

work-010

PCAの簡単な導出を行う。データ集合 $\{\mathbf{x}_n\} \ n = 1, \dots, N$ を考える。 \mathbf{x} は D 次元のデータ(特徴量が D 個)であり、 $D < N$ を仮定する。

$$\mathbf{x}_i = x_{i1}\mathbf{e}_1 + x_{i2}\mathbf{e}_2 + \dots + x_{iD}\mathbf{e}_D$$

ここで、 $\{\mathbf{e}_j\}$ は正規直交基底とする。以下で第一主成分を求める。

まず、 D 次元単位ベクトル \mathbf{u}_1 を考える。第一主成分とは、「データの分散が最大となる方向」であるので、 $\mathbf{u}_1^T \mathbf{x}$ の分散が最大となるような \mathbf{u}_1 を求めればよい。

1. $\mathbf{u}_1^T \mathbf{x}$ はどのような量か。説明せよ。

分散は

$$\text{Var}(\mathbf{u}_1^T \mathbf{x}) = \frac{1}{N} \sum_i^N \{\mathbf{u}_1^T \mathbf{x}_i - \mathbf{u}_1^T \bar{\mathbf{x}}\}^2 = \mathbf{u}_1^T \boldsymbol{\Sigma} \mathbf{u}_1$$

となる。

ここで、 $\bar{\mathbf{x}}$ 、 $\boldsymbol{\Sigma}$ はそれぞれ \mathbf{x} の平均、共分散行列を表す。

2. $\bar{\mathbf{x}}$ 、 $\boldsymbol{\Sigma}$ を書け。
3. $\text{Var}(\mathbf{u}_1^T \mathbf{x}) = \mathbf{u}_1^T \boldsymbol{\Sigma} \mathbf{u}_1$ の計算過程を書け。

$\text{Var}(\mathbf{u}_1^T \mathbf{x})$ を最大にする \mathbf{u}_1 を求める。ラグランジュの未定乗数法を用いると \mathbf{u}_1 の満たすべき条件式が以下になることがわかる。

$$\boldsymbol{\Sigma} \mathbf{u}_1 = \lambda_1 \mathbf{u}_1$$

4. ラグランジュの未定乗数法を用いて上式を導け。
5. 求めた \mathbf{u}_1 の各成分はどのような意味を持つか？
6. \mathbf{u}_1 方向の分散を求めよ。

以上が、第一主成分を求める手続きである。上記のような方法を繰り返すことで、第二、第三主成分も同様に求めることができる。

求めた主成分を用いて各 \mathbf{u}_i 方向にデータ \mathbf{x}_n を射影する。具体的には以下のように表現される。

$$y_{n1} = \boldsymbol{u}_1^T \boldsymbol{x}_n$$

$$y_{n2} = \boldsymbol{u}_2^T \boldsymbol{x}_n$$

$$\vdots$$

$$y_{ni} = \boldsymbol{u}_i^T \boldsymbol{x}_n$$

$$\vdots$$

$$y_{nd} = \boldsymbol{u}_d^T \boldsymbol{x}_n$$

行列とベクトルで表現すれば、

$$\boldsymbol{y}_n = \boldsymbol{U}^T \boldsymbol{x}_n$$

となる。

7. 上式はどのような意味を持つか。説明せよ。