

# クラスタリングと非負行列の因数分解

情報科学類 3 年 江畑 拓哉 (201611350)

# Outline

## 1 非負値行列因子分解

ある行列  $\mathbf{A} \in \mathbb{R}^{m \times n}$  が与えられ、非負の要素を持よように制限された  $k$  ランクの近似を行いたいと考えた時、言い換えれば  $\mathbf{W} \in \mathbb{R}^{m \times k}$  と  $\mathbf{H} \in \mathbb{R}^{k \times n}$  を仮定して以下の式を解きたい場合について考える。

$$\min_{\mathbf{W} \geq 0, \mathbf{H} \geq 0} \|\mathbf{A} - \mathbf{WH}\|_F$$

$$\text{where } \|\mathbf{A}\|_F = \sqrt{\sum_i \sum_j |a_{ij}|^2} \text{ means Frobenius norm} \quad (9.3)$$

同時に  $\mathbf{W}$  と  $\mathbf{H}$  を最適化しようとした時、この問題は非線形になる。

しかし、行列の一つが既にわかっている場合、例えば  $\mathbf{W}$  がわかっている場合で言えば、この  $\mathbf{H}$  を求める問題は、非負の制限がついた右側の行列についての最小二乗問題であると言える。

従って元の問題に対する最も一般的な解決法は、交互最小二乗法 (ALS) を使うことである。

# Alternating nonnegative least squares のアルゴリズム

## Alternating nonnegative least squares algorithm

- 1 初期値  $\mathbf{W}^{(1)}$  を与える。
- 2  $k = 1, 2, \dots$  と収束するまで繰り返す。
  - $\min_{\mathbf{H} \geq 0} \|\mathbf{A} - \mathbf{W}^{(k)} \mathbf{H}\|_F$  を解き、 $\mathbf{H}^{(k)}$  を得る
  - $\min_{\mathbf{W} \geq 0} \|\mathbf{A} - \mathbf{W} \mathbf{H}^{(k)}\|_F$  を解き、 $\mathbf{W}^{(k+1)}$

しかし、近似  $\mathbf{WH}$  は単一のものではない。ここで正の対角要素を持つ任意の対角行列  $\mathbf{D}$  とその逆行列を要素間に適用しすることが出来る。

$$\mathbf{WH} = (\mathbf{WD})(\mathbf{D}^{-1}\mathbf{H})$$

ある要素の増大と他の要素の減衰を防ぐために、毎反復ごとにそれらの一つを正規化する必要がある。一般的な正規化は、 $\mathbf{W}$  の各列の最大要素が 1 になるように  $\mathbf{W}$  の列をスケーリングすることである。

$\mathbf{A}$  と  $\mathbf{H}$  の列要素として、 $\mathbf{a}_j$  と  $\mathbf{h}_j$  を置く。各行の要素を一つずつ書き下すと、この最小二乗問題は以下の  $n$  個の独立したベクトルの最小二乗問題と等しいものとみなすことが出来る。

$$\min_{\mathbf{h}_j \geq 0} \|\mathbf{a}_j - \mathbf{W}^{(k)} \mathbf{h}_j\|_2 \text{ where } j = 1, 2, \dots, n$$

これらは 23 章で紹介されるアクティブセットアルゴリズムで解決される。この行列を転置することで、 $\mathbf{W}$  を決定するの最小二乗問題は、独立な  $m$  のこのベクトルの最小二乗問題に変換される。即ち ALS アルゴリズムのコア部分は擬似的な MATLAB コードで以下のように表すことが出来る。

## 疑似 MATLAB コード

```
1 while (not converged)
2     [W] = normalize(W);
3     for i = 1:n
4         H(:,i) = lsqnonneg(W, A(:,i));
5     end
6     for i = 1:m
7         w = lsqnonneg(W, A(:, i));
8         W(i,:) = w';
9     end
10 end
```

非負値行列因子分解のためのアルゴリズムは多くの種類がある。前頁のアルゴリズムは、非負な最小二乗法のためのアクティブセット法にかなり時間がかかってしまうという欠点がある。より簡易な代替手段として、部分 QR 分解  $\mathbf{W} = \mathbf{QR}$  を用いることで非負という制約がない最小二乗解を得ることが出来る。そして  $\mathbf{H}$  におけるすべての負の要素はゼロと等しいとみなすことが出来る。これは  $\mathbf{W}$  の計算においても同様の議論をすることが出来る。

$$\mathbf{H} = \mathbf{R}^{-1} \mathbf{Q}^T \mathbf{A}$$

疑似 MATLAB コード (上の例に基づけば  $\mathbf{V} = \mathbf{A}$ )

```
1 while (not converged)
2     W = W.*(W >= 0);
3     H = H.*(W'*V)./((W'*W)*H+epsilon);
4     H = H.*(H>=0);
5     W = W.*(V*H')./(W*(H*H')+epsilon);
6     [W,H] = normalize(W, H);
7 end
```

$\epsilon$  は極小の値であり、これはゼロ除算を避けるために用いられている。  $\cdot *$  や  $\cdot /$  で表される行列操作はそれぞれの構成要素についての演算で、

$$H_{ij} := H_{ij} \frac{(W^T A)_{ij}}{(W^T W H)_{ij} + \epsilon}, \quad W_{ij} := W_{ij} \frac{(A H^T)_{ij}}{(W H H^T)_{ij} + \epsilon}$$

尚このアルゴリズムは勾配降下法と考えることが出来る。



非負値行列分解には非常に多くの重要な用途があるため、このアルゴリズム開発は活発に行われている。例えば、反復法を用いた終了基準を見つける問題は未だ良い解決策を見つけたとは言い難い。

非負値行列分解  $\mathbf{A} \approx \mathbf{WH}$  はクラスタリングにも用いられている。各データを表すベクトル  $\mathbf{a}_j$  は、もし  $\mathbf{h}_{ij}$  が  $\mathbf{H}$  の  $j$  列の最大の要素であるなら、それはクラスタ  $i$  に割り当てられる。

さらにこの分解法は以下のような分野でも用いられている。

- 文書分類
- 電子メールの監視
- 音楽の記譜 (楽譜にすること)
- バイオインフォマティクス
- スペクトル分析

# 初期化

非負値行列分解のアルゴリズムにはいくつかの問題がある。それは全体での最適解が求まる保証がないということである。このアルゴリズムではしばしば収束が遅いことや順最適解 (厳密解ではない) になってしまうことがある。良好な初期の近似を計算するための効率的な手法として、 $A$  の SVD に基づいて行うというものがある。最初の  $k$  個の特異の三つの組  $(\sigma_i, \mathbf{u}_i, \mathbf{v}_i)_{i=1}^k$  は、フロベニウスノルムにおいて  $A$  の最適なランク  $k$  の近似を与える。

もし  $A$  が非負な行列であったならば、 $u_i$  や  $v_1$  が非負であることは明らか。(Section 6.4)

つまりもし  $A = U\Sigma V^T$  が  $A$  の SVD であるならば、特異ベクトル  $u_1$  を  $W^{(1)}$  の最初の列であるとする事が出来る。(同様に以降のアルゴリズムのため、 $v_1^T$  を初期近似  $H^{(1)}$  の第1行であるとする。)

次の最良なベクトル  $u_2$  は直交性が満たされているために負の成分を有する可能性が非常に高い。しかし行列  $C^{(2)} = u_2 v_2^T$  を計算しすべての負の成分をゼロにすることで非負な行列  $C_+^{(2)}$  を得ると、この行列の最初の特異ベクトルは非負であることがわかる。さらにそれは、これが  $u_2$  の合理的で良い近似であると考えることが出来るので、 $W^{(1)}$  の第2列として取り上げることが出来る。

前頁の手続きを MATLAB を使って簡潔に書き下すと以下のような  
る。

\*  $[U, S, V] = svds(A, k)$  は Lanczos 法を用いることで、 $k$  個の最大  
特異値及び対応する特異ベクトルとを計算する。標準な SVD 関数であ  
る  $svd(A)$  はすべての分解を計算するがこれはかなり遅く、特に行列が  
大きな疎行列のときはより遅くなる。

## MATLAB

```
1 [U,S,V] = svds (A, k) % Compute only the k largest singular
2 W (:,1) = U (:,1);    % values and the corresponding vectors
3 for j = 2:k
4     C = U (:,j)*V (:,j)';
5     C = C .* (C>=0);
6     [u, s, v] = svds (C, 1);
7     W (:,j) = u;
8 end
```

## 例 9.4

ランク 2 の行列  $\mathbf{A}$  の非負値分解の例を図 9.3 に示す。ここでは初期化をランダムな値で行ったものと、SVD ベースで行ったものを比較している。ランダムな値で初期化したものは収束がより遅くなっており、10 回反復させても収束したとは言い難い。これに対して SVD ベースで初期化したものの相対近似誤差は 0.574 であることがわかる。(尚 k-means 法においてこの誤差は 0.596 であった) 更にいくつかのケースでランダムに初期化したものは最適でない準最適な値に収束してしまった。

## 図 9.3

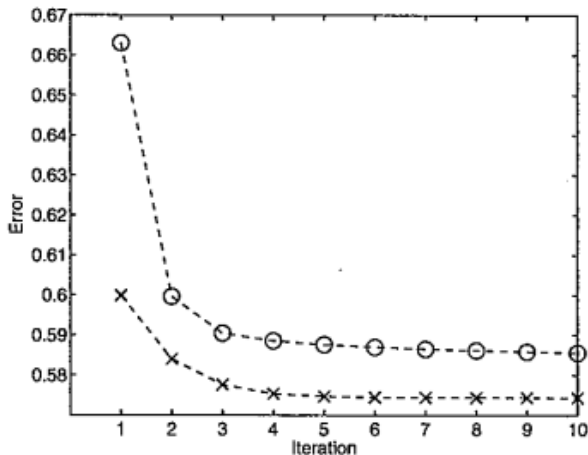


Figure: 9.3 反復回数を横軸とした相対近似誤差のグラフ。上のカーブは初期化をランダムに行ったもの、下のカーブは SVD ベースで初期化を行ったものである。

SVD ベースの初期化をするこの分解は具体的に以下ようになった。

$$WH = \begin{pmatrix} 0.3450 & 0 \\ 0.1986 & 0 \\ 0.1986 & 0 \\ 0.6039 & 0.1838 \\ 0.2928 & 0 \\ 0 & 0.5854 \\ 1.0000 & 0.0141 \\ 0.0653 & 1.0000 \\ 0.8919 & 0.0604 \\ 0.0653 & 1.0000 \end{pmatrix} \begin{pmatrix} 0.7740 & 0 & 0.9687 & 0.9120 & 0.5251 \\ 0 & 1.0863 & 0.8214 & 0 & 0 \end{pmatrix}$$

前頁のそれは分解の処理を打ち切ることが出来る。最初の四つの文書は基底ベクトルによって表されており、これは Google-related keywords のための大きな要素を持っている。これに対して最後の文書は 1 つめの基底ベクトルによって表されているが、この座標値は先述の四つの文書に比べて小さくなっている。

この手法では、ランク 2 の近似は Google-related contents を強調するが、"football-document" は強調しない。

11 章では、他の低ランクの近似について学ぶが、(例えば SVD をベースとしたもの) これらも類似の効果を確認することが出来る。



対して、ランク 3 の近似を計算した時には以下の値を得ることが出来る。

$W$  の三番目のベクトルは、本質的に "football" についての基底であり、その一方で他の 2 つのベクトルは Google-related document を示している基底である。

$$WH = \begin{pmatrix} 0.2516 & 0 & 0.1633 \\ 0 & 0 & 0.7942 \\ 0 & 0 & 0.7942 \\ 0.6924 & 0.1298 & 0 \\ 0.3786 & 0 & 0 \\ 0 & 0.5806 & 0 \\ 1.0000 & 0 & 0.0444 \\ 0.0589 & 1.0000 & 0.0007 \\ 0.4237 & 0.1809 & 1.0000 \\ 0.0589 & 1.0000 & 0.0007 \end{pmatrix} \begin{pmatrix} 1.1023 & 0 & 1.0244 & 0.8045 & 0 \\ 0 & 1.0815 & 0.8315 & 0 & 0 \\ 0 & 0 & 0.1600 & 0.3422 & 1.1271 \end{pmatrix}$$