

V-3-3. Mahalanobis' Distance

In the multivariate analysis, we often discuss similarity of among data. Generally, we consider similarity as distance between data in multidimensional space. If the similarity of two data is high the point of the data will be plotted in the neighborhood in multidimensional space. However, when they has correlation and variance we have to modify the data considering the correlation and variance.

When the data a and b is expressed as vector \mathbf{a} and \mathbf{b} ,

$$\mathbf{a} = \begin{pmatrix} a_1 \\ a_2 \\ a_3 \end{pmatrix}, \quad \mathbf{b} = \begin{pmatrix} b_1 \\ b_2 \\ b_3 \end{pmatrix}$$

The distance is $|\mathbf{a} - \mathbf{b}|$

$$\begin{aligned} \mathbf{a} - \mathbf{b} &= \begin{pmatrix} a_1 - b_1 \\ a_2 - b_2 \\ a_3 - b_3 \end{pmatrix} \\ |\mathbf{a} - \mathbf{b}| &= \sqrt{(\mathbf{a} - \mathbf{b})^T (\mathbf{a} - \mathbf{b})} \\ &= \sqrt{\begin{pmatrix} a_1 - b_1 & a_2 - b_2 & a_3 - b_3 \end{pmatrix} \begin{pmatrix} a_1 - b_1 \\ a_2 - b_2 \\ a_3 - b_3 \end{pmatrix}} \\ &= \sqrt{(a_1 - b_1)^2 + (a_2 - b_2)^2 + (a_3 - b_3)^2} \end{aligned}$$

When the vectors are transformed by matrix \mathbf{A} . The distance between \mathbf{Aa} and \mathbf{Ab} is as follow.

$$\begin{aligned} |\mathbf{Aa} - \mathbf{Ab}| &= |\mathbf{A}(\mathbf{a} - \mathbf{b})| \\ &= \sqrt{(\mathbf{A}(\mathbf{a} - \mathbf{b}))^T \mathbf{A}(\mathbf{a} - \mathbf{b})} \\ &= \sqrt{\begin{pmatrix} a_1 - b_1 & a_2 - b_2 & a_3 - b_3 \end{pmatrix} \mathbf{A}^T \mathbf{A} \begin{pmatrix} a_1 - b_1 \\ a_2 - b_2 \\ a_3 - b_3 \end{pmatrix}} \\ |\mathbf{A}(\mathbf{a} - \mathbf{b})|^2 &= \begin{pmatrix} a_1 - b_1 & a_2 - b_2 & a_3 - b_3 \end{pmatrix} \mathbf{A}^T \mathbf{A} \begin{pmatrix} a_1 - b_1 \\ a_2 - b_2 \\ a_3 - b_3 \end{pmatrix} \\ |(\mathbf{a} - \mathbf{b})\mathbf{A}|^2 &= (\mathbf{a} - \mathbf{b})^T \mathbf{A}^T \mathbf{A} (\mathbf{a} - \mathbf{b}) \end{aligned}$$

Here we denote $\mathbf{a} - \mathbf{b} = \mathbf{x}$

$$\mathbf{x}^T \mathbf{A}^T \mathbf{A} \mathbf{x}$$

When

$$\mathbf{A} = \boldsymbol{\Sigma}^{-\frac{1}{2}}$$

$\boldsymbol{\Sigma}$: variance covariance matrix

$$\begin{aligned}
|\mathbf{Ax}| &= \sqrt{\left(\boldsymbol{\Sigma}^{-\frac{1}{2}}\mathbf{x}\right)^T \boldsymbol{\Sigma}^{-\frac{1}{2}}\mathbf{x}} \\
&= \sqrt{\mathbf{x}^T \left(\boldsymbol{\Sigma}^{-\frac{1}{2}}\right)^T \boldsymbol{\Sigma}^{-\frac{1}{2}}\mathbf{x}} = \sqrt{\mathbf{x}^T \boldsymbol{\Sigma}^{-\frac{1}{2}}\boldsymbol{\Sigma}^{-\frac{1}{2}}\mathbf{x}} = \sqrt{\mathbf{x}^T \boldsymbol{\Sigma}^{-1}\mathbf{x}} \\
&\quad (\because \boldsymbol{\Sigma}^{-\frac{1}{2}} \text{ is symmetric. } \left(\boldsymbol{\Sigma}^{-\frac{1}{2}}\right)^T = \boldsymbol{\Sigma}^{-\frac{1}{2}})
\end{aligned}$$

This is Mahalanobis' Distance $D_{\mathbf{a}-\mathbf{b}} = \sqrt{(\mathbf{a} - \mathbf{b})^T \boldsymbol{\Sigma}^{-1}(\mathbf{a} - \mathbf{b})}$
($\boldsymbol{\Sigma}$: variance covariance matrix)

Formula 70

Mahalanobis' Distance is compensated distance by variance and covariance and Mahalanobis distance is useful tool to consider similarity of data when the data has correlation.