

学会名	The Conference of Data Science, Statistics & Visualisation (DSSV 2019) August 13-15, 2019 in Kyoto, Japan.
演題名	New Sparse Modeling of Sample Mahalanobis Distance
発表者	Yasuyuki Kobayashi
内容	<p>Sparse modeling, such as Least Absolute Selection and Shrinkage Operator (LASSO) for regression has gained interest in variable selection to extract the essential data variables and prevent over-learning problems. Therefore, sparse modeling has also been applied to study the anomaly distance (AD). Thus far, only a sample covariance matrix S of learning samples x has been made sparse, for example, by applying graphical LASSO. However, the AD, such as the sample Mahalanobis distance (MD), of test sample y was not made sparse. Hence, this study was focused on making the AD of test sample y sparse.</p> <p>In principle, ordinal sample MD D^2 is given by $D^2 = (y - \bar{x})^T S^{-1} (y - \bar{x}) = z^T z$, where \bar{x} is the mean of the learning samples, and z is the studentized score vector (SSV) of y, i.e., z is the solution of linear equation $y - \bar{x} = S^{1/2} z$.</p> <p>I propose a new kind of sparse MD, \hat{D}^2, given by $\hat{D}^2 = \hat{z}^T \hat{z}$, where \hat{z} is the sparse solution of the equation obtained by applying the coordinate-descent method to solve LASSO. This sparse MD cancels the unstable effect of numerical error on the sample MD as follows.</p> <p>When learning samples x follow the p-variate normal distribution with population eigenvalues λ such that one $\lambda_0 = 0$ and the other $\lambda > 0$ at the Monte Carlo simulation, sample eigenvalue l_0 of S corresponding to λ_0 becomes slightly positive under the influence of the numerical error, and D^2 becomes unstable owing to l_0. Subsequently, distributions of the element corresponding to l_0 of the SSV of test sample y were simulated as $a(y)$, $b(y)$, and $c(y)$ for the ordinal, ridge, and sparse MDs, respectively. Here, $a(y) = ((y - \bar{x}) \cdot v_0) / l_0$, $b(y) = ((y - \bar{x}) \cdot v_0) / (\sqrt{l_0} + \rho)$, and $c(y) = \hat{z}_l(0)$, where both x and y follow the same normal distribution with dimensionality $p=7$, v_0 is the sample eigenvector corresponding to l_0, regularizing constant $\rho = 10^{-30} \cong l_0$, and $\hat{z}_l(0)$ is the element corresponding to l_0 of \hat{z}. The result shows that the sparse MD is useful for numerical computing. $a(y)$ has a broad distribution because of a numerical error, and $b(y)$ has a narrower distribution around $y=0$ than $a(y)$. However, $c(y)$ degenerates at $y=0$ correctly as $\lambda_0=0$, i.e., the effect of the numerical error is removed.</p> <p>However, for each of the other elements of the SSV of y, all the three distances show the same distribution.</p>