

Multivariate data analysis based on two location vectors and two scatter matrices

Hannu Oja
Tampere School of Public Health
FI-33014 University of Tampere, FINLAND
Email: hannu.oja@uta.fi

The regular sample mean vector and covariance matrix are popular tools to describe location and scatter of a multivariate data cloud. Weighted mean vectors and covariance matrices, \mathbf{M} and \mathbf{S} location and scatter estimates, for example, yield often less efficient but more robust estimates of the population mean vector and covariance matrix in the multivariate normal and elliptic case. In case of skew or nonelliptical distributions, different location and scatter estimates usually estimate different population quantities.

In this talk we propose that two location vectors and two scatter matrices, suitably chosen, should be used together in the analysis of non-elliptical data. We show how these can be used for a description of multivariate data cloud (location, scatter, skewness and kurtosis). Tests for multivariate normality and ellipticity may then be based on skewness and kurtosis statistics based on pairs of location and scatter statistics. Two scatter matrices together give an invariant coordinate system (ICS) and then the transformation-retransformation (TR) technique can be used for further analyses. Two scatter matrices with the so called independence property can be used to find independent components in the independent component analysis (ICA). Invariant coordinate system helps in hunting for clusters and outliers, and it can be used in dimension reduction as well. Several examples are given.

The talk is based on joint work and numerous discussions with Frank Critchley, Jan Erikson, Anna-Maija Kankainen, Visa Koivunen, Klaus Nordhausen, Esa Ollila, Davy Paindaveine, Seija Sirkiä, Sara Taskinen and David Tyler (among others).