

# 特徴抽出

カーネル主成分分析, 非負値行列因子分解, オートエンコーダ, t-SNE

中田和秀

東京科学大学 工学院 経営工学系

機械学習入門

<https://www.nakatalab.iee.e.titech.ac.jp/text/nakata.html>

## 概要

ここでは教師なし学習の一つである特徴抽出について説明する。特徴抽出によって得られた低次元ベクトルは、教師あり学習の特徴量として使ったり、可視化して特徴を理解することなどに利用できる。

目次：

1. カーネル主成分分析
  - 1.1 主成分分析
  - 1.2 カーネル主成分分析
2. 非負値行列因子分解 (MNF)
3. オートエンコーダー
4. t-SNE
  - 4.1 SNE (Stochastic Neighbor Embedding)
  - 4.2 t-SNE (t-Distributed Stochastic Neighbor Embedding)

記号の使い方：

- $A := B$  は、 $B$  で  $A$  を定義する、 $B$  を  $A$  に代入することを意味する
- $[n]$  は  $n$  までのインデックスの集合を表し  $[n] := \{1, 2, \dots, n\}$

# 特徴抽出

データ：  $\{\mathbf{x}_d\}_{d \in [D]}$ ,  $\mathbf{x}_d \in \mathbb{R}^n$

目的：  $n$  次元のベクトルを  $P$  次元 ( $P \ll n$ ) のベクトルで表現

- 無駄な情報（ノイズ？）を削り、本質的な情報を取り出すため
- 各データ点を理解するため

見方の違いによって様々な名前で呼ばれる。

- 特徴抽出： 重要な特徴を低次元のベクトルで表現する
- 次元圧縮： 高次元ベクトルを低次元ベクトルに変換する
- 埋め込み： 各データを低次元な線形空間の中へ埋め込む
- 非可逆データ圧縮： 元の情報を低次元なベクトルで出来るだけ保持する
- 可視化： 2次元（3次元）のベクトルに変換して図示する

ここでは、以下の手法を紹介する

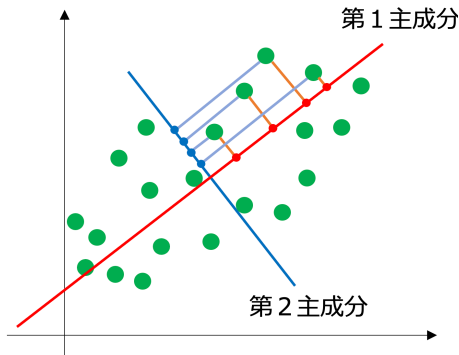
- カーネル主成分分析
- 非負値行列因子分解 (NMF)
- オートエンコーダー
- t-SNE

# 主成分分析

- $\mathbf{x} = x_1, x_2, \dots, x_D$  に対して線形変換をして、一つの実数で特徴付け

$$s := \mathbf{w}^T \mathbf{x} \in \mathbb{R}$$

- アイデア：新しい特徴量  $s$  が各データ点の違いを表現するためには、 $s$  の分散が大きくなほうがよい（情報量が大きい）。



# 分散共分散行列

$$\bar{\mathbf{x}} := \frac{1}{D} \sum_{d \in [D]} \mathbf{x}_d, \quad \overline{\mathbf{x}}_d := \mathbf{x}_d - \bar{\mathbf{x}}, \quad \mathbf{X}_0 := \begin{pmatrix} \overline{\mathbf{x}}_1^T \\ \overline{\mathbf{x}}_2^T \\ \vdots \\ \overline{\mathbf{x}}_D^T \end{pmatrix} \in \mathbb{R}^{D \times n}$$

とする。(標本) 分散共分散行列は

$$\mathbb{V}[\mathbf{x}] := \frac{1}{D} \mathbf{X}_0^T \mathbf{X}_0 \in \mathbb{R}^{n \times n}$$

$\Sigma := \mathbb{V}[\mathbf{x}] = \frac{1}{D} \mathbf{X}_0^T \mathbf{X}_0$  とすると、

$$\mathbb{V}[\mathbf{s}] = \mathbb{V}[\mathbf{w}^T \mathbf{x}] = \mathbf{w}^T \mathbb{V}[\mathbf{x}] \mathbf{w} = \mathbf{w}^T \Sigma \mathbf{w}$$

# 第一主成分

次の最適化問題を解く

$$\begin{array}{ll}\max_{\boldsymbol{w}} & \boldsymbol{w}^T \boldsymbol{\Sigma} \boldsymbol{w} \\ \text{s.t.} & \boldsymbol{w}^T \boldsymbol{w} = 1.\end{array}$$

この問題の最適解は  $\boldsymbol{\Sigma}$  の最大固有値に対応する固有ベクトル

→ 説明は次スライド

- 最適解  $\boldsymbol{w}_1^*$  を使って計算された特徴量  $s_1 := \boldsymbol{w}_1^{*T} \bar{\boldsymbol{x}}$  を第1主成分得点という。
- 全てのデータ  $\{\boldsymbol{x}_d\}_{d \in [D]}$  に対する第1主成分得点を並べたベクトル

$$\boldsymbol{s}_1 := \begin{pmatrix} \boldsymbol{w}_1^{*T} \bar{\boldsymbol{x}}_1 \\ \boldsymbol{w}_1^{*T} \bar{\boldsymbol{x}}_2 \\ \vdots \\ \boldsymbol{w}_1^{*T} \bar{\boldsymbol{x}}_D \end{pmatrix} = \begin{pmatrix} \bar{\boldsymbol{x}}_1^T \boldsymbol{w}_1^* \\ \bar{\boldsymbol{x}}_2^T \boldsymbol{w}_1^* \\ \vdots \\ \bar{\boldsymbol{x}}_D^T \boldsymbol{w}_1^* \end{pmatrix} = \boldsymbol{X}_0 \boldsymbol{w}_1^* \in \mathbb{R}^D$$

# 説明

等式制約付き非線形最適化問題なのでラグランジュの未定乗数法で解く。  
ラグランジュ関数は、ラグランジュ乗数  $\alpha$  を導入して

$$L(\boldsymbol{w}, \alpha) := \boldsymbol{w}^T \boldsymbol{\Sigma} \boldsymbol{w} + \alpha(1 - \boldsymbol{w}^T \boldsymbol{w})$$

となる。最適解であるための必要条件は、

$$\nabla_{\boldsymbol{w}} L(\boldsymbol{w}, \alpha) = 2\boldsymbol{\Sigma} \boldsymbol{w} - 2\alpha \boldsymbol{w} = \mathbf{0},$$

$$\nabla_{\alpha} L(\boldsymbol{w}, \alpha) = 1 - \boldsymbol{w}^T \boldsymbol{w} = 0$$

1つ目の式より、

$$\boldsymbol{\Sigma} \boldsymbol{w} = \alpha \boldsymbol{w}$$

つまり、最適解の候補  $\boldsymbol{w}$  は  $\boldsymbol{\Sigma}$  の固有ベクトルであることを示している。固有値  $\lambda$  に対する固有ベクトルを  $\boldsymbol{w}$  とすると、目的関数値は

$$\boldsymbol{w}^T \boldsymbol{\Sigma} \boldsymbol{w} = \boldsymbol{w}^T (\lambda \boldsymbol{w}) = \lambda \boldsymbol{w}^T \boldsymbol{w} = \lambda$$

よって、最大固有値に対する固有ベクトルが最適解になる。

# 他の主成分

## 第二主成分

- 第一主成分の影響を取り除いた上で、分散が最大となる線形変換を考える
- 第二主成分  $w_2^*$  は  $\Sigma$  の 2 番目に大きな固有値に対する固有ベクトル  
→ 詳細は次ページ
- 第二主成分得点を並べたものは  $s_2 := X_0 w_2^* \in \mathbb{R}^D$

## 第 $p$ 主成分 ( $p \in [P]$ )

- 第  $p-1$  主成分までの影響を取り除いた上、分散が最大となる線形変換を考える
- 第  $p$  主成分は  $\Sigma$  の  $p$  番目に大きな固有値に対応する固有ベクトル  $w_p^*$
- 第  $p$  主成分得点を並べたものは  $s_p := X_0 w_p^* \in \mathbb{R}^D$

$x_d$  の次元圧縮をして得られた新たな特徴ベクトルは

$$([s_1]_d, [s_2]_d, \dots, [s_P]_d)^T \in \mathbb{R}^P$$



## 第二主成分の計算

- $\overline{x_d}$  に対して、既に特徴量化した成分を引いておく（直交射影）。

$$\overline{x_d} - (w_1^T \overline{x_d}) w_1$$

- 2 番目の特徴量  $s_2 := w_2^T x$  を作るための重み  $w_2$  を考える。

$$\begin{aligned}\overline{x_d} &\rightarrow \overline{x_d} - (w_1^T \overline{x_d}) w_1 \\ X_0^T &\rightarrow X_0^T - w_1 w_1^T X_0^T \\ \Sigma &\rightarrow \Sigma - \lambda_1 w_1 w_1^T\end{aligned}$$

式変形は次スライド

最適化問題は次のようになる

$$\begin{aligned}\max_{\boldsymbol{w}} \quad & \boldsymbol{w}^T (\Sigma - \lambda_1 w_1 w_1^T) \boldsymbol{w} \\ \text{s.t.} \quad & \boldsymbol{w}^T \boldsymbol{w} = 1.\end{aligned}$$

この問題の最適解は  $\Sigma$  の 2 番目に大きな固有値に対する固有ベクトルである。

# 式変形

平均

$$\begin{aligned}\frac{1}{D} \sum_{d \in [D]} \bar{x}_d - (\mathbf{w}_1^T \bar{x}_d) \mathbf{w}_1 &= \frac{1}{D} \sum_{d \in [D]} \bar{x}_d - \mathbf{w}_1 \mathbf{w}_1^T \bar{x}_d \\ &= \frac{1}{D} (\mathbf{I} - \mathbf{w}_1 \mathbf{w}_1^T) \sum_{d \in [D]} \bar{x}_d = \mathbf{0}\end{aligned}$$

分散共分散行列

$$\begin{aligned}& \frac{1}{D} (\mathbf{X}_0^T - \mathbf{w}_1 \mathbf{w}_1^T \mathbf{X}_0^T) (\mathbf{X}_0 - \mathbf{X}_0 \mathbf{w}_1 \mathbf{w}_1^T) \\ &= \frac{1}{D} \left( \mathbf{X}_0^T \mathbf{X}_0 - \mathbf{w}_1 \mathbf{w}_1^T \mathbf{X}_0^T \mathbf{X}_0 - \mathbf{X}_0^T \mathbf{X}_0 \mathbf{w}_1 \mathbf{w}_1^T + \mathbf{w}_1 \mathbf{w}_1^T \mathbf{X}_0^T \mathbf{X}_0 \mathbf{w}_1 \mathbf{w}_1^T \right) \\ &= \frac{1}{D} \mathbf{X}_0^T \mathbf{X}_0 - \lambda_1 \mathbf{w}_1 \mathbf{w}_1^T - \lambda_1 \mathbf{w}_1 \mathbf{w}_1^T + \lambda_1 \mathbf{w}_1 \mathbf{w}_1^T \\ &= \frac{1}{D} \mathbf{X}_0^T \mathbf{X}_0 - \lambda_1 \mathbf{w}_1 \mathbf{w}_1^T = \Sigma - \lambda_1 \mathbf{w}_1 \mathbf{w}_1^T\end{aligned}$$

固有値分解

$$\Sigma = \sum_{i=1}^n \lambda_i \mathbf{w}_i \mathbf{w}_i^T \quad \text{のとき、} \quad \Sigma - \lambda_1 \mathbf{w}_1 \mathbf{w}_1^T = \sum_{i=2}^n \lambda_i \mathbf{w}_i \mathbf{w}_i^T$$

# 寄与率

各特徴量の分散の和は固有値の和と一致する。

$$\sum_{i \in [n]} \mathbb{V}[x_i] = \sum_{i \in [n]} \Sigma_{ii} = \text{Tr}(\Sigma) = \sum_{i \in [n]} \lambda_i$$

分散全体に対する主成分得点  $s_p$  の分散の割合を寄与率という。

$$\text{寄与率: } c_p := \frac{\lambda_p}{\sum_{i \in [n]} \lambda_i} = \frac{\lambda_p}{\sum_{i \in [n]} \Sigma_{ii}} \quad (p \in [P])$$

$$\text{累積寄与率: } r_p := \sum_{i \in [p]} c_i = \frac{\sum_{i \in [p]} \lambda_i}{\sum_{i \in [n]} \lambda_i} = \frac{\sum_{i \in [p]} \lambda_i}{\sum_{i \in [n]} \Sigma_{ii}} \quad (p \in [P])$$

$$0 \leq c_P \leq c_{P-1} \leq \cdots \leq c_1 = r_1 \leq r_2 \leq \cdots r_P \leq 1$$

累積寄与率によって、どの程度の情報を抽出できているか判断する

# 可視化

- 1 番目のスコアを x 軸、2 番目のスコアを y 軸にとると可視化が可能になる。

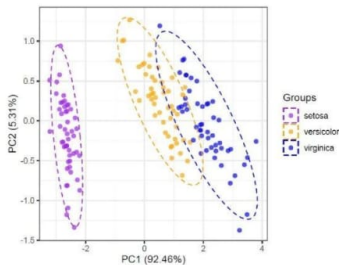
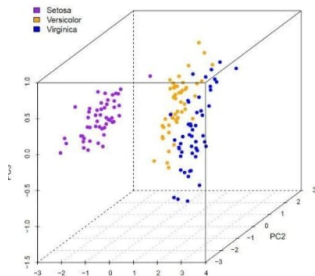


Figure: <https://www.binarydevelop.com/article/rpca-28546>

# カーネル主成分分析

- 主成分分析では、線形変換によって新しい特徴量（主成分得点）を計算
- 非線形な変換を行いたい
- 入力空間を高次元の特徴空間に写像して、その特徴空間上で主成分分析を行う

$\phi: \mathbb{R}^n \rightarrow V$  として、 $\mathbf{x}_d$  の代わりに  $\phi(\mathbf{x}_d)$  を使う<sup>1</sup>。

$$\overline{\phi_{\mathbf{x}}} := \frac{1}{D} \sum_{d \in [D]} \phi(\mathbf{x}_d), \quad \overline{\phi(\mathbf{x}_d)} := \phi(\mathbf{x}_d) - \overline{\phi_{\mathbf{x}}}, \quad \widetilde{\mathbf{X}}_0 := \begin{pmatrix} \overline{\phi(\mathbf{x}_1)}^T \\ \overline{\phi(\mathbf{x}_2)}^T \\ \vdots \\ \overline{\phi(\mathbf{x}_D)}^T \end{pmatrix} \in \mathbb{R}^{D \times |V|}$$

分散共分散行列

$$\Sigma := \frac{1}{D} \widetilde{\mathbf{X}}_0^T \widetilde{\mathbf{X}}_0 \in \mathbb{R}^{|V| \times |V|}$$

---

<sup>1</sup>写像先の高次元空間  $V$  の次元  $|V|$  は無限かもしれないが、有限のイメージで話を進める

## カーネルトリック<sup>2</sup>

$\frac{1}{D} \widetilde{\mathbf{X}}_0^T \widetilde{\mathbf{X}}_0 \in \mathbb{R}^{|V| \times |V|}$  の固有値分解が必要

- $\phi(x)$  を使わず、 $\phi(x_i)^T \phi(x_j)$  に相当する  $\text{Ker}(x_i, x_j)$  を使いたい
- つまり、 $\mathbf{K} := \widetilde{\mathbf{X}} \widetilde{\mathbf{X}}^T \in \mathbb{R}^{D \times D}$  を使って計算する

$$\widetilde{\mathbf{X}} := \begin{pmatrix} \phi(x_1)^T \\ \phi(x_2)^T \\ \vdots \\ \phi(x_D)^T \end{pmatrix} \in \mathbb{R}^{D \times |V|}$$

### 問題点

- $\widetilde{\mathbf{X}}$  でなく  $\widetilde{\mathbf{X}}_0$  を使った計算が必要（平均を  $\mathbf{0}$  にする平行移動がある）
- $\widetilde{\mathbf{X}}_0^T \widetilde{\mathbf{X}}_0 \in \mathbb{R}^{|V| \times |V|}$  だが、計算できるのはカーネル行列  $\mathbf{K} := \widetilde{\mathbf{X}} \widetilde{\mathbf{X}}^T \in \mathbb{R}^{D \times D}$

---

<sup>2</sup>カーネルトリックの詳細については、カーネル法のスライドを参照されたい

# 修正カーネル行列

計算可能なカーネル行列  $K := \widetilde{X} \widetilde{X}^T$  から修正カーネル行列  $K_0 := \widetilde{X}_0 \widetilde{X}_0^T$  を作る。

$$\text{平均ベクトル} \quad \overline{\phi_{\mathbf{x}}} := \frac{1}{D} \sum_{d \in [D]} \phi(\mathbf{x}_d) = \frac{1}{D} \widetilde{X}^T \mathbf{e}$$

であるので、

$$\widetilde{X}_0 := \begin{pmatrix} \phi(\mathbf{x}_1)^T - \overline{\phi_{\mathbf{x}}}^T \\ \phi(\mathbf{x}_2)^T - \overline{\phi_{\mathbf{x}}}^T \\ \vdots \\ \phi(\mathbf{x}_D)^T - \overline{\phi_{\mathbf{x}}}^T \end{pmatrix} = \widetilde{X} - \mathbf{e} \overline{\phi_{\mathbf{x}}}^T = \widetilde{X} - \frac{1}{D} \mathbf{e} \mathbf{e}^T \widetilde{X}$$

という関係が成り立つ。このとき、

$$K_0 = K - \frac{1}{D} \mathbf{e} \mathbf{e}^T K - \frac{1}{D} K \mathbf{e} \mathbf{e}^T + \frac{\mathbf{e}^T K \mathbf{e}}{D^2} \mathbf{e} \mathbf{e}^T$$

# 固有値の計算

$\widetilde{\mathbf{X}}_0^T \widetilde{\mathbf{X}}_0$  の固有値  $\lambda$  と固有ベクトル  $\mathbf{w} \in \mathbb{R}^{|V|}$  は次の関係を満たす。

$$\widetilde{\mathbf{X}}_0^T \widetilde{\mathbf{X}}_0 \mathbf{w} = \lambda \mathbf{w}$$

このとき、次の関係が成り立つ。

$$\exists \mathbf{a} \in \mathbb{R}^D, \quad \mathbf{w} = \widetilde{\mathbf{X}}_0^T \mathbf{a}$$

これを代入すると、

$$\widetilde{\mathbf{X}}_0^T \widetilde{\mathbf{X}}_0 \widetilde{\mathbf{X}}_0^T \mathbf{a} = \lambda \widetilde{\mathbf{X}}_0^T \mathbf{a}$$

左から  $\widetilde{\mathbf{X}}_0$  をかけると

$$\widetilde{\mathbf{X}}_0 \widetilde{\mathbf{X}}_0^T \widetilde{\mathbf{X}}_0 \widetilde{\mathbf{X}}_0^T \mathbf{a} = \lambda \widetilde{\mathbf{X}}_0 \widetilde{\mathbf{X}}_0^T \mathbf{a}$$

これは、 $\mathbf{K}_0^2 \mathbf{a} = \lambda \mathbf{K}_0 \mathbf{a}$  なので、次の関係を満たせば上の式も成り立つ。

$$\mathbf{K}_0 \mathbf{a} = \lambda \mathbf{a}$$

この固有値問題を解き、求まった  $\mathbf{a}$  に対し  $\widetilde{\mathbf{X}}_0^T \mathbf{a}$  が必要な固有ベクトル  $\mathbf{w}$  となる。



# 主成分得点の計算

既に固有ベクトル  $\mathbf{a}_p$  ( $p \in [P]$ ) は計算されているとする。

第  $p$  主成分得点

$$s_p := \mathbf{w}_p^{*T} \overline{\phi(\mathbf{x})} = \overline{\phi(\mathbf{x})}^T \mathbf{w}_p^* = \overline{\phi(\mathbf{x})}^T \widetilde{\mathbf{X}}_0^T \mathbf{a}_p$$

全データ  $\{\mathbf{x}_d\}_{d \in [D]}$  に対する得点を縦に並べると

$$\begin{aligned} \mathbf{s}_p &:= \begin{pmatrix} \overline{\phi(\mathbf{x}_1)}^T \\ \overline{\phi(\mathbf{x}_2)}^T \\ \vdots \\ \overline{\phi(\mathbf{x}_D)}^T \end{pmatrix} \widetilde{\mathbf{X}}_0^T \mathbf{a}_p \\ &= \widetilde{\mathbf{X}}_0 \widetilde{\mathbf{X}}_0^T \mathbf{a}_p \\ &= \mathbf{K}_0 \mathbf{a}_p \\ &= \lambda_p \mathbf{a}_p \quad \in \mathbb{R}^D \end{aligned}$$

# 学習アルゴリズム

## 学習アルゴリズム

ステップ1  $K_{ij} := \text{Ker}(\mathbf{x}_i, \mathbf{x}_j) \quad i, j \in [D]$

ステップ2  $\mathbf{K}_0 := \mathbf{K} - \frac{1}{D} \mathbf{e} \mathbf{e}^T \mathbf{K} - \frac{1}{D} \mathbf{K} \mathbf{e} \mathbf{e}^T + \frac{\mathbf{e}^T \mathbf{K} \mathbf{e}}{D^2} \mathbf{e} \mathbf{e}^T$

ステップ3  $\mathbf{K}_0$  の固有値が大きな  $p$  個の固有ベクトルの計算。  $\mathbf{a}_p \quad (p \in [P])$

ステップ4  $p \in [P]$  に対して、第  $p$  主成分得点の計算

$$\mathbf{s}_p := \lambda_p \mathbf{a}_p \in \mathbb{R}^D$$

$\phi(\mathbf{x}_d)$  の次元圧縮をして得られた新たな特徴ベクトルは

$$([\mathbf{s}_1]_d, [\mathbf{s}_2]_d, \dots, [\mathbf{s}_P]_d)^T \in \mathbb{R}^P$$

- $\{\phi(\mathbf{x}_d)\}_{d \in [D]}$  に対する線形変換
- $\{\mathbf{x}_d\}_{d \in [D]}$  に対する非線形変換

# 計算例

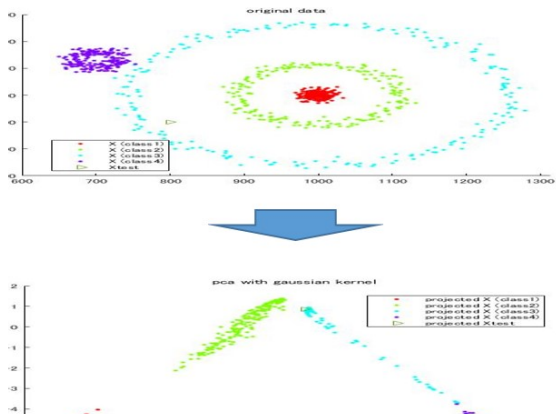


Figure: [https://jp.mathworks.com/matlabcentral/fileexchange/71647-matlab-kernel-pca?s\\_tid=prof\\_contriblnk](https://jp.mathworks.com/matlabcentral/fileexchange/71647-matlab-kernel-pca?s_tid=prof_contriblnk)

# カーネル主成分分析のまとめ

## 特徴

- データ数  $D$  の大きさの行列の固有値分解が必要

計算量： $O(D^3)$

計算量のオーダーは減らすこともできるが、色々難しい・・・

- データ数が増えると固有ベクトルの計算が困難
- カーネルを導入すると、特徴量（軸）の解釈が困難になる
- 主成分の数  $P$  は累積寄与率をみて決めることも可能

# 非負値行列因子分解 (NMF)

データ  $T \in \mathbb{R}^{I \times J}$

$T_{ij} \geq 0$  : ユーザ  $i$  のアイテム  $j$  に対する評価値

目的

各ユーザと各アイテムの特徴量を構築したい

モデル:

- 評価値  $T_{ij}$  は、 $K$  個の要因の和で構成される
- ユーザやアイテムは、各要因に対する感度を持つ
  - ▶  $\mathbf{x}_i \in \mathbb{R}^K$  : ユーザ  $i$  の各要因に対する感度
  - ▶  $\mathbf{y}_j \in \mathbb{R}^K$  : アイテム  $j$  の各要因に対する感度
- $T_{ij} \simeq \mathbf{x}_i^T \mathbf{y}_j = \sum_{k \in [K]} [\mathbf{x}_i]_k [\mathbf{y}_j]_k$

$$\mathbf{x}_i \geq \mathbf{0}$$

$$\mathbf{y}_j \geq \mathbf{0}$$

$\mathbf{x}_i, \mathbf{y}_j$  をユーザやアイテムの特徴量とする

# 最適化問題

2乗誤差  $(T_{ij} - \mathbf{x}_i^T \mathbf{y}_j)^2$  が小さくなるように特徴ベクトル  $\mathbf{x}_i, \mathbf{y}_j$  を決める

## 非負値行列因子分解

変数:  $\mathbf{x}_i \in \mathbb{R}^K$  ( $i \in [I]$ ),  $\mathbf{y}_j \in \mathbb{R}^K$  ( $j \in [J]$ )

$$\begin{aligned} \min \quad & \sum_{i \in [I]} \sum_{j \in [J]} (T_{ij} - \mathbf{x}_i^T \mathbf{y}_j)^2 \\ \text{s.t.} \quad & \mathbf{x}_i \geq \mathbf{0} \quad (i \in [I]), \\ & \mathbf{y}_j \geq \mathbf{0} \quad (j \in [J]). \end{aligned}$$

$$\mathbf{X} := \begin{pmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_I^T \end{pmatrix} \in \mathbb{R}^{I \times K}, \quad \mathbf{Y} := \begin{pmatrix} \mathbf{y}_1 & \mathbf{y}_2 & \cdots & \mathbf{y}_J \end{pmatrix} \in \mathbb{R}^{K \times J}$$

とする

# 行列表現

## 非負値行列因子分解

変数:  $\mathbf{X} \in \mathbb{R}^{I \times K}$ ,  $\mathbf{Y} \in \mathbb{R}^{K \times J}$

$$\begin{aligned} \min \quad & \|\mathbf{T} - \mathbf{XY}\|_F^2 \\ \text{s.t.} \quad & \mathbf{X} \geq \mathbf{O}, \mathbf{Y} \geq \mathbf{O}. \end{aligned}$$

フロベニウスノルム  $\|\mathbf{A}\|_F := \sqrt{\text{Tr}(\mathbf{A}^T \mathbf{A})} = \sqrt{\text{Tr}(\mathbf{A} \mathbf{A}^T)} = \sqrt{\sum_{i \in [I]} \sum_{j \in [J]} (A_{ij})^2}$

$$\mathbf{T} \simeq \mathbf{XY}$$

- 行列を低ランク行列で近似

$$\text{rank}(\mathbf{T}) = \min\{I, J\}, \quad \text{rank}(\mathbf{XY}) = K \quad (I, J \gg K)$$

- 行列を2つの行列に分解

非負値行列因子分解 (NMF: Nonnegative Matrix Factorization) と呼ぶ

# 行列因子分解 (MF)

まず、非負条件の無い行列因子分解は簡単であることをみておく。

## 行列因子分解

変数:  $\mathbf{X} \in \mathbb{R}^{I \times K}$ ,  $\mathbf{Y} \in \mathbb{R}^{K \times J}$

$$\min \|\mathbf{T} - \mathbf{XY}\|_F^2$$

$\mathbf{T}$  の特異値分解: 特異値は降順に  $\Sigma$  の左上から対角部分に並んでいるとする

$$\mathbf{T} = \mathbf{P}\Sigma\mathbf{Q}^T$$

- $\Sigma_K \in \mathbb{R}^{K \times K}$ :  $\Sigma$  の左上  $K \times K$  を抜き出した対角行列
- $\mathbf{P}_K \in \mathbb{R}^{I \times K}$ :  $\mathbf{P}$  の左側  $K$  列を抜き出した行列
- $\mathbf{Q}_K \in \mathbb{R}^{J \times K}$ :  $\mathbf{Q}$  の左側  $K$  列を抜き出した行列

このとき、最適解  $\mathbf{X}^*, \mathbf{Y}^*$  は次のように表せる。

$$\mathbf{X}^* := \mathbf{P}_K \Sigma_K^{1/2} \mathbf{H}, \quad \mathbf{Y}^* := \mathbf{H}^{-1} \Sigma_K^{1/2} \mathbf{Q}_K^T$$

ただし、 $\mathbf{H} \in \mathbb{R}^{K \times K}$  は任意の正則行列



## Eckart–Young の定理

$Z$  を変数とする次の最適化問題

$$\begin{aligned} \min \quad & \|T - Z\|_F \\ \text{s.t.} \quad & \text{rank}(Z) \leq K. \end{aligned}$$

の最適解は  $Z^* = P_K \Sigma_K Q_K^T$  ( $P_K, \Sigma_K, Q_K$  は前ページで定義したもの)

証明は [https://en.wikipedia.org/wiki/Low-rank\\_approximation](https://en.wikipedia.org/wiki/Low-rank_approximation)

この定理より、 $X^* Y^* = P_K \Sigma_K Q_K^T$  と分かる。このとき、最適解  $X^*, Y^*$  は

$$X^* := P_K \Sigma_K^{1/2} H, \quad Y^* := H^{-1} \Sigma_K^{1/2} Q_K^T$$

と表せる。ただし、 $H \in \mathbb{R}^{K \times K}$  は任意の正則行列

- この形式で表された  $X^*$  と  $Y^*$  が最適解であるのは明らか。
- 全ての最適解がこの形式で表されることも、特異値を 0 と 0 以外のもので分ければ証明できる。

$H$  の取り方に任意性があるため、最適解が一意に決まらないのが欠点

# 非負値行列因子分解

## 最適化問題

$$\begin{aligned} \min_{\mathbf{X}, \mathbf{Y}} \quad & \|\mathbf{T} - \mathbf{XY}\|_F^2 \\ \text{s.t.} \quad & \mathbf{X} \geq \mathbf{O}, \mathbf{Y} \geq \mathbf{O}. \end{aligned}$$

$\mathbf{X}$  と  $\mathbf{Y}$  を交互に最適化することを考える。

$$\begin{aligned} \min_{\mathbf{X}} \quad & \|\mathbf{T} - \mathbf{XY}\|_F^2 \\ \text{s.t.} \quad & \mathbf{X} \geq \mathbf{O}. \end{aligned}$$

$$\begin{aligned} \min_{\mathbf{Y}} \quad & \|\mathbf{T} - \mathbf{XY}\|_F^2 \\ \text{s.t.} \quad & \mathbf{Y} \geq \mathbf{O}. \end{aligned}$$

目的関数は  $\mathbf{X}, \mathbf{Y}$  それぞれの変数に関する凸 2 次関数であるが、非負制約があるため、最適解を求めることは難しい。

→ 補助関数法を使う。

# 補助関数 1

$$F(\mathbf{X}, \mathbf{Y}) := \|\mathbf{T} - \mathbf{X}\mathbf{Y}\|_F^2 = \sum_{i \in [I]} \sum_{j \in [J]} \left( T_{ij} - \sum_{k \in [K]} X_{ik} Y_{kj} \right)^2$$

$i \in [I], j \in [J]$  に対し、

$$\left( T_{ij} - \sum_{k \in [K]} X_{ik} Y_{kj} \right)^2 = \left( \sum_{k \in [K]} X_{ik} Y_{kj} \right)^2 - 2 T_{ij} \sum_{k \in [K]} X_{ik} Y_{kj} + T_{ij}^2$$

$\left( \sum_{k \in [K]} X_{ik} Y_{kj} \right)^2$  は、補助関数法の  $f_0 \left( \sum_{k \in [K]} f_k(\mathbf{X}, \mathbf{Y}) \right)$  というケース

- $f_0(x) := x^2$  は凸関数
- $f_k(\mathbf{X}, \mathbf{Y}) := X_{ik} Y_{kj} \geq 0$

## 補助関数 2

任意の  $\sum_{k \in [K]} r_{ijk} = 1, r_{ijk} > 0$  に対し、

$$\left( \sum_{k \in [K]} X_{ik} Y_{kj} \right)^2 = \left( \sum_{k \in [K]} r_{ijk} \frac{X_{ik} Y_{kj}}{r_{ijk}} \right)^2 \leq \sum_{k \in [K]} r_{ijk} \left( \frac{X_{ik} Y_{kj}}{r_{ijk}} \right)^2 = \sum_{k \in [K]} \frac{X_{ik}^2 Y_{kj}^2}{r_{ijk}}$$

そして、 $r_{ijk}^* = \frac{X_{ik} Y_{kj}}{\sum_{k' \in [K]} X_{ik'} Y_{k'j}}$  のとき等号が成り立つ。

## 補助関数

$$G(\mathbf{X}, \mathbf{Y}, \mathbf{R}) := \sum_{i \in [I]} \sum_{j \in [J]} \left( \sum_{k \in [K]} \frac{X_{ik}^2 Y_{kj}^2}{r_{ijk}} - 2 T_{ij} \sum_{k \in [K]} X_{ik} Y_{kj} + T_{ij}^2 \right)$$

は、 $F(\mathbf{X}, \mathbf{Y})$  の補助関数

( $r_{ijk}$  は 3 つの添字からなる変数だが、まとめて行列風の  $\mathbf{R}$  という記号を使う)

# 補助関数法

このとき、補助関数法のステップ2で解く最適化問題は

$$\begin{aligned} \min_{\mathbf{X}, \mathbf{Y}} \quad & G(\mathbf{X}, \mathbf{Y}, \mathbf{R}) \\ \text{s.t.} \quad & \mathbf{X} \geq \mathbf{O}, \mathbf{Y} \geq \mathbf{O}. \end{aligned}$$

である。これを  $\mathbf{X}$  と  $\mathbf{Y}$  の問題に分割して順に解く。

$$\begin{aligned} \min_{\mathbf{X}} \quad & G(\mathbf{X}, \mathbf{Y}, \mathbf{R}) \\ \text{s.t.} \quad & \mathbf{X} \geq \mathbf{O}. \end{aligned}$$

$$\begin{aligned} \min_{\mathbf{Y}} \quad & G(\mathbf{X}, \mathbf{Y}, \mathbf{R}) \\ \text{s.t.} \quad & \mathbf{Y} \geq \mathbf{O}. \end{aligned}$$

これらの問題の最適解は  $\hat{\mathbf{T}} := \mathbf{X}\mathbf{Y}$  としたとき、次のように計算できる<sup>3</sup>。

$$X_{ik}^* = X_{ik} \frac{\sum_{j \in [J]} T_{ij} Y_{kj}}{\sum_{j \in [J]} \hat{T}_{ij} Y_{kj}}$$

$$Y_{kj}^* = \frac{\sum_{i \in [I]} X_{ik} T_{ij}}{\sum_{i \in [I]} X_{ik} \hat{T}_{ij}}$$

証明は2ページ後

---

<sup>3</sup> $\mathbf{R}$  は  $\mathbf{X}, \mathbf{Y}$  で置き換えているので消えている

# アルゴリズム

## 非負値行列因子分解のアルゴリズム

ステップ0 初期解  $X, Y > O$  を適当に決める。

$X$  の最適化

ステップ1  $\hat{T} := XY$

ステップ2  $A := TY^T, \hat{A} := \hat{T}Y^T$

ステップ3  $X_{ik} := X_{ik} \frac{A_{ik}}{\hat{A}_{ik}} \quad (i \in [I], k \in [K])$

$Y$  の最適化

ステップ4  $\hat{T} := XY$  (更新した  $X$  を使うので、ステップ1の  $\hat{T}$  とは異なる)

ステップ5  $B := X^T T, \hat{B} := X^T \hat{T}$

ステップ6  $Y_{kj} := Y_{kj} \frac{B_{kj}}{\hat{B}_{kj}} \quad (k \in [K], j \in [J])$

ステップ7 終了条件を満たさなければ、ステップ1に戻る。

※ 局所最適解が求める。初期点を替えて何度か実行するとよい。

# 証明 1

$X$  に関する最適化問題を考える。

非負条件  $X \geq O$  を無視すると、目的関数は  $X$  に関する凸 2 次関数。

$$\begin{aligned} G(X, Y, R) &:= \sum_{i,j} \left( \sum_k \frac{X_{ik}^2 Y_{kj}^2}{r_{ijk}} - 2T_{ij} \sum_k X_{ik} Y_{kj} + T_{ij}^2 \right) \\ &= \sum_{i,k} \left( \left( \sum_j \frac{Y_{kj}^2}{r_{ijk}} \right) X_{ik}^2 - 2 \left( \sum_j T_{ij} Y_{kj} \right) X_{ik} \right) + \sum_{i,j} T_{ij}^2 \end{aligned}$$

$X_{ik}$  ごとに、2 次関数の最小化を考えればよいので、最適解は次のようになる。

$$X_{ik}^* = \frac{\sum_j T_{ij} Y_{kj}}{\sum_j \frac{Y_{kj}^2}{r_{ijk}}}$$

$$r_{ijk} := \frac{X_{ik} Y_{kj}}{\sum_{k'} X_{ik'} Y_{k'j}} \text{ を代入する }^4$$

---

<sup>4</sup> $r_{ijk}$  は定数であり ( $X, Y$  は現在の反復点)、最適化問題の変数ではないことに注意

## 証明 2

$$\begin{aligned} X_{ik}^* &= \frac{\sum_j T_{ij} Y_{kj}}{\sum_j \frac{Y_{kj}}{X_{ik}} \sum_{k'} X_{ik'} Y_{k'j}} = X_{ik} \frac{\sum_j T_{ij} Y_{kj}}{\sum_j Y_{kj} \sum_{k'} X_{ik'} Y_{k'j}} \\ &= X_{ik} \frac{\sum_j T_{ij} Y_{kj}}{\sum_j \hat{T}_{ij} Y_{kj}} \quad (\text{ただし } \hat{T} := XY) \end{aligned}$$

- $X, Y, T \geq O$  のとき、 $X^* \geq 0$  となり、自動的に非負条件制約は満たされる。よって、これが最適解である。
- 現在の反復点  $X, Y$  からの予測値  $\tilde{T}$  と正解  $T$  との比から次の反復点  $X$  が作られる。

$Y$  の最適化も同様に、 $G(X, Y, R)$  を  $Y_{kj}$  に関する2次関数をつくれば  $Y$  に関する最小解を求めることができる。

$$Y_{kj}^* := Y_{kj} \frac{\sum_i T_{ij} X_{ik}}{\sum_i X_{ik} \sum_{k'} X_{ik'} Y_{k'j}} = Y_{kj} \frac{\sum_i X_{ik} T_{ij}}{\sum_i X_{ik} \tilde{T}_{ij}}$$



# 非負値行列因子分解のまとめ

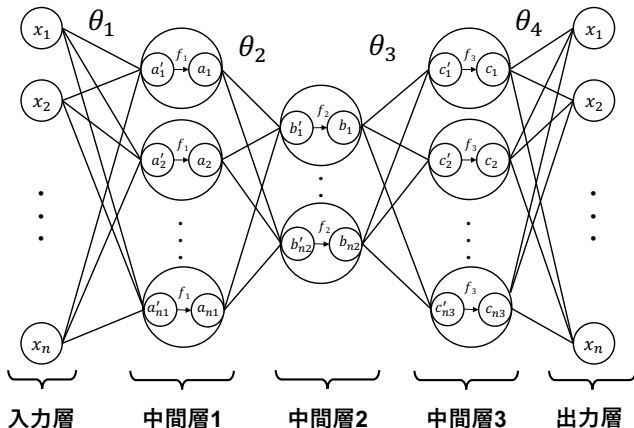
## 特徴

- 行列の掛け算を繰り返す計算  
1 反復当たりの計算量:  $O(IJK)$
- 1 反復当たりの計算量はデータの 1 乗に比例  
※ データ数を  $D := IJ$  として  $K$  を固定すると  $O(D)$
- 多くの手法よりも高速に計算できる。
- ユーザとアイテムの両側で特徴量化を行うことができる
- 非負の値を取るため解釈がしやすい。

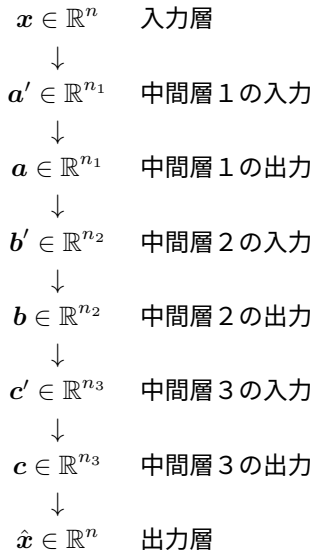
# オートエンコーダー（自己符号化）

ニューラルネットワークを使った特徴抽出

- 入力  $x_d$  から出力  $x_d$  を作るニューラルネットワーク
- 中間層の次元を小さくしておき、その情報を復元して出力  $x_d$  を作る
- 中間層の情報は  $x_d$  の主要な情報となるので、これが特徴抽出となる



# ベクトルの変換



$$\mathbf{a}' = f_{\theta_1}(\mathbf{x}) := \mathbf{M}_1 \mathbf{x} + \mathbf{v}_1$$

$$\mathbf{a} = g_1(\mathbf{a}') \quad \text{活性化関数}$$

$$\mathbf{b}' = f_{\theta_2}(\mathbf{a}) := \mathbf{M}_2 \mathbf{a} + \mathbf{v}_2$$

$$\mathbf{b} = g_2(\mathbf{b}') \quad \text{活性化関数}$$

$$\mathbf{c}' = f_{\theta_3}(\mathbf{b}) := \mathbf{M}_3 \mathbf{b} + \mathbf{v}_3$$

$$\mathbf{c} = g_3(\mathbf{c}') \quad \text{活性化関数}$$

$$\hat{\mathbf{x}} = f_{\theta_4}(\mathbf{c}) := \mathbf{M}_4 \mathbf{c} + \mathbf{v}_4$$

# 自己符号化関数と誤差関数

## 自己符号化関数

自己符号化関数  $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$

$$f(\boldsymbol{x}) = f_{\theta_4}(g_3(f_{\theta_3}(g_2(f_{\theta_2}(g_1(f_{\theta_1}(\boldsymbol{x})))))))$$

パラメタは  $\boldsymbol{\theta} := \{\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \boldsymbol{\theta}_3, \boldsymbol{\theta}_4\}$

データ  $\{\boldsymbol{x}_d\}_{d \in [D]}$  に対する予測誤差を定義する

## 誤差関数

$\boldsymbol{x}$  : 目標、 $\hat{\boldsymbol{x}} := f(\boldsymbol{x})$  : 予測

$$L(\boldsymbol{x}, \hat{\boldsymbol{x}}) := \|\boldsymbol{x} - \hat{\boldsymbol{x}}\|^2$$

# 学習

学習データ： $\{\boldsymbol{x}_d\}_{d \in [D]}$

## 学習

パラメタ  $\boldsymbol{\theta} := \{\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \boldsymbol{\theta}_3, \boldsymbol{\theta}_4\}$

$$\min_{\boldsymbol{\theta}} \frac{1}{D} \sum_{d \in [D]} L(\boldsymbol{x}_d, f_{\boldsymbol{\theta}_4}(g_3(f_{\boldsymbol{\theta}_3}(g_2(f_{\boldsymbol{\theta}_2}(g_1(f_{\boldsymbol{\theta}_1}(\boldsymbol{x}_d))))))))))$$

$F_d(\boldsymbol{\theta}) := L(\boldsymbol{x}_d, f_{\boldsymbol{\theta}_4}(g_3(f_{\boldsymbol{\theta}_3}(g_2(f_{\boldsymbol{\theta}_2}(g_1(f_{\boldsymbol{\theta}_1}(\boldsymbol{x}_d))))))))))$  とすると、

$$\min_{\boldsymbol{\theta}} \frac{1}{D} \sum_{d \in [D]} F_d(\boldsymbol{\theta})$$

## 学習アルゴリズム

確率的勾配降下法やその改良アルゴリズムで最適なパラメタ  $\boldsymbol{\theta}$  を求める。

# 誤差逆伝播 1,2

$$F_d \xleftarrow{L} \hat{\mathbf{x}} \xleftarrow{f_{\theta_4}} \mathbf{c} \xleftarrow{g_3} \mathbf{c}' \xleftarrow{f_{\theta_3}} \mathbf{b} \xleftarrow{g_2} \mathbf{b}' \xleftarrow{f_{\theta_2}} \mathbf{a} \xleftarrow{g_1} \mathbf{a}' \xleftarrow{f_{\theta_1}} \mathbf{x}_d$$

連鎖律

$$\frac{\partial F_d}{\partial \theta_1} = \frac{\partial F_d}{\partial \hat{\mathbf{x}}} \frac{\partial \hat{\mathbf{x}}}{\partial \mathbf{c}} \frac{\partial \mathbf{c}}{\partial \mathbf{c}'} \frac{\partial \mathbf{c}'}{\partial \mathbf{b}} \frac{\partial \mathbf{b}}{\partial \mathbf{b}'} \frac{\partial \mathbf{b}'}{\partial \mathbf{a}} \frac{\partial \mathbf{a}}{\partial \mathbf{a}'} \frac{\partial \mathbf{a}'}{\partial \theta_1}$$

$$\frac{\partial F_d}{\partial \theta_2} = \frac{\partial F_d}{\partial \hat{\mathbf{x}}} \frac{\partial \hat{\mathbf{x}}}{\partial \mathbf{c}} \frac{\partial \mathbf{c}}{\partial \mathbf{c}'} \frac{\partial \mathbf{c}'}{\partial \mathbf{b}} \frac{\partial \mathbf{b}}{\partial \mathbf{b}'} \frac{\partial \mathbf{b}'}{\partial \theta_2}$$

$$\frac{\partial F_d}{\partial \theta_3} = \frac{\partial F_d}{\partial \hat{\mathbf{x}}} \frac{\partial \hat{\mathbf{x}}}{\partial \mathbf{c}} \frac{\partial \mathbf{c}}{\partial \mathbf{c}'} \frac{\partial \mathbf{c}'}{\partial \theta_3}$$

$$\frac{\partial F_d}{\partial \theta_4} = \frac{\partial F_d}{\partial \hat{\mathbf{x}}} \frac{\partial \hat{\mathbf{x}}}{\partial \theta_4}$$

各パーツの微分は簡単に計算できる

## 誤差逆伝播 3

$$\frac{\partial F_d}{\partial \theta_1} = \frac{\partial F_d}{\partial \hat{x}} \frac{\partial \hat{x}}{\partial c} \frac{\partial c}{\partial c'} \frac{\partial c'}{\partial b} \frac{\partial b}{\partial b'} \frac{\partial b'}{\partial a} \frac{\partial a}{\partial a'} \frac{\partial a'}{\partial \theta_1}$$

ノード数を全て  $n$  としてバイアス項を除いたとき、

$$\frac{\partial F_d}{\partial \hat{x}} \in \mathbb{R}^{1 \times n}, \quad \frac{\partial \hat{x}}{\partial c}, \frac{\partial c}{\partial c'}, \frac{\partial c'}{\partial b}, \frac{\partial b}{\partial b'}, \frac{\partial b'}{\partial a}, \frac{\partial a}{\partial a'} \in \mathbb{R}^{n \times n}, \quad \frac{\partial a'}{\partial \theta_1} \in \mathbb{R}^{n \times n^2}$$

### 計算量

- 入力の方から計算（連鎖律の右から計算）  
 $n \times n$  のベクトルと  $n \times n^2$  の行列の積：  $O(n^4)$
  - 出力の方から計算（連鎖律の左から計算）  
 $1 \times n$  のベクトルと  $n \times n$  の行列の積：  $O(n^2)$
- ※ 最後の  $\frac{\partial a'}{\partial \theta_1}$  の掛け算は疎性を使うと簡単に計算可能

## 誤差逆伝播法 4

行列  $\frac{\partial \mathbf{a}'}{\partial \theta_1}$  の疎性（非ゼロ要素の割合）

行と列の意味

- $\theta_1$  は  $n \times n$  の行列を適当な順で並び替えたベクトル
- $\mathbf{a}'$  は中間層 1 の入力となるベクトル

$\theta_1$  の一つの成分に対して、影響を与える  $\mathbf{a}'$  は一つの成分

$\Rightarrow \frac{\partial \mathbf{a}'}{\partial \theta_1}$  の各列に非ゼロ要素は一つ

$\Rightarrow \mathbf{v} \in \mathbb{R}^n$  に対し、 $\mathbf{v}^T \frac{\partial \mathbf{a}'}{\partial \theta_1}$  の計算に必要な計算量は  $O(n^2)$   
(非ゼロの部分だけ計算すればよい)

疎性があるため行列とベクトルの掛け算は簡単



## 誤差逆伝播法 5

$$\frac{\partial F_d}{\partial \theta_1} = \frac{\partial F_d}{\partial \hat{x}} \frac{\partial \hat{x}}{\partial c} \frac{\partial c}{\partial c'} \frac{\partial c'}{\partial b} \frac{\partial b}{\partial b'} \frac{\partial b'}{\partial a} \frac{\partial a}{\partial a'} \frac{\partial a'}{\partial \theta_1}$$

$$\frac{\partial F_d}{\partial \theta_2} = \frac{\partial F_d}{\partial \hat{x}} \frac{\partial \hat{x}}{\partial c} \frac{\partial c}{\partial c'} \frac{\partial c'}{\partial b} \frac{\partial b}{\partial b'} \frac{\partial b'}{\partial \theta_2}$$

$$\frac{\partial F_d}{\partial \theta_3} = \frac{\partial F_d}{\partial \hat{x}} \frac{\partial \hat{x}}{\partial c} \frac{\partial c}{\partial c'} \frac{\partial c'}{\partial \theta_3}$$

$$\frac{\partial F_d}{\partial \theta_4} = \frac{\partial F_d}{\partial \hat{x}} \frac{\partial \hat{x}}{\partial \theta_4}$$

出力の方から順に計算すると、同じ計算が出てくるため省略できる

勾配の計算に必要な計算量（ $s$ ：中間層の層数、 $n$ ：各層の（平均）ノード数）

- 入力の方から順に計算（連鎖律の右から計算）：  $O(s^2 n^4)$
- 出力の方から順に計算（連鎖律の左から計算）：  $O(sn^2)$

出力の方から計算する手法を誤差逆伝播法（バックプロパゲーション）という。  
パラメタの数  $O(sn^2)$  個に比例した計算量で勾配の計算が可能で高速。

# オートエンコーダーのまとめ

## 特徴

- データ数よりも、中間層の層数やノード数の方が計算時間に影響  
計算量はパラメタ数と反復回数に比例
- パラメタ数を増やすと勾配の計算が困難
- 中間層の数やノード数などのハイパーパラメータを適切に決める必要がある
- 深層学習のフレームワークで計算できるため、ソフトウェアが充実

高次元のベクトル： $\{\mathbf{x}_d\}_{d \in [D]}$      $\mathbf{x}_d \in \mathbb{R}^n$



低次元のベクトル： $\{\mathbf{y}_d\}_{d \in [D]}$      $\mathbf{y}_d \in \mathbb{R}^2$

※ 特徴量の可視化を想定しているため  $\mathbf{y}_d \in \mathbb{R}^2$  を考える。  
アルゴリズムとしては任意の次元のベクトルを構築することが可能。

アイデア：

ベクトルの近傍が保存されることを目指し、高次元空間上の2つのベクトル間の類似度と対応する低次元空間上の2つのベクトルの類似度が変わらないように変換。

論点

- 高次元空間と低次元空間での類似度の定義
- 「高次元空間での類似性」と「低次元空間での類似性」間の近さを示す距離

# 学習モデル

- $P_{ij} : x_i$  と  $x_j$  の類似度
- $Q_{ij} : y_i$  と  $y_j$  の類似度
- $f(P, Q) : \text{類似度行列 } P, Q \in \mathbb{R}^{D \times D} \text{ 間の距離}$

## 学習モデル

$y_d$  ( $d \in [D]$ ) を変数とする制約なし最適化：

$$\min_{\{y_d\}_{d \in [D]}} f(P, Q)$$

※  $Q$  は  $y_d$  から計算される行列

次の順に説明を行う。

- SNE(Stochastic Neighbor Embedding)
- t-SNE(t-Distributed Stochastic Neighbor Embedding)

以下では、 $\sum$  で使う添字  $i, j, k, l \in [D]$  に対し、 $[D]$  は省略することにする。

例： $\sum_{i(i \neq j)}$  は  $\sum_{i \in [D] - \{j\}}$  の意味、 $\sum_{i, j(i \neq j)}$  は  $\sum_{i \in [D]} \sum_{j \in [D] - \{i\}}$  の意味

類似度：正規分布を利用した条件付き確率

$$P_{ij} : \quad p_{j|i} := \frac{\exp\{-\|\mathbf{x}_j - \mathbf{x}_i\|^2/2\sigma^2\}}{\sum_{k(k \neq i)} \exp\{-\|\mathbf{x}_k - \mathbf{x}_i\|^2/2\sigma^2\}} \quad i, j \in [D] \quad (i \neq j)$$

$$Q_{ij} : \quad q_{j|i} := \frac{\exp\{-\|\mathbf{y}_j - \mathbf{y}_i\|^2\}}{\sum_{k(k \neq i)} \exp\{-\|\mathbf{y}_k - \mathbf{y}_i\|^2\}} \quad i, j \in [D] \quad (i \neq j)$$

- データのばらつきを表す  $\sigma$  は、データから適切に決める。
- $\mathbf{p}_i := \{p_{j|i}\}_{j \in [D] - \{i\}}$ ,  $\mathbf{q}_i := \{q_{j|i}\}_{j \in [D] - \{i\}}$  は確率分布

類似度間の距離：KL ダイバージェンス<sup>5</sup> の和

$$C := \sum_i KL(\mathbf{p}_i || \mathbf{q}_i) = \sum_i \sum_{j(j \neq i)} p_{j|i} \log \frac{p_{j|i}}{q_{j|i}}$$

<sup>5</sup>KL ダイバージェンスは確率分布間の距離の一つ

# 学習アルゴリズム

勾配：

計算過程は次スライド

$$\nabla_{\mathbf{y}_i} C = 2 \sum_{j(j \neq i)} (p_{j|i} + p_{i|j} - q_{j|i} - q_{i|j})(\mathbf{y}_i - \mathbf{y}_j)$$

最急降下法にモメンタム項を加えた反復解法で計算をする。

## 学習アルゴリズム

ステップ0  $x_i$  と  $x_j$  の類似度  $p_{j|i}$  の計算  $i, j \in [D], (i \neq j)$

$t = 0, \mathbf{y}_i^t \in \mathbb{R}^2 (i \in [D])$  の初期値を設定

ステップ1  $\mathbf{y}_i^t$  と  $\mathbf{y}_j^t$  の類似度  $q_{j|i}$  の計算  $i, j \in [D], (i \neq j)$

ステップ2 勾配  $\nabla_{\mathbf{y}_i} C$  の計算  $i \in [D]$

ステップ3  $\mathbf{y}_i^t$  の更新 ( $\alpha, \beta^t$  はステップサイズ)  
$$\mathbf{y}_i^{t+1} := \mathbf{y}_i^t - \alpha \nabla_{\mathbf{y}_i} C + \beta^t (\mathbf{y}_i^t - \mathbf{y}_i^{t-1})$$
  $i \in [D]$

ステップ4  $t := t + 1$  として、ステップ1へ。

※  $\beta^t = 0$  とすると最急降下法

# 式変形 1

$$\begin{aligned} C &:= \sum_{i,j(i \neq j)} p_{j|i} \log \frac{p_{j|i}}{q_{j|i}} \\ &= \sum_{i,j(i \neq j)} p_{j|i} (\log p_{j|i} - \log q_{j|i}) \\ &= \sum_{i,j(i \neq j)} p_{j|i} \left( \log p_{j|i} + \|\mathbf{y}_j - \mathbf{y}_i\|^2 + \log \sum_{k(k \neq i)} \exp\{-\|\mathbf{y}_k - \mathbf{y}_i\|^2\} \right) \\ &= \sum_{i,j(i \neq j)} p_{j|i} \log p_{j|i} + \sum_{i,j(i \neq j)} p_{j|i} \|\mathbf{y}_j - \mathbf{y}_i\|^2 \\ &\quad + \sum_i \log \sum_{k(k \neq i)} \exp\{-\|\mathbf{y}_k - \mathbf{y}_i\|^2\} \quad \left( \because \sum_{j(j \neq i)} p_{j|i} = 1 \right) \end{aligned}$$

## 式変形 2

$\mathbf{y}_s$  での微分  $\frac{\partial C}{\partial \mathbf{y}_s}$  を考える。 $C$  の中で  $\mathbf{y}_s$  に依存する部分は次の4つ。

- $C_1 := \sum_{j(j \neq s)} p_{j|s} \|\mathbf{y}_j - \mathbf{y}_s\|^2$  ... 2 項目の  $j = s$  の場合
- $C_2 := \sum_{i(i \neq s)} p_{s|i} \|\mathbf{y}_s - \mathbf{y}_i\|^2$  ... 2 項目の  $i = s$  の場合
- $C_3 := \log \sum_{k(k \neq s)} \exp\{-\|\mathbf{y}_k - \mathbf{y}_s\|^2\}$  ... 3 項目の  $i = s$  の場合
- $C_4 := \sum_{i(i \neq s)} \log \sum_{k(k \neq i)} \exp\{-\|\mathbf{y}_k - \mathbf{y}_i\|^2\}$  ... 3 項目の  $i \neq s$  の場合

それぞれを  $\mathbf{y}_s$  で微分する。

$$\frac{\partial C_1}{\partial \mathbf{y}_s} = 2 \sum_{j(j \neq s)} p_{j|s} (\mathbf{y}_s - \mathbf{y}_j)^T, \quad \frac{\partial C_2}{\partial \mathbf{y}_s} = 2 \sum_{i(i \neq s)} p_{s|i} (\mathbf{y}_s - \mathbf{y}_i)^T$$



## 式変形 3

$$\begin{aligned}\frac{\partial C_3}{\partial \mathbf{y}_s} &= \frac{1}{\sum_{k(k \neq s)} \exp\{-\|\mathbf{y}_k - \mathbf{y}_s\|^2\}} \sum_{k(k \neq s)} \exp\{-\|\mathbf{y}_k - \mathbf{y}_s\|^2\} (-2) (\mathbf{y}_s - \mathbf{y}_k)^T \\ &= -2 \sum_{k(k \neq s)} q_{k|s} (\mathbf{y}_s - \mathbf{y}_k)^T \\ \frac{\partial C_4}{\partial \mathbf{y}_s} &= \sum_{i(i \neq s)} \frac{1}{\sum_{k(k \neq i)} \exp\{-\|\mathbf{y}_k - \mathbf{y}_i\|^2\}} \exp\{-\|\mathbf{y}_s - \mathbf{y}_i\|^2\} (-2) (\mathbf{y}_s - \mathbf{y}_i)^T \\ &= -2 \sum_{i(i \neq s)} q_{s|i} (\mathbf{y}_s - \mathbf{y}_i)^T\end{aligned}$$

よって、

$$\begin{aligned}\nabla_{\mathbf{y}_s} C &= \left( \frac{\partial C}{\partial \mathbf{y}_s} \right)^T = \left( \frac{\partial C_1}{\partial \mathbf{y}_s} + \frac{\partial C_2}{\partial \mathbf{y}_s} + \frac{\partial C_3}{\partial \mathbf{y}_s} + \frac{\partial C_4}{\partial \mathbf{y}_s} \right)^T \\ &= 2 \sum_{j(j \neq s)} (p_{j|s} + p_{s|j} - q_{j|s} - q_{s|j}) (\mathbf{y}_s - \mathbf{y}_j)\end{aligned}$$

# t-SNE

t-SNE(t-Distributed Stochastic Neighbor Embedding) では、

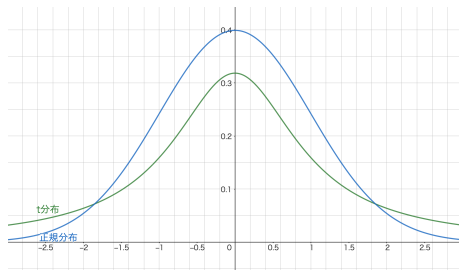
- (1) 低次元空間での確率分布を正規分布でなく  $t$  分布にする
- (2) 同時確率を考えることにより対称化を行う

## (1) $t$ 分布の導入

高次元のベクトルを低次元に埋め込むと crowding という問題が生じる。  
正規分布よりも裾が厚い  $t$  分布を使うことにより、類似していないデータを遠くに置くことを許容できるようになる。

自由度 1 の  $t$  分布：

$$f(x) := \frac{1}{\pi(1+x^2)}$$



# t 分布の導入

類似度：t 分布による条件付き確率

$$Q_{ij} : \quad q_{j|i} := \frac{(1 + \|\mathbf{y}_j - \mathbf{y}_i\|^2)^{-1}}{\sum_{k(k \neq i)} (1 + \|\mathbf{y}_k - \mathbf{y}_i\|^2)^{-1}}$$

この場合の勾配は次のようになる。

$$\nabla_{\mathbf{y}_i} C = 2 \sum_{j(j \neq i)} \frac{p_{j|i} + p_{i|j} - q_{j|i} - q_{i|j}}{1 + \|\mathbf{y}_j - \mathbf{y}_i\|^2} (\mathbf{y}_i - \mathbf{y}_j)$$

導出過程は省略するが、計算の手順は SNE や t-SNE の場合と同じ。

## (2) 同時確率による対称化

条件付き確率  $p_{j|i}$  の代わりに同時確率  $p_{ij}$  を考える

$$P_{ij} : \quad p_{ij} := \frac{\exp\{-\|\mathbf{x}_i - \mathbf{x}_j\|^2/2\sigma^2\}}{\sum_{k,l(k \neq l)} \exp\{-\|\mathbf{x}_k - \mathbf{x}_l\|^2/2\sigma^2\}} \quad i, j \in [D], (i \neq j)$$

$$Q_{ij} : \quad q_{ij} := \frac{(1 + \|\mathbf{y}_i - \mathbf{y}_j\|^2)^{-1}}{\sum_{k,l(k \neq l)} (1 + \|\mathbf{y}_k - \mathbf{y}_l\|^2)^{-1}} \quad i, j \in [D], (i \neq j)$$

- $p_{ij} = p_{ji}$ ,  $q_{ij} = q_{ji}$  となり類似度に対称性がある。
- $P, Q$  は確率分布

類似度間の距離：KL ダイバージェンス

$$C := KL(P||Q) = \sum_{i,j(i \neq j)} p_{ij} \log \frac{p_{ij}}{q_{ij}}$$

# 工夫を加えた対称化

$\mathbf{x}_i$  が他の点から遠く外れていた場合、同時確率

$$p_{ij} := \frac{\exp\{-\|\mathbf{x}_i - \mathbf{x}_j\|^2/2\sigma^2\}}{\sum_{k,l(k \neq l)} \exp\{-\|\mathbf{x}_k - \mathbf{x}_l\|^2/2\sigma^2\}}$$

は小さくなり過ぎる。そのため、代わりに次の確率を使うとよい（らしい）。

$$p_{j|i} := \frac{\exp\{-\|\mathbf{x}_j - \mathbf{x}_i\|^2/2\sigma^2\}}{\sum_{k(k \neq i)} \exp\{-\|\mathbf{x}_k - \mathbf{x}_i\|^2/2\sigma^2\}}$$

条件付き確率

$$p_{ij} := \frac{p_{j|i} + p_{i|j}}{2D}$$

対称化

- $p_{ij} = p_{ji}$  であるため、類似度に対称性がある。
- $P$  を確率分布にするため  $D$  で割っている。
- $q_{ij}$  は変えない（変えると勾配の計算が困難）。

# 勾配

この場合の勾配は次のようになる。

$$\nabla \mathbf{y}_i C = 4 \sum_{j(j \neq i)} \frac{p_{ij} - q_{ij}}{1 + \|\mathbf{y}_i - \mathbf{y}_j\|^2} (\mathbf{y}_i - \mathbf{y}_j)$$

目的関数：

$$\begin{aligned} C &:= \sum_{i,j(i \neq j)} p_{ij} \log \frac{p_{ij}}{q_{ij}} \\ &= \sum_{i,j(i \neq j)} p_{ij} \left( \log p_{ij} + \log(1 + \|\mathbf{y}_i - \mathbf{y}_j\|^2) + \log \sum_{k,l(k \neq l)} (1 + \|\mathbf{y}_k - \mathbf{y}_l\|^2)^{-1} \right) \\ &= \sum_{i,j(i \neq j)} p_{ij} \log p_{ij} + \sum_{i,j(i \neq j)} p_{ij} \log(1 + \|\mathbf{y}_i - \mathbf{y}_j\|^2) \\ &\quad + \log \sum_{k,l(k \neq l)} (1 + \|\mathbf{y}_k - \mathbf{y}_l\|^2)^{-1} \quad (\because \sum_{i,j(i \neq j)} p_{ij} = 1) \end{aligned}$$

# 式変形 1

$\mathbf{y}_s$  で微分することを考える。 $\mathbf{y}_s$  に依存する部分は次の3つ。

- $C_1 := \sum_{j(j \neq s)} p_{sj} \log(1 + \|\mathbf{y}_s - \mathbf{y}_j\|^2)$  ... 2項目の  $i = s$  の場合
- $C_2 := \sum_{i(i \neq s)} p_{is} \log(1 + \|\mathbf{y}_i - \mathbf{y}_s\|^2)$  ... 2項目の  $j = s$  の場合
- $C_3 := \log \sum_{k,l(k \neq l)} (1 + \|\mathbf{y}_k - \mathbf{y}_l\|^2)^{-1}$  ... 3項目

それぞれを微分する。

$$\frac{\partial C_1}{\partial \mathbf{y}_s} = \sum_{j(j \neq s)} p_{sj} \frac{1}{1 + \|\mathbf{y}_s - \mathbf{y}_j\|^2} 2(\mathbf{y}_s - \mathbf{y}_j)^T$$
$$\frac{\partial C_2}{\partial \mathbf{y}_s} = \sum_{i(i \neq s)} p_{is} \frac{1}{1 + \|\mathbf{y}_i - \mathbf{y}_s\|^2} 2(\mathbf{y}_s - \mathbf{y}_i)^T$$

## 式変形 2

$$\begin{aligned}
 \frac{\partial C_3}{\partial \mathbf{y}_s} &= \frac{1}{\sum_{k,l(k \neq l)} (1 + \|\mathbf{y}_k - \mathbf{y}_l\|^2)^{-1}} \frac{\partial}{\partial \mathbf{y}_s} \sum_{k,l(k \neq l)} (1 + \|\mathbf{y}_k - \mathbf{y}_l\|^2)^{-1} \\
 &= \frac{1}{\sum_{k,l(k \neq l)} (1 + \|\mathbf{y}_k - \mathbf{y}_l\|^2)^{-1}} 2 \frac{\partial}{\partial \mathbf{y}_s} \sum_{l(l \neq s)} (1 + \|\mathbf{y}_s - \mathbf{y}_l\|^2)^{-1} \\
 &= \frac{1}{\sum_{k,l(k \neq l)} (1 + \|\mathbf{y}_k - \mathbf{y}_l\|^2)^{-1}} 2 \sum_{l(l \neq s)} -(1 + \|\mathbf{y}_s - \mathbf{y}_l\|^2)^{-2} 2(\mathbf{y}_s - \mathbf{y}_l)^T \\
 &= -4 \sum_{l(l \neq s)} \frac{q_{sl}}{1 + \|\mathbf{y}_s - \mathbf{y}_l\|^2} (\mathbf{y}_s - \mathbf{y}_l)^T
 \end{aligned}$$

よって、

$$\begin{aligned}
 \nabla_{\mathbf{y}_s} C &= \left( \frac{\partial C}{\partial \mathbf{y}_s} \right)^T = \left( \frac{\partial C_1}{\partial \mathbf{y}_s} + \frac{\partial C_2}{\partial \mathbf{y}_s} + \frac{\partial C_3}{\partial \mathbf{y}_s} \right)^T \\
 &= 4 \sum_{j(j \neq s)} \frac{p_{sj} - q_{sj}}{1 + \|\mathbf{y}_s - \mathbf{y}_j\|^2} (\mathbf{y}_s - \mathbf{y}_j)
 \end{aligned}$$



## 学習アルゴリズム

ステップ0  $p_{j|i} := \frac{\exp\{-\|\mathbf{x}_i - \mathbf{x}_j\|^2/2\sigma^2\}}{\sum_{k(k \neq i)} \exp\{-\|\mathbf{x}_i - \mathbf{x}_k\|^2/2\sigma^2\}}$   $i, j \in [D], (i \neq j)$

$$p_{ij} := \frac{p_{j|i} + p_{i|j}}{2D} \quad i, j \in [D], (i \neq j)$$

$t = 0, \mathbf{y}_i^t \in \mathbb{R}^2 (i \in [D])$  の初期値を設定

ステップ1 類似度の計算

$$q_{ij} := \frac{(1 + \|\mathbf{y}_i^t - \mathbf{y}_j^t\|^2)^{-1}}{\sum_{k, l(k \neq l)} (1 + \|\mathbf{y}_k^t - \mathbf{y}_l^t\|^2)^{-1}} \quad i, j \in [D], (i \neq j)$$

ステップ2 勾配の計算

$$\nabla_{\mathbf{y}_i} C = 4 \sum_{j(j \neq i)} \frac{p_{ij} - q_{ij}}{1 + \|\mathbf{y}_i^t - \mathbf{y}_j^t\|^2} (\mathbf{y}_i^t - \mathbf{y}_j^t) \quad i \in [D]$$

ステップ3  $\mathbf{y}_i^t$  の更新 ( $\alpha, \beta^t$  はステップサイズ)

$$\mathbf{y}_i^{t+1} := \mathbf{y}_i^t - \alpha \nabla_{\mathbf{y}_i} C + \beta^t (\mathbf{y}_i^t - \mathbf{y}_i^{t-1}) \quad i \in [D]$$

ステップ4  $t := t + 1$  として、ステップ1へ。

# t-SNE のまとめ

## 特徴

- 人間の直感に合うような可視化となることが多い。  
もちろん、必ず良いものが出来るということではない。
- 非線形な変換を伴った可視化ツールとしてよく使われる。  
データの特徴抽出法としても時々使われる。
- 1反復あたりの計算量はデータ数の2乗に比例  
計算量： $O(D^2)$   
※  $\mathbf{y} \in \mathbb{R}^k$  のとき、 $O(D^2k)$  となる。
- データ量が増えると計算が困難になるが、カーネル主成分分析ほどではない