

NMF

Non Negative Matrix Factorization

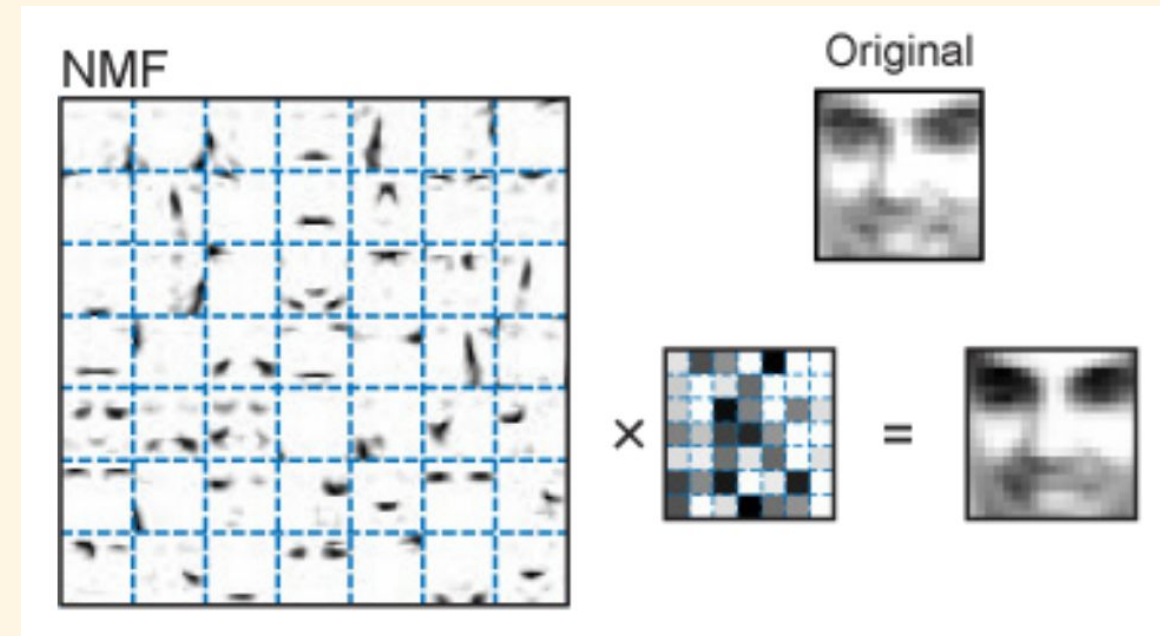
非負值行列因子分解

NMFの背景

画像処理や音声処理分野で活用

→画素値、スペクトル、頻度(単語の頻度)、個数(語彙数)etc

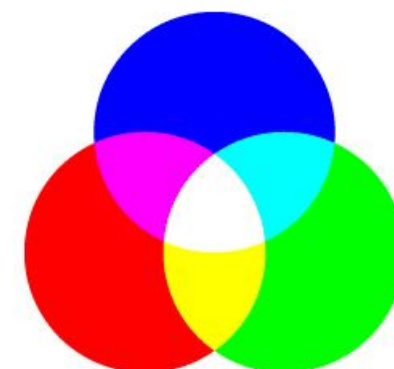
→実世界データは非負値であることが多い



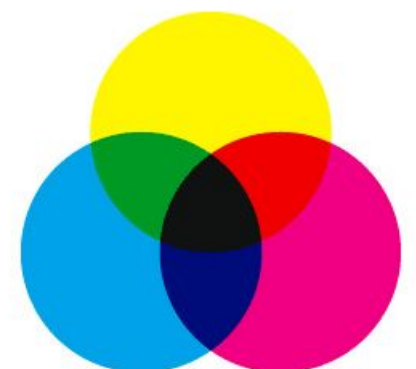
観測データの性質

- ・観測データの非負性→RGB、単語頻度、音の周波数は全て非負
- ・基本となる特徴の非負性→波長が負の値である光は観測不可
- ・重ね合わせにおける非負性→観測データは特徴の重ね合わせで表現

eg)人間が知覚可能な色は光の三原色の重ね合わせによって表現される



光の三原色(RGB)



色の三原色(CMYK)

NMFの仮定

仮定: 観測情報 y は基底 h の重み u 付き和で表現可能(加法性)
→ NMFの目的は基底 H と重み U を求めること

仮定: 基底 H および重み U は非負

eg) 任意の色 = 赤 * A + 青 * B + 緑 * C (A, B, C は重み)

※画素や周波数等は近似的に加法性が成立としてNMFを適用している

N 個の観測ベクトル

$\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_N$

基底が M 個あると仮定し、基底ベクトル \mathbf{h}_m 、結合係数 $u_{m,n}$ で表すと、

$$\mathbf{y}_n \simeq \sum_{m=1}^M \mathbf{h}_m u_{m,n} \quad (n = 1, \dots, N)$$

NMFの行列表現

観測ベクトルを並べたデータ行列を

$$\mathbf{Y} = [\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_N] = (y_{k,n})_{K \times N}$$

基底ベクトルを並べて基底行列を

$$\mathbf{H} = [\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_M] = (h_{k,m})_{K \times M}$$

結合係数を並べて係数行列を

$$\mathbf{U} = (u_{m,n})_{M \times N}$$

→観測データ行列は基底行列と結合係数行列の積で表現可能

$$\mathbf{Y} \simeq \mathbf{H}\mathbf{U}$$

$$(2 \times 3 \text{ 行列}) \times (3 \times 1 \text{ 行列}) = (2 \times 1 \text{ 行列})$$

$$\begin{pmatrix} a1 & a2 & a3 \\ b1 & b2 & b3 \end{pmatrix} \times \begin{pmatrix} c \\ d \\ e \end{pmatrix} = \begin{pmatrix} a1 \cdot c & a2 \cdot d & a3 \cdot e \\ b1 \cdot c & b2 \cdot d & b3 \cdot e \end{pmatrix}$$

一般に

$$(A \times B) \text{ 行列} \times (B \times C) \text{ 行列} \\ = (A \times C) \text{ 行列}$$

NMFの原理図

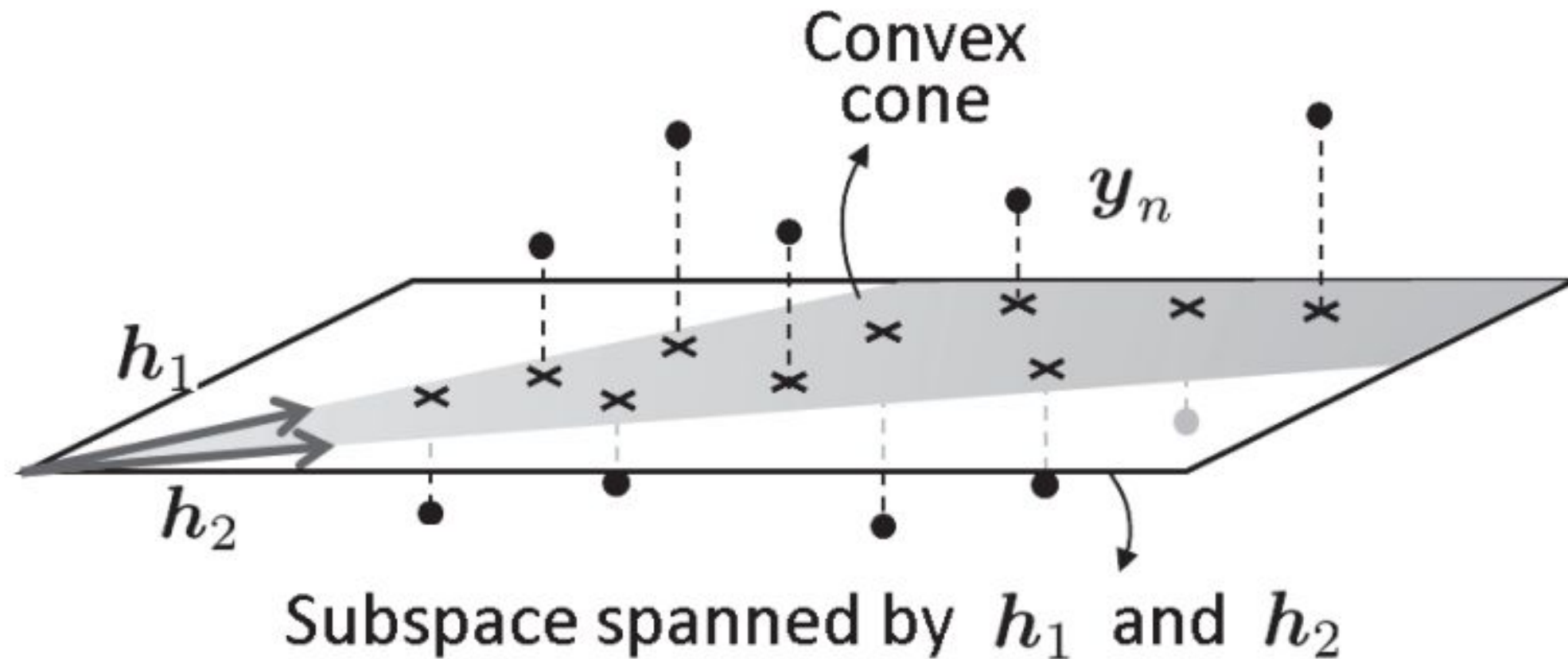


図2 非負値基底の非負結合

任意の観測データが基底ベクトル h_1 、 h_2 から構成される凸錐上に射影
→次元を落とした近似として解釈可能

$$y_n \simeq u_1 \mathbf{h}_1 + u_2 \mathbf{h}_2$$

閑話：凸錐

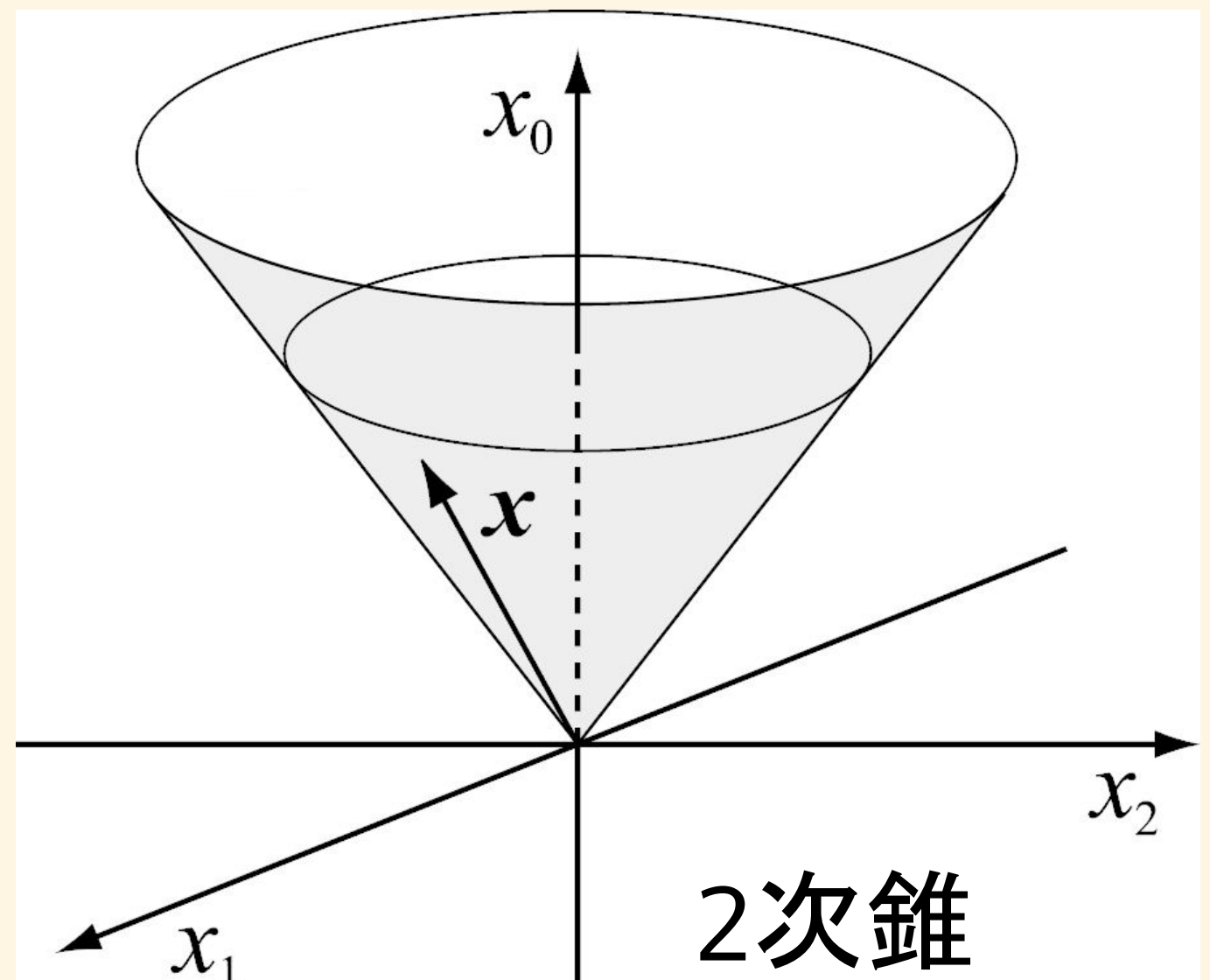
錐 (*cone*) とは次のような集合を指す

$$\boldsymbol{x} \in C, \quad \alpha \in [0, \infty) \rightarrow \alpha \boldsymbol{x} \in C$$

錐が凸集合であるならば、それは凸錐である。

$$eg) C = \{\boldsymbol{x} \in \boldsymbol{R}^n \mid \boldsymbol{x} = \sum_{i=1}^m \alpha_i \boldsymbol{\alpha}^i, \alpha \geq 0\}$$

直交ベクトルによって形成される錐状
の領域の例



NMFの性質

性質

- ・低次元近似

→基底ベクトルの次元Mは観測ベクトルの次元Kやデータ数Nに対して一般に小さく設定(M=Kの時H=E単位行列、M=Nの時U=E)

$$M < \min(K, N)$$

- ・観測データの共起成分をグルーピングする傾向(ここ謎)

→共起する成分をひとまとめにして基底と捉えるっぽい

- ・係数行列のスパース性

→あるベクトルの近似に有効でない基底の影響を除外したい

→有効でない基底の係数は0

- ・分解の一意性

→一般的にH、Uは一意に定まらない($HSS^{\top}(-1)U$ も解になりえる)

→最適化の際、次のような対策を取る

- ・スパースネス(行列の全要素に対する零要素の割合)の条件を付ける

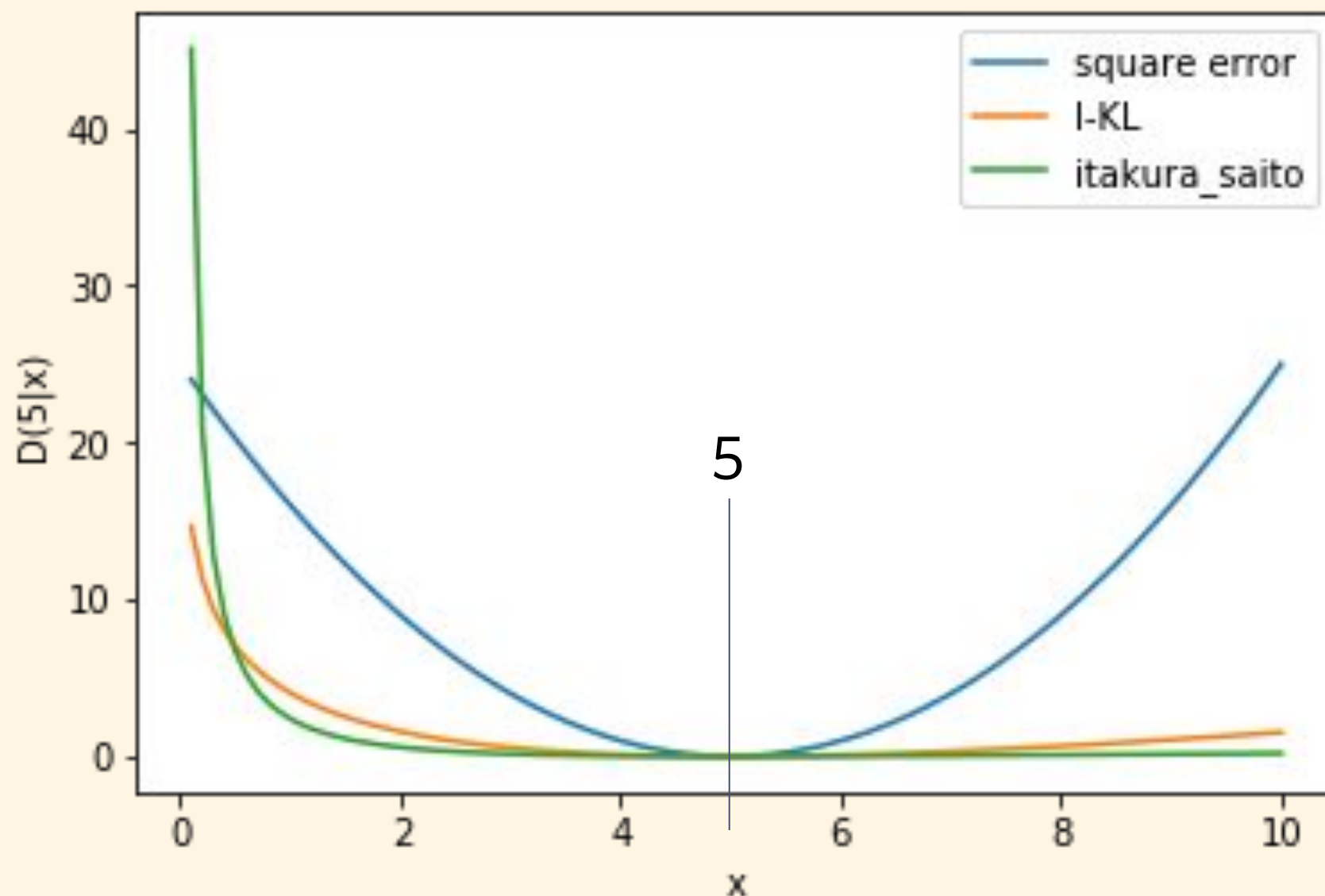
- ・初期値をランダムに設定し、複数回アルゴリズム回す

NMFのアルゴリズム

基底行列Hと係数行列Uを求める際の方針

→データ行列Yと行列積HUの乖離度を最小化する最適化問題

→乖離度：二乗誤差、一般化KLダイバージェンス、板倉斎藤距離



二乗誤差

$$D(y|x) = (y - x)^2$$

一般化KLダイバージェンス

$$D(y|x) = y \ln\left(\frac{y}{x}\right) - (y - x)$$

板倉 - 斎藤距離

$$D(y|x) = \frac{y}{x} - \ln\left(\frac{y}{x}\right) - 1$$

NMFのアルゴリズム

基底行列Hと係数行列Uを求める際の方針

→データ行列Yと行列積HUの乖離度を最小化する最適化問題

→乖離度: 二乗誤差、一般化KLダイバージェンス、板倉斎藤距離

→行列の要素計算: フロベニウスノルム

$$||A||_F = \sqrt{\sum_{i,j} a_{i,j}^2}$$

→フロベニウスノルムを使い、データ行列と行列積の差分を乖離度と仮定

$$\begin{aligned} D(\mathbf{Y}|\mathbf{HU}) &= ||\mathbf{Y} - \mathbf{HU}||_F^2 \\ &= \sum_{k,n} ||y_{k,n} - \sum_m h_{k,m} u_{m,n}||^2 \\ &= \sum_{k,n} \left(|y_{k,n}|^2 - 2y_{k,n} \sum_m h_{k,m} u_{m,n} - \left| \sum_m h_{k,m} u_{m,n} \right|^2 \right) \end{aligned}$$

閑話：イエンセンの不等式

$$\begin{aligned} D(\mathbf{Y}|\mathbf{H}\mathbf{U}) &= \|\mathbf{Y} - \mathbf{H}\mathbf{U}\|_F^2 \\ &= \sum_{k,n} \left\| y_{k,n} - \sum_m h_{k,m} u_{m,n} \right\|^2 \\ &= \sum_{k,n} \left(|y_{k,n}|^2 - 2y_{k,n} \sum_m h_{k,m} u_{m,n} - \left| \sum_m h_{k,m} u_{m,n} \right|^2 \right) \end{aligned}$$

和の絶対値を取り扱うことは難しい→**イエンセンの不等式**

$f(\cdot)$ が凸関数の時、

$$\sum_{i=1}^n \lambda_i f(x_i) \geq f\left(\sum_{i=1}^n \lambda_i x_i\right)$$

$$\begin{cases} \sum_i \lambda_i = 1 \\ \lambda_i \geq 0 \end{cases}$$

凸関数の性質

- 性質1：任意の $x_1, x_2, \lambda (0 \leq \lambda \leq 1)$ に対して、 $\lambda f(x_1) + (1 - \lambda)f(x_2) \geq f(\lambda x_1 + (1 - \lambda)x_2)$ を満たす。
- 性質2：任意の x_1, x_2 に対して、2点 $(x_1, f(x_1)), (x_2, f(x_2))$ を結ぶ線分が関数の上側にある。
- 性質3：二階微分 $f''(x)$ が存在して0以上である

NMFのアルゴリズム

$$\left| \sum_m h_{k,m} u_{m,n} \right|^2 \leq \sum_i \frac{h_{k,m}^2 u_{m,n}^2}{\lambda_{k,m,n}}$$

$$\begin{aligned} D(\mathbf{Y}|\mathbf{H}\mathbf{U}) &= \|\mathbf{Y} - \mathbf{H}\mathbf{U}\|_F^2 \\ &= \sum_{k,n} \left\| \mathbf{y}_{k,n} - \sum_m h_{k,m} u_{m,n} \right\|^2 \\ &= \sum_{k,n} \left(\|\mathbf{y}_{k,n}\|^2 - 2\mathbf{y}_{k,n} \sum_m h_{k,m} u_{m,n} + \underbrace{\left\| \sum_m h_{k,m} u_{m,n} \right\|^2}_{\leq \sum_m \frac{h_{k,m}^2 u_{m,n}^2}{\lambda_{k,m,n}}} \right) \\ &\leq \sum_{k,n} \left(\|\mathbf{y}_{k,n}\|^2 - 2\mathbf{y}_{k,n} \sum_m h_{k,m} u_{m,n} + \sum_m \frac{h_{k,m}^2 u_{m,n}^2}{\lambda_{k,m,n}} \right) \end{aligned}$$

変形手順

$$\left(\sum_i x_i \right)^2 = \left(\sum_i \lambda_i \frac{x_i}{\lambda_i} \right)^2 \leq \sum_i \lambda_i \left(\frac{x_i}{\lambda_i} \right)^2 = \sum_i \frac{x_i^2}{\lambda_i}$$

閑話：補助関数法

目的関数の上界を定める関数 G が求まった

→上界を定める関数 G を反復的に降下させることで解を求める

→反復降下の補助を行う関数(補助関数)として上界を定める関数 G を使う

→補助関数の反復降下による目的関数の降下方法を補助関数法

H と U の行列要素 $h_{k,1}, \dots, h_{k,M}$, $u_{1,n}, \dots, u_{M,n}$ を含んだ非線形関数項であることに気付く. この項に対し, 行列要素 $h_{k,m}$, $u_{m,n}$ ごとの関数の和に分離した形をした上限関数を設けたい. 2 次関数は凸関数であるため,

NMFのアルゴリズム

$$D(\mathbf{Y}|\mathbf{H}\mathbf{U}) \leq \sum_{k,n} \left(\|\mathbf{y}_{k,n}\|^2 - 2y_{k,n} \sum_m h_{k,m} u_{m,n} + \sum_m \frac{h_{k,m}^2 u_{m,n}^2}{\lambda_{k,m,n}} \right) \quad \mathbf{G}$$

→ 補助変数 $\lambda_{k,m,n}$ (イエンセンの不等式における係数)

→ 係数の条件より
$$\lambda_{k,m,n} = \frac{h_{k,m} u_{m,n}}{\sum_l h_{k,l} u_{l,n}}$$

→ 補助関数 \mathbf{G} を \mathbf{H} 、 \mathbf{U} に関して偏微分を行い停留点を探す

$$\frac{\partial G}{\partial h_{k,m}} = -2y_{k,n} \sum_m u_{m,n} + \sum_m 2 \frac{h_{k,m} u_{m,n}^2}{\lambda_{k,m,n}} = 0$$

$$\frac{\partial G}{\partial u_{m,n}} = -2y_{k,n} \sum_m h_{k,m} + \sum_m 2 \frac{h_{k,m}^2 u_{m,n}}{\lambda_{k,m,n}} = 0$$

NMFのアルゴリズム

$$\frac{\partial G}{\partial h_{k,m}} = -2y_{k,n} \sum_m u_{m,n} + \sum_m 2 \frac{h_{km} u_{m,n}^2}{\lambda_{k,m,n}} = 0$$
$$\frac{\partial G}{\partial u_{m,n}} = -2y_{k,n} \sum_m h_{k,m} + \sum_m 2 \frac{h_{k,m}^2 u_{m,n}}{\lambda_{k,m,n}} = 0$$

整理途中

参考文献

- ・非負値行列因子分解/亀岡弘和/計測と制御 第 51 巻 第 9 号 2012 年 9 月号
- ・チュートリアル: 非負値行列因子分解/亀岡弘和
- ・非負値行列因子分解/亀岡弘和/
<http://www.kecl.ntt.co.jp/people/kameoka.hirokazu/publications/Kameoka2012SICE09published.pdf>
- ・NMFアルゴリズムの導出/<https://r9y9.github.io/blog/2013/07/27/nmf-euclid/>
- ・NMFでMovieLensレコメンド/<https://ohke.hateblo.jp/entry/2017/12/26/230000>
- ・非負値行列因子分解を改めてやり直してみた/
https://qiita.com/sumita_v09/items/d22850f41257d07c45ea

ABOUT SANOGRAPHIX.NET

—

SANOGRAPHIX.NETは、佐野章核の制作物をただ並べただけの個人サイトです。

PROFILE (->[EN](#))

—

佐野章核 Showkaku Sano

グラフィックデザイナー。

2007年頃、他人の同人誌の装丁を請け負う活動を開始。その後、同人サークル「[jadda](#)」および「[kone1](#)」を立ち上げ、情報誌を不定期に発行しています。また、2012年あたりからTumblrテーマ制作を趣味としており、イラストポートフォリオ用の

「[Illustfolio](#)」、日記書きたい人用の「[ZEN](#)」「[Apollo](#)」などのTumblrテーマを今までに制作しました。休日は寺社仏閣巡りをしています。写真[このへん](#)で見れます。

<http://www.sanographix.net/>