

Goodfellow : 深層学習

Part1 : 応用数学と機械学習の基礎

第2章 : 線形代数

ノルムの整理

- L^p ノルム ($p < \infty$) の一般的な定義は以下の通り。

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

- 機械学習でよく使用されるのは、 L^2, L^1, L^∞ ノルム。

$$\|\mathbf{x}\|_2 = \sqrt{\sum_i |x_i|^2}$$

ユークリッドノルム

$$\|\mathbf{x}\|_1 = \sum_i |x_i|$$

$$\|\mathbf{x}\|_\infty = \max_i |x_i|$$

最大値ノルム

- フロベニウスノルム：行列の各成分の2乗の総和の正の平方根

$$\|\mathbf{A}\|_F = \sqrt{\sum_{i,j} (A_{i,j})^2}$$

固有値分解

- **固有値と固有ベクトル**

- 正方行列 A に対して、次の式を満たす非ゼロベクトル v を A の**固有ベクトル**、スカラー λ を**固有値**という。
- v が A の固有ベクトルなら、その定数倍も同じ固有値をもつ固有ベクトルである。
→ 固有ベクトルを考える際はノルムが1の単位固有ベクトルを基本的には考える。

$$Av = \lambda v$$

- **固有値分解** (Goodfellow本で考える固有値分解)

- 実対称行列 A の固有値と固有ベクトルの組み合わせを $(\lambda_n, v^{(n)})$ とおく。
- 行列 V と Λ を $V = [v^{(1)}, \dots, v^{(n)}]$, $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$ とおく。

$$\begin{aligned} A &= \lambda_1 v^{(1)} (v^{(1)})^\top + \dots + \lambda_n v^{(n)} (v^{(n)})^\top \\ &= V \Lambda V^\top \end{aligned}$$

特異値分解

- $m \times n$ の行列 A は $m \neq n$ の時には固有値分解ができない。
しかし、特異値分解は任意の実行列に対して可能である。
- $A \in \mathbb{R}^{m \times n}$ に対して、 $U \in \mathbb{R}^{m \times m}$, $D \in \mathbb{R}^{m \times n}$, $V \in \mathbb{R}^{n \times n}$ を用いて、
以下のような積で表すことを**特異値分解**という。

$$A = U D V^{\top}$$

- D の対角成分を行列 A の**特異値**といい、
 U, V の各列をそれぞれ**左特異ベクトル**、**右特異ベクトル**という。
- 行列 A の左特異ベクトルは行列 AA^{\top} の固有ベクトルであり、
行列 AA^{\top} の固有値の平方根は行列 A の特異値と一致する。

トレース演算子

- 行列 A の対角成分の総和として、 $\text{Tr}(A)$ が定義される。
このトレース演算子について以下の関係式が成立する。

① 行列 A のフロベニウスノルム：
$$\|A\|_F = \sqrt{\text{Tr}(AA^\top)} = \sqrt{\text{Tr}(A^\top A)}$$

② 行列 A の転置行列のトレース：
$$\text{Tr}(A^\top) = \text{Tr}(A)$$

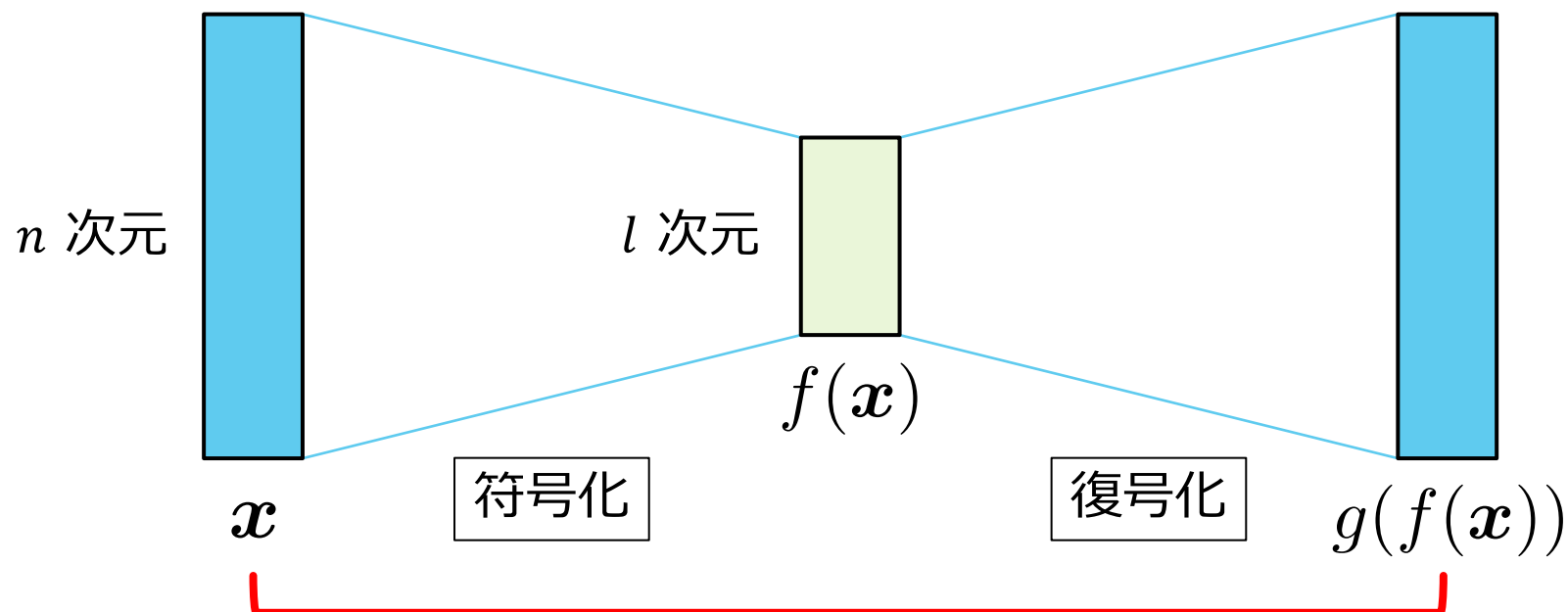
③ 3つの正方行列の積のトレース：
$$\text{Tr}(ABC) = \text{Tr}(CAB) = \text{Tr}(BCA)$$

④ 2つの行列 A, B が正方行列でない場合の積のトレース：

$$\text{Tr}(AB) = \text{Tr}(BA)$$

主成分分析の導出 (1/11)

- \mathbb{R}^n における m 個の点の集まり $x^{(1)}, \dots, x^{(m)}$ の精度はそのままに低次元に圧縮することを考える。
- ここでは、上記の次元圧縮を符号化と復号化の枠組みで考える。



両者がおおよそ一致するような変換 f と g が得ることができれば、適切な次元圧縮が可能になる。

主成分分析の導出 (2/11)

- **主成分分析 (PCA)** : 復号化の変換として線形変換を使用。

$$g(\mathbf{c}) = D\mathbf{c}$$

- 問題の簡単化のため、いくつかの制約を設ける。
 - Dの列が互いに直交することを仮定する。
 - Dとcを適当に定数倍したのも次元圧縮先としては適切な場合がありうるため、圧縮先の候補は無数に存在する。
→ Dの全ての列ベクトルが単位ベクトルであることを仮定する。

主成分分析の導出 (3/11)

- 主成分分析の定式化

→ 入力 \mathbf{x} と復号化器の出力 $g(\mathbf{c})$ の値が十分近くなるような \mathbf{x} の符号化先 $\mathbf{c} = f(\mathbf{x})$ を探索すること

$$\mathbf{c}^* = \operatorname{argmax}_{\mathbf{c}} \|\mathbf{x} - g(\mathbf{c})\|_2^2$$

- $g(\mathbf{c}) = \mathbf{D}\mathbf{c}$ なので、上の最適化問題は以下のようにになる。

$$\begin{aligned}\mathbf{c}^* &= \operatorname{argmax}_{\mathbf{c}} \|\mathbf{x} - \mathbf{D}\mathbf{c}\|_2^2 \\ &= \operatorname{argmax}_{\mathbf{c}} (\mathbf{x}^\top \mathbf{x} - 2\mathbf{x}^\top \mathbf{D}\mathbf{c} + \mathbf{c}^\top \mathbf{D}^\top \mathbf{D}\mathbf{c}) \\ &= \operatorname{argmax}_{\mathbf{c}} (-2\mathbf{x}^\top \mathbf{D}\mathbf{c} + \mathbf{c}^\top \mathbf{c})\end{aligned}$$

\mathbf{c} に依存しない

\mathbf{D} の各列が互いに直交し、かつ単位ベクトルであるので、 $\mathbf{D}^\top \mathbf{D}$ は単位行列となる。

主成分分析の導出 (4/11)

- c^* を求める。
 - c に関する勾配が 0 になる点を求める。

$$c^* = \operatorname{argmax}_c (c^\top c - 2x^\top Dc)$$

$$\Rightarrow \nabla_c (c^\top c - 2x^\top Dc) = 2c - 2D^\top x = 0$$

$$\Rightarrow c^* = D^\top x$$

【ベクトルの微分の公式】

$$\frac{\partial}{\partial x} (a^\top x) = a$$

$$\frac{\partial}{\partial x} (x^\top x) = 2x$$

- 以上より、復号化器が $g(c) = Dc$ の時に符号化器を $f(x) = D^\top x$ と設計することが最適であることがわかった。
- この最適設計の下では $g(f(x)) = DD^\top x$ であるため、
行列 D は L2 ノルム $\|x - DD^\top x\|_2^2$ が最小となるように設計することが望ましい。

主成分分析の導出 (5/11)

- 行列 D の設計

- m 個のデータの変換に対して共通の D を使用するため、全てのデータの処理に対して有効な行列を考える。
- 「 D の列が互いに直交する」という制約条件を考慮する。

$$D^* = \operatorname{argmin}_D \sqrt{\sum_{i=1}^m \sum_{j=1}^n \left(x_j^{(i)} - (DD^\top x^{(i)})_j \right)^2}, \quad \text{subject to } \underline{D^\top D = I}$$

入出力ベクトルの
各次元の一致

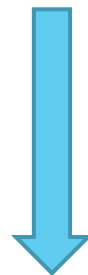
m 個のデータ全ての総和をとることで、
全データに効果的な D を選択できるようにする。

D の列が互いに直交することに
対応する条件

主成分分析の導出 (6/11)

- 以下では $l = 1$ の場合を考えて、行列 D の設計方法を導出する。
 - 行列 D のサイズは $n \times l$ である。 $l = 1$ の場合、 D は n 次元ベクトルとなる。
(そこで D の代わりに d を用いて記すことにする。)
 - ベクトルの $L2$ ノルムを使用すると以下の式と等価になる。
→ この式を変形することで主成分分析の導出法を示す。

$$\mathbf{d}^* = \operatorname{argmin}_{\mathbf{d}} \sum_i \left(\|\mathbf{x}^{(i)} - \mathbf{d}\mathbf{d}^\top \mathbf{x}^{(i)}\|_2 \right)^2, \quad \text{subject to } \mathbf{d}^\top \mathbf{d} = 1$$



$\mathbf{d}^\top \mathbf{x}^{(i)}$ はスカラー：前に移動可能

スカラー：自身と転置したものが同一

$$\mathbf{d}^* = \operatorname{argmin}_{\mathbf{d}} \sum_i \left(\|\mathbf{x}^{(i)} - (\mathbf{x}^{(i)})^\top \mathbf{d}\mathbf{d}\|_2 \right)^2, \quad \text{subject to } \mathbf{d}^\top \mathbf{d} = 1$$

主成分分析の導出 (7/11)

- L2ノルムの部分を変形。

$$\left(\|\mathbf{x}^{(i)} - (\mathbf{x}^{(i)})^\top \mathbf{d} \mathbf{d}^\top\|_2 \right)^2 = \sum_{j=1}^n \left[x_j^{(i)} - \left(\sum_{k=1}^n x_k^{(i)} d_k \right) d_j \right]^2 = \sum_{j=1}^n \left[x_j^{(i)} - \underbrace{\left(\sum_{k=1}^n x_k^{(i)} (d_k d_j) \right)} \right]^2$$

- 緑の下線部は以下の行列積の第 j 成分になる。

$$\begin{pmatrix} x_1^{(i)} & x_2^{(i)} & \cdots & x_n^{(i)} \end{pmatrix} \begin{pmatrix} d_1 d_1 & d_1 d_2 & \cdots & d_1 d_n \\ d_2 d_1 & d_2 d_2 & \cdots & d_2 d_n \\ \vdots & \vdots & \ddots & \vdots \\ d_n d_1 & d_n d_2 & \cdots & d_n d_n \end{pmatrix}$$

- \mathbf{X} を $\mathbf{X}_{i,:} = (\mathbf{x}^{(i)})^\top$ により定義すると、最適化問題は以下の式に帰着する。

$$\mathbf{d}^* = \underset{\mathbf{d}}{\operatorname{argmin}} \left(\|\mathbf{X} - \mathbf{X} \mathbf{d} \mathbf{d}^\top\|_F \right)^2, \quad \text{subject to } \mathbf{d}^\top \mathbf{d} = 1$$

主成分分析の導出 (8/11)

- フロベニウスノルムの部分を式変形する。

$$\operatorname{argmin}_d (\| \mathbf{X} - \mathbf{X} \mathbf{d} \mathbf{d}^\top \|_F)^2$$

$$= \operatorname{argmin}_d \operatorname{Tr} ((\mathbf{X} - \mathbf{X} \mathbf{d} \mathbf{d}^\top)^\top (\mathbf{X} - \mathbf{X} \mathbf{d} \mathbf{d}^\top))$$

$$= \operatorname{argmin}_d \operatorname{Tr} ((\mathbf{X}^\top - \mathbf{d} \mathbf{d}^\top \mathbf{X}^\top)(\mathbf{X} - \mathbf{X} \mathbf{d} \mathbf{d}^\top))$$

$$= \operatorname{argmin}_d \operatorname{Tr} (\underbrace{\mathbf{X}^\top \mathbf{X}}_{\text{d に依存しない}} - \underbrace{\mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top}_{\text{転置行列と元の行列のトレースは同一}} - \underbrace{\mathbf{d} \mathbf{d}^\top \mathbf{X}^\top \mathbf{X}}_{\text{転置行列と元の行列のトレースは同一}} + \mathbf{d} \mathbf{d}^\top \mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top)$$

$$= \operatorname{argmin}_d (-2\operatorname{Tr}(\mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top) + \operatorname{Tr}(\mathbf{d} \mathbf{d}^\top \mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top))$$

$$= \operatorname{argmin}_d (-2\operatorname{Tr}(\mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top) + \operatorname{Tr}(\mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top \mathbf{d} \mathbf{d}^\top))$$

フロベニウスノルムは
トレースを使って書き換えられる

転置行列と元の行列の
トレースは同一

d に依存しない

中身の行列を巡回させても
トレースは不変

主成分分析の導出 (9/11)

- 前ページの式変形より、 d に関する最適化問題は以下の式で表せる。

$$\mathbf{d}^* = \operatorname{argmin}_{\mathbf{d}} \left(-2\operatorname{Tr}(\mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top) + \operatorname{Tr}(\mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top \mathbf{d} \mathbf{d}^\top) \right), \quad \text{subject to } \mathbf{d}^\top \mathbf{d} = 1$$

- 制約条件とトレースの性質から次の式と等価である。

$$\begin{aligned} \mathbf{d}^* &= \operatorname{argmin}_{\mathbf{d}} \left(-\operatorname{Tr}(\mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top) \right), \quad \text{subject to } \mathbf{d}^\top \mathbf{d} = 1 \\ &= \operatorname{argmax}_{\mathbf{d}} \left(\operatorname{Tr}(\mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top) \right), \quad \text{subject to } \mathbf{d}^\top \mathbf{d} = 1 \\ &= \operatorname{argmax}_{\mathbf{d}} \left(\operatorname{Tr}(\mathbf{d}^\top \mathbf{X}^\top \mathbf{X} \mathbf{d}) \right), \quad \text{subject to } \mathbf{d}^\top \mathbf{d} = 1 \end{aligned}$$

主成分分析の導出 (10/11)

- Lagrangeの未定乗数法を利用する。

$$d^* = \operatorname{argmax}_d (\operatorname{Tr}(d^\top X^\top X d)), \quad \text{subject to } d^\top d = 1$$



Lagrange関数

$$d^* = \operatorname{argmax}_d (\operatorname{Tr}(d^\top X^\top X d) - \lambda(d^\top d - 1))$$

- ベクトルの微分の公式を利用する。

- 停留点 d^* は次の式を満たす。

$$2(X^\top X d^* - \lambda d^*) = 0 \implies X^\top X d^* = \lambda d^*$$

【ベクトルの微分の公式】

$$\frac{\partial}{\partial A} (\operatorname{Tr}(A B A^\top)) = A(B + B^\top)$$

$$\frac{\partial}{\partial x} (x^\top x) = 2x$$

$$\frac{\partial}{\partial A^\top} f(A) = \left(\frac{\partial}{\partial A} f(A) \right)^\top$$

主成分分析の導出 (11/11)

- 以上より、最適な d を求める問題は $X^T X$ の固有ベクトルに対応する。

$$d^* = \operatorname{argmax}_d (\operatorname{Tr}(d^T X^T X d) - \lambda(d^T d - 1))$$

$$X^T X d^* = \lambda d^*$$

- $X^T X d^* = \lambda d^*$ の条件をLagrange関数に代入すると、その値は λ となる。
→ Lagrange関数が最大となるのは λ が最大の時。
- つまり、最適な d を求める問題は $X^T X$ の最大固有値に対応する固有ベクトルを求める固有値問題と等価である。