

第10回 主成分分析

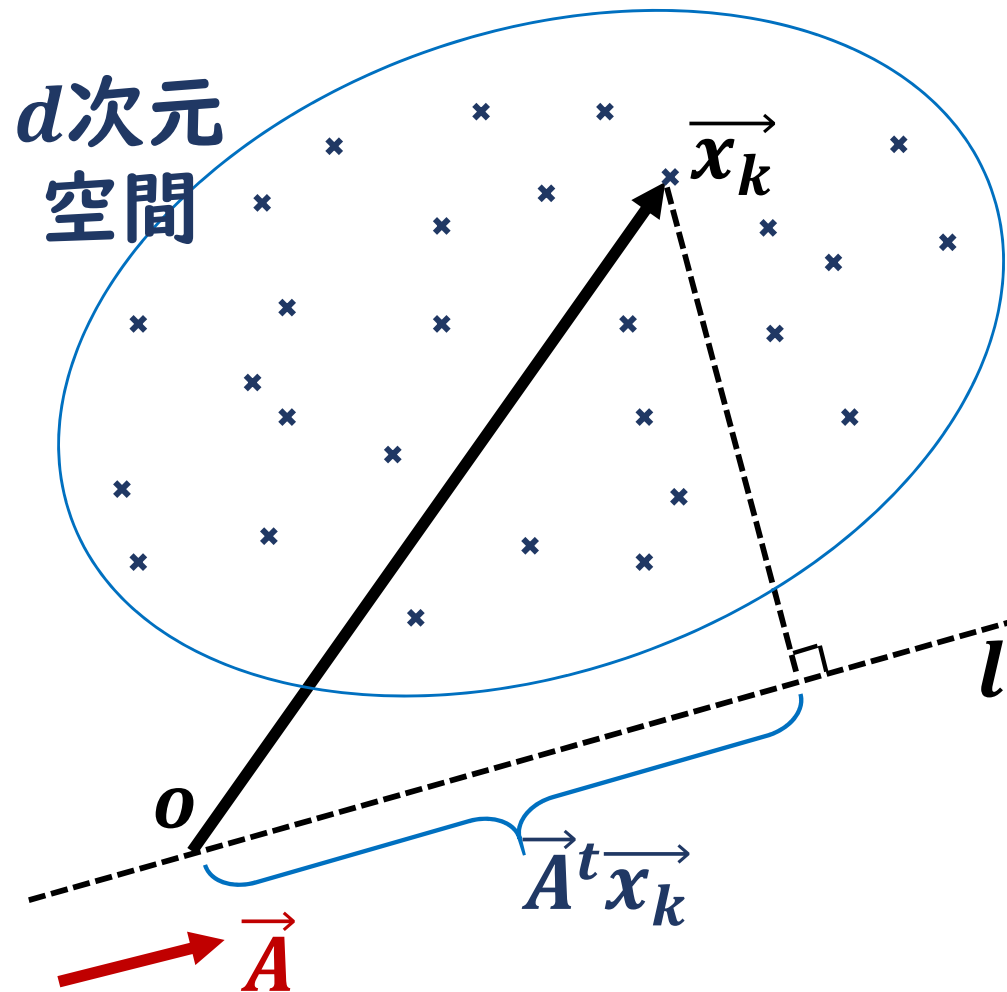
画像認識工学

情報工学科4年

科目担当 鈴木

主成分分析 (PCA: Principal Component Analysis)

d 次元空間の多数のデータ $\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n$ の分布において、分散最大の軸 \vec{A} (主成分) を求める手法を主成分分析という。



\vec{A} は直線 l の方向ベクトルなので単位ベクトルと仮定しても一般性は失われない。

直線 l は原点を通る直線と仮定しても一般性は失われない。

\vec{x} を直線 l に正射影するとき、原点からその点までの長さは、 $\vec{A}^t \vec{x}_k$ である。

正射影後の分散

$\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n$ の平均（平均ベクトル） \vec{m} は次式で表される。

$$\vec{m} = \frac{1}{n} \sum_{k=1}^n \vec{x}_k$$

\vec{x}_k を直線 l に正射影したときの値は $\vec{A}^t \vec{x}_k$ だから、その平均値 m_A は次のようになる。

$$m_A = \frac{1}{n} \sum_{k=1}^n \vec{A}^t \vec{x}_k = \vec{A}^t \left[\frac{1}{n} \sum_{k=1}^n \vec{x}_k \right] = \vec{A}^t \vec{m}$$

直線 l 上のデータの分散 σ^2 は次のようになる。

$$\sigma^2 = \frac{1}{n} \sum_{k=1}^n \left(\vec{A}^t \vec{x}_k - \vec{A}^t \vec{m} \right)^2$$

σ^2 を最大にする
 \vec{A} を求める!

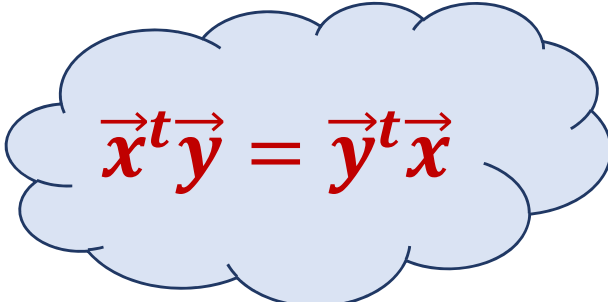
分散の解析

$$\sigma^2 = \frac{1}{n} \sum_{k=1}^n \left(\vec{A}^t \vec{x}_k - \vec{A}^t \vec{m} \right)^2 = \frac{1}{n} \sum_{k=1}^n \left\{ \vec{A}^t (\vec{x}_k - \vec{m}) \right\}^2$$

$$= \frac{1}{n} \sum_{k=1}^n \underbrace{\left\{ \vec{A}^t (\vec{x}_k - \vec{m}) \cdot \vec{A}^t (\vec{x}_k - \vec{m}) \right\}}_{\text{スカラー}}$$

$$= \frac{1}{n} \sum_{k=1}^n \left\{ \vec{A}^t (\vec{x}_k - \vec{m}) \cdot (\vec{x}_k - \vec{m})^t \vec{A} \right\}$$

$$= \vec{A}^t \left[\frac{1}{n} \sum_{k=1}^n \{ (\vec{x}_k - \vec{m}) (\vec{x}_k - \vec{m})^t \} \right] \vec{A} = \vec{A}^t \Sigma \vec{A}$$


$$\vec{x}^t \vec{y} = \vec{y}^t \vec{x}$$

Work①

d 次元空間の多数のデータ $\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n$ を考える。

$\vec{x}_k = (x_{k1}, x_{k2}, \dots, x_{kd})^t$ とするとき、下記の行列について次の値を求めよ。ただし、 $\vec{m} = (m_1, m_2, \dots, m_d)^t$ は $\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n$ の平均ベクトルとする。

$$\Sigma = \frac{1}{n} \sum_{k=1}^n \{(\vec{x}_k - \vec{m})(\vec{x}_k - \vec{m})^t\}$$

- ① Σ の1行1列成分を求めよ。
- ② Σ の1行2列成分を求めよ。
- ③ Σ の*i*行*j*列成分を求めよ。

共分散行列 (Covariance Matrix)

d 次元空間の多数のデータ $\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n$ に対し、行列の i, j 成分が、データの第 i 番目の要素と第 j 番目の要素の共分散になっているような行列 Σ を分散・共分散行列 (Variance-Covariance Matrix) または共分散行列 (Covariance Matrix) という。

共分散行列 Σ は d 次の実対称行列であり次式で計算できる。

$$\Sigma = \frac{1}{n} \sum_{k=1}^n \{(\vec{x}_k - \vec{m})(\vec{x}_k - \vec{m})^t\}$$

分散最大性

$$\sigma^2 = \vec{A}^t \Sigma \vec{A} \quad \Sigma = \frac{1}{n} \sum_{k=1}^n \{(\vec{x}_k - \vec{m})(\vec{x}_k - \vec{m})^t\}$$

↓
最大化

□ Σ はデータから求まる定数行列。

□ $\|\vec{A}\|$ を大きくすればいくらでも大きくなる。
(\vec{A} を単位ベクトルと仮定していた。。。)

↓
 $\|\vec{A}\| = 1$ という制約のもとで $\sigma^2 = \vec{A}^t \Sigma \vec{A}$ を最大にする問題を解けばよい。このような問題を **変分問題** という。

ラグランジュの未定係数法

$g(\vec{x}) = 0$ という制約のもとで $f(\vec{x})$ が最大(最小)となる \vec{x} を求めるためには、 λ を定数として $F(\vec{x}) = f(\vec{x}) - \lambda g(\vec{x})$ とおき、

$$\frac{\partial F(\vec{x})}{\partial x_k} = \frac{\partial F(\vec{x})}{\partial \lambda} = 0$$

$$k = 1, 2, \dots, d \quad \vec{x}^t = (x_1, x_2, \dots, x_d)$$

を解けばよい。この手法をラグランジュ (Lagrange) の未定係数 (乗数) 法という。

Work②

- ① 二次曲線 $x^2 + y^2 - 2xy - \sqrt{2}x - \sqrt{2}y + 2 = 0$ 上の点のうち、原点に一番近い点を求めよ。
- ② 三次元曲面 $z = 2x^2 + 3y^2 - 12y + 13$ において、点 $(0, 3, 4)$ における接平面の方程式を求めよ。

分散最大性

$\|\vec{A}\| = 1$ 、つまり $g(\vec{a}) = \vec{A}^t \vec{A} - 1 = 0$ という制約のもとで
 $f(\vec{A}) = \sigma^2$ が最大になる $\vec{A} = (A_1, A_2, \dots, A_d)^t$ を求めれば良い。

$$\sigma^2 = \vec{A}^t \Sigma \vec{A} \quad \Sigma = \frac{1}{n} \sum_{k=1}^n \{(\vec{x}_k - \vec{m})(\vec{x}_k - \vec{m})^t\}$$



$$F(\vec{A}) = \vec{A}^t \left[\frac{1}{n} \sum_{k=1}^n \{(\vec{x}_k - \vec{m})(\vec{x}_k - \vec{m})^t\} \right] \vec{A} - \lambda (\vec{A}^t \vec{A} - 1)$$

$$\frac{\partial F}{\partial \vec{A}} = \left(\frac{\partial F}{\partial A_1}, \frac{\partial F}{\partial A_2}, \dots, \frac{\partial F}{\partial A_d} \right)^t = \vec{0} \quad \frac{\partial F}{\partial \lambda} = 0$$