf{X} \mathbf{WH}\|_{\phi}$, $\phi \in \{1 \le p \le 2, 2-1\}$, Iterative reweighted least squares
- Experiments : some rough illustrations
- Some open / unsolved problems and further refinments

Slides (and code) avaliable at angms.science

- END OF PRESENTATION -

Extra - 1. What if some elements of \mathbf{H} is very very small

If some ${f H}$ of the data are very small, convex hull looks like conical hull



Extra - 1. What if some elements of \mathbf{H} is very very small

Solution-1 See them as outliers

Their \mathbf{H}_{ij} small \implies norm small \implies error small \implies ok if only a few of them

Solution-2 Augmenting W = [W' a] where a is a small vector. Note : r changes, and a cannot be 0 as matrix [W 0] is not full rank



Extra - 2. What if data are clustered

What if :



Extra - 2. What if data are clustered

Still works, but other method will be better (e.g. some clustering method such as K-means)

