

平行移動不変な非負値行列因子分解とその分析

Shift-Invariant Non-negative Matrix Factorization and its Empirical Evaluation

鈴木 慶介[†] 今井 英幸[†] 張 若霓[†] 瀧川 一学^{†‡} 湊 真一[†]
 Keisuke Suzuki Hideyuki Imai Ruoni Zhang Ichigaku Takigawa Shin-ichi Minato

1. はじめに

NMF(非負値行列因子分解, Non-negative Matrix Factorization) は Lee と Seung[1] により, 要素が非負値のデータ行列を二つの非負値行列の積に分解する手法として提案され, 様々な応用で広く利用されている. NMF は非負値制約により, 分解前のデータを分解後の行列で“加法的”に表現できるため, 分解結果がより直感的に解釈し易くなる. この特性により, ソフトクラスタリングや特徴抽出の手法としてよく用いられてきた. しかし, NMF はデータの平行移動に対して不変ではないことが分かっている. これはデータの平行移動によってクラスタリング結果が変わりうることを意味する. 例えば, k -means 法など, 通常のクラスタリング手法では, データの相対的な分布形状により結果が定まるため, 平行移動に対して結果は不変である. 一方, NMF ではデータを分解行列の列ベクトルの錐結合で表現することに対応するため, 原点の取り方に依存しない相対的な分布形状の表現とならず, 平行移動不変性が成り立たない. 本研究では, 錐結合による分解を凸結合による分解とすることで, 平行移動に対して不変な結果を出力する NMF を提案する.

2. NMF(非負値行列因子分解)

データとして, 非負の要素を持つ p 次元のベクトル $\mathbf{x} = (x_1, x_2, \dots, x_p)^T, x_i \geq 0$ が n 個与えられ, その全体をデータ行列 $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$ で表す. また, 行列 \mathbf{X} の要素が全て非負であることを $\mathbf{X} \geq 0$ と表記する. NMF は行列 $\mathbf{X} \in \mathbb{R}^{p \times n}, \mathbf{X} \geq 0$ を

$$\mathbf{X} = \mathbf{F}\mathbf{G}^T, \quad \mathbf{F} \geq 0, \mathbf{G} \geq 0$$

となるような二つの非負値行列 $\mathbf{F} \in \mathbb{R}^{p \times k}$ と $\mathbf{G} \in \mathbb{R}^{n \times k}$ の積に分解しようとする計算に対応する. ここで, k はデータ全体を表現するのに必要な代表ベクトルの数を指し, プログラム実行者があらかじめ決めておく. クラスタリングとして見れば, p 次元のデータベクトル n 個を, p 次元の代表ベクトル k 個の錐結合で表現することに対応する. 従って, k はデータの次元 p , データの個数 n よりも小さい値とする. このとき, 一般には等号は成り立たず上記のような分解は得られない. そこで, NMF では \mathbf{X} と $\mathbf{F}\mathbf{G}^T$ の二乗誤差を目的関数 $J(\mathbf{F}, \mathbf{G})$ とし,

$$J(\mathbf{F}, \mathbf{G}) = \|\mathbf{X} - \mathbf{F}\mathbf{G}^T\|_F^2$$

を最小化する二つの行列に分解する. ここで $\|\cdot\|_F$ は行列のフロベニウスノルムである. なお, この二乗誤差の最小化は x_{wi} が, $\sum_{t=1}^k f_{wt}g_{ti}$ を平均とした正規分布により生成されると仮定した場合の統計的推定問題に対応する. \mathbf{F}, \mathbf{G} の \mathbf{X} からの乖離度を表現する基準に応じて異なる最適化問題に帰着するが, 本研究では上記の $J(\mathbf{F}, \mathbf{G})$ を対象とする.

NMF では, 観測データの各々 \mathbf{x}_i を

$$\mathbf{x}_i \approx \sum_{t=1}^k \mathbf{f}_t g_{ti} \quad g_{ti} \geq 0, i \in \{1, \dots, n\} \quad (1)$$

という形で, 行列 \mathbf{G} の要素 g_{ti} を係数, \mathbf{F} の列ベクトルを基底ベクトルとした \mathbf{f}_t の錐結合で最小二乗近似しようとしていると解釈できる. 従って, NMF による分解は元の観測データを最も良く近似する錐包を張る k 個の基底ベクトルを探索する作業と言える. また, 各基底ベクトル \mathbf{f}_t の係数 g_{ti} の比をソフトクラスタリングの結果として解釈することが可能である. 目的関数 $J(\mathbf{F}, \mathbf{G})$ は両変数に対し同時に凸ではないので解析的に解を定めることは出来ないが, 以下の二つの更新式の反復で求められることが知られている [2].

$$\mathbf{F} \leftarrow \mathbf{F} \odot \frac{\mathbf{X}\mathbf{G}}{\mathbf{F}\mathbf{G}^T\mathbf{G}}, \quad \mathbf{G} \leftarrow \mathbf{G} \odot \frac{\mathbf{X}^T\mathbf{F}}{\mathbf{G}\mathbf{F}^T\mathbf{F}}$$

ここで, \odot は要素同士の乗算を行う演算子である. また行列同士の除算に関しても, 本文では要素同士の除算とする. この更新について目的関数は単調減少し局所最適解に収束するが, 必ずしも大域最適解には収束しない.

3. Shift-Invariant NMF

係数 g_{ti} をデータ点 \mathbf{x}_i の代表ベクトル \mathbf{f}_t への寄与度と解釈し, ソフトクラスタリングを行うことを考える. データ点 \mathbf{x}_i が平行移動 \mathbf{a} を受け $\mathbf{x}_i + \mathbf{a}$ となると, (1) の右辺は

$$\sum_{t=1}^k (\mathbf{f}_t + \mathbf{a})g_{ti} = \sum_{t=1}^k \mathbf{f}_t g_{ti} + (\sum_{t=1}^k g_{ti})\mathbf{a} \quad (2)$$

となり, 第二項分の乖離が生じるため, 実際には基底ベクトルとして $\mathbf{f}_t + \mathbf{a}$ は選択されず, クラスタへの寄与度 g_{ti} も当然変化してしまう. 従って, 一般に NMF には平行移動による不変性がないことが分かる.

しかし, クラスタリングにおいては, 移動前, 移動後に関わらず観測データに対し基底との位置的な関係が不変であることが望ましい. つまり, データ行列 \mathbf{X} に対し, 列ベクトル $\mathbf{a} = (a_1, \dots, a_p)^T$ だけ平行移動を行ったデータを NMF で分解したとき, 平行移動前のデータが $\mathbf{X} \approx \mathbf{F}\mathbf{G}^T$ と分解されているとすると,

$$\mathbf{X} + \mathbf{a}\mathbf{1}_n^T \approx \tilde{\mathbf{F}}\tilde{\mathbf{G}}^T = (\mathbf{F} + \mathbf{a}\mathbf{1}_n^T)\mathbf{G}^T$$

の様に, 基底行列 \mathbf{F} に含まれる基底ベクトル \mathbf{f}_i を同じくベクトル \mathbf{a} だけ平行移動させた行列と, \mathbf{G} による積で表現されなければならない. ここで $\mathbf{1}_n = (1, 1, \dots, 1)^T$ であり, $\mathbf{1}$ を n 個要素を持つような列ベクトルとして定義する. 式 (2) から示唆される通り, 以上の条件を満たす Shift-Invariant NMF (SINMF) は以下の形で定義できる.

$$\min J(\mathbf{F}, \mathbf{G}) \text{ subject to } \mathbf{F} \geq 0, \mathbf{G} \geq 0, \sum_{t=1}^k g_{ti} = 1$$

ここで g_{ti} は \mathbf{G} の t 行 i 列目の要素を示す. 以上の制約により, データ行列の非負値性を崩さない範囲での平行移動 $\mathbf{a}, \mathbf{X} + \mathbf{a}\mathbf{1}_n^T \geq 0$ の下で, SINMF は目的関数が不変である. 両手法の差として, 基底による空間の表現方法の違いが挙げられる. 通常の NMF は基底ベクトルの錐結合により張られる錐包でデータを近似するが, SINMF は k 本の基底の凸結合により張られる $k-1$ 次元単体でデータを表現する. このとき, 単体によって最も誤差なくデータを表現するような

[†]北海道大学大学院 情報科学研究科

[‡]科学技術振興機構 (JST) さきがけ

頂点(基底)の配置は原点の取り方に依らない。また、NMFとSINMFの求める基底 \mathbf{F} は一般に異なる。SINMFの拡張として、アフィン変換に対して不変で、負値データを扱える手法(Affine-Invariant Semi-NMF)[3]も構成できることが分かっている。

図1にNMFとSINMFのある平行移動に対する基底の変化の違いの様子を示す。また、図2に、その平行移動に対して、どれほどの基底のずれが連続的に生じるかを示す。ここで「ずれ」の尺度として $\text{gap}(\mathbf{F}, \tilde{\mathbf{F}}_j, \mathbf{a}_j) = \|\mathbf{F} + \mathbf{a}_j \mathbf{1}_k^T - \tilde{\mathbf{F}}_j\|_F^2$ で定義される関数を用いる。ここで、 \mathbf{F} は平行移動前の基底、 \mathbf{a}_j は \mathbf{a} を始点から終点に向けて100分割した時の j 番目の点で表されるベクトル、 $\tilde{\mathbf{F}}_j$ は $\mathbf{X} + \mathbf{a}_j \mathbf{1}_n^T$ を分解して得られた基底である。

図1はデータ行列 $\mathbf{X} \in \mathbb{R}^{3 \times 100}$ に対して、 $k=2$ としてNMFとSINMFにより得られた2本の基底ベクトルを二次元表示したものである。それぞれ、(a)NMF:平行移動前、(b)NMF:平行移動後のデータと基底、(c)SINMF:平行移動前、(d)SINMF:平行移動後を表す。NMFは(b)(d)に矢印で示されている平行移動ベクトル \mathbf{a} に対し基底とデータとの相対的な距離が大きく変化しているが、SINMFの場合平行移動に対して基底とデータとの距離の変化は非常に少なく、そのデータの特徴を一貫して保持している。また図2より、平行移動に対して基底のずれが増加している様子が見てとれる。NMFでは、一般にずれは0にならないが、平行移動の大きさに対して単調に増加するとは限らない。

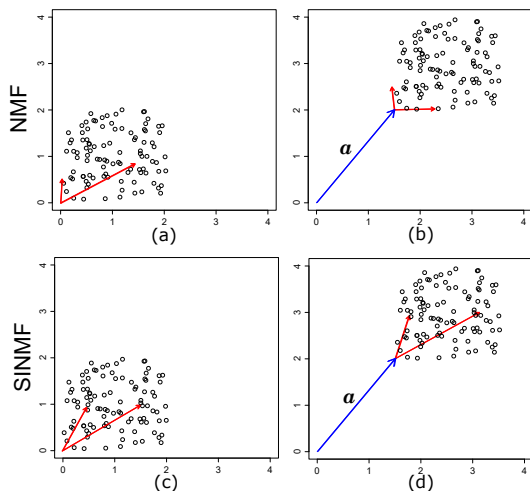


図1: 平行移動 \mathbf{a} の前後での基底ベクトルの変化の違い

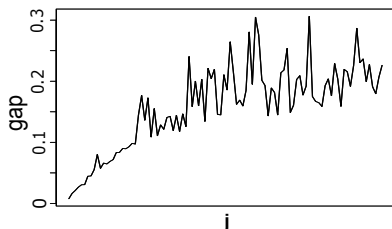


図2: 平行移動 $j \cdot \mathbf{a}/100$ に対するNMFの基底のずれ

4. 数値例

実データを用いてNMFとSINMFの比較例を示す。今回は文書クラスタリングを比較の方法に採用する。データセットとして、文書クラスタリングの数値実験に汎用的に用いられているニュース文書のデータセット Reuters21578を使用した。結果として得られる分解の解釈を容易にするため、全ての21578本の記事のうち、その中から原油に関する20本

の文書、企業の合併・買収に関する50本の文書を選択した(この場合実質的なクラスタは2つとなる)。このサブセットは、文書数70、出現した単語の種類2010(Stop Wordを除く)のデータであり、「原油関連」と「企業の合併・買収関連」の2つのクラスタを持つ。

このデータ行列 $\mathbf{X} \in \mathbb{R}^{2010 \times 70}$ に対して、基底数 $k=2$ として、 $\mathbf{F} \in \mathbb{R}^{2010 \times 2}$ 、 $\mathbf{G} \in \mathbb{R}^{70 \times 2}$ を求め、既知の2つのクラスタとの関係を見る。収束条件は、更新前後の目的関数値の誤差が定数内 $\epsilon = 10^{-7}$ かどうかで判定を行った。初めに、通常のNMF、SINMFで \mathbf{X} の分解を行って得られた第一基底における出現頻度上位10語を示す。また基底はノルムが1になるよう正規化した。

表1: NMF, SINMFにおけるそれぞれの \mathbf{f}_1

移動前	NMF \mathbf{f}_1	移動前	SINMF \mathbf{f}_1
dlrs	0.432	dlrs	0.438
company	0.256	pct	0.302
mln	0.223	mln	0.274
express	0.218	company	0.266
american	0.217	inc	0.215
stock	0.212	shares	0.214
pct	0.212	reuter	0.209
shearson	0.207	stock	0.187
analysis	0.175	share	0.140
inc	0.163	offer	0.140

次に、平行移動 $\mathbf{a} = (5, 5, \dots, 5)^T \in \mathbb{R}^{2010}$ (全ての単語の出現回数を5回増加)を加えたデータセット $\mathbf{X} + \mathbf{a} \mathbf{1}_n^T$ をそれぞれのNMFで分解した結果得られた第一基底を以下に示す。

表2: NMF, SINMFにおける平行移動後の $\tilde{\mathbf{f}}_1$

移動後	NMF $\tilde{\mathbf{f}}_1$	$\mathbf{f}_1 + \mathbf{a}$	移動後	SINMF $\tilde{\mathbf{f}}_1$	$\mathbf{f}_1 + \mathbf{a}$
dlrs	0.031	0.032	dlrs	0.032	0.031
pct	0.029	0.027	pct	0.029	0.029
company	0.028	0.028	company	0.028	0.028
mln	0.024	0.027	mln	0.027	0.028
shares	0.027		shares	0.027	0.027
inc	0.027	0.027	inc	0.027	0.027
reuter	0.026		reuter	0.027	0.027
stock	0.026	0.027	stock	0.026	0.026
offer	0.025		offer	0.025	0.025
share	0.025		share	0.025	0.025

ここで、ベクトルの要素の値は特徴の中でその単語(次元)がどれほどの重要度を占めているかを測る相対的指標にあたる。表1, 2より通常のNMFとSINMFにより生成された基底は、平行移動前、後、どちらもデータの特徴をよく表す単語が上位にとられているが、SINMFによる移動後の基底は、 $\mathbf{f}_1 + \mathbf{a}$ で表される真に望まれる基底に近い。ここで空白のセルは、その単語が真に望まれる基底の頻度上位10位に入らなかったことを意味する(順位が入れ替わっている)。これにより、通常のNMFに比べ特徴が平行移動に対し変化しない傾向が見てとれる。

5. まとめ

通常のNMFが平行移動に対して基底が変化する特性を実験を通して示し、今回新たに提案したSINMFがデータの平行移動に対して特徴を保存する点を分析した。今後の課題としては、テンソル分解への応用などが考えられる。

参考文献

- [1] D. D. Lee and S. H. Sebastian, Learning the parts of objects by non-negative matrix factorization. *Nature*, vol.401, pp.788-791, 1999.
- [2] D. D. Lee and S. H. Sebastian, Algorithms for Non-negative Matrix Factorization. *Advances in Neural Information Processing Systems*, 13, pp.556-562, 2001.
- [3] 張若霓, 今井英幸, Designing affine transformation based semi-nonnegative matrix factorization. 日本計算機統計学会第30回大会, 京都市, 2016年5月19日-20日.