



	The Conference of Data Science, Statistics & Visualisation (DSSV 2019) August 13-15, 2019 in Kyoto, Japan.
	New Sparse Modeling of Sample Mahalanobis Distance
	Yasuyuki Kobayashi
	<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 <math>S</math> of learning samples <math>x</math> has been made sparse, for example, by applying graphical LASSO. However, the AD, such as the sample Mahalanobis distance (MD), of test sample <math>y</math> was not made sparse. Hence, this study was focused on making the AD of test sample <math>y</math> sparse.</p> <p>In principle, ordinal sample MD <math>D^2</math> is given by <math>D^2 = (y - \bar{x})^T S^{-1} (y - \bar{x}) = z^T z</math>, where <math>\bar{x}</math> is the mean of the learning samples, and <math>z</math> is the studentized score vector (SSV) of <math>y</math>, i.e., <math>z</math> is the solution of linear equation <math>y - x = S^{1/2} z</math>.</p> <p>I propose a new kind of sparse MD, <math>\hat{D}^2</math>, given by <math>\hat{D}^2 = \hat{z}^T \hat{z}</math>, where <math>\hat{z}</math> 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 <math>x</math> follow the <math>p</math>-variate normal distribution with population eigenvalues such that one <math>\lambda_0 = 0</math> and the other <math>&gt; 0</math> at the Monte Carlo simulation, sample eigenvalue <math>l_0</math> of <math>S</math> corresponding to <math>\lambda_0</math> becomes slightly positive under the influence of the numerical error, and <math>D^2</math> becomes unstable owing to <math>l_0</math>. Subsequently, distributions of the element corresponding to <math>l_0</math> of the SSV of test sample <math>y</math> were simulated as <math>a(y)</math>, <math>b(y)</math>, and <math>c(y)</math> for the ordinal, ridge, and sparse MDs, respectively. Here, <math>a(y) = ((y - \bar{x})^T v_0) / (l_0)</math>, <math>b(y) = ((y - \bar{x})^T v_0) / (l_0 + \epsilon)</math>, and <math>c(y) = \hat{z}_0^T((y - \bar{x})^T v_0)</math>, where both <math>x</math> and <math>y</math> follow the same normal distribution with dimensionality <math>p = 7</math>, <math>v_0</math> is the sample eigenvector corresponding to <math>l_0</math>, regularizing constant <math>\epsilon = 10^{-30} l_0</math>, and <math>\hat{z}_0^T((y - \bar{x})^T v_0)</math> is the element corresponding to <math>l_0</math> of <math>\hat{z}</math>. The result shows that the sparse MD is useful for numerical computing. <math>a(y)</math> has a broad distribution because of a numerical error, and <math>b(y)</math> has a narrower distribution around <math>y = 0</math> than <math>a(y)</math>. However, <math>c(y)</math> degenerates at <math>y = 0</math> correctly as <math>\lambda_0 = 0</math>, i.e., the effect of the numerical error is removed.</p> <p>However, for each of the other elements of the SSV of <math>y</math>, all the three distances show the same distribution.</p>