

Unsupervised Multivariate Analyses for NMR-based Metabolic Profiling

Product used: Nuclear Magnetic Resonance (NMR)

Recently, NMR metabolic profiling (NMR-MP) techniques, which are useful for systematically analyzing metabolites in an organism (metabolome), have been widely used in post-genome studies and quality controls for agricultural and biotechnological products. NMR-MP is characterized as a non-target metabolome analysis. It is difficult to interpret results directly from data in non-target analysis; therefore, multivariate analyses are applied to explore characteristic information from multivariate data. This is referred to as “data-mining”.

It is important to choose appropriate multivariate analyses that correspond to the goal of your study. A few example analyses are dimension reduction, discrimination analysis, regression, modeling, and prediction. In this note, NMR-MP with principal component analysis (PCA) and hierarchical cluster analysis (HCA) that are unsupervised multivariate analyses* are introduced.

Materials and Methods

Polar metabolites in six commercially available sprouts (radish, broccoli, mustard, alfalfa, pea, and okra) were extracted by a modified Bligh-Dyer method. The extracts were then dried with a centrifugal evaporator. Dried metabolite mixtures were resuspended with the NMR solvent (a 100 mM deuterated potassium phosphate buffer, pH = 7.0). ¹H homonuclear 2D *J*-resolved spectra were recorded and their projections were used for multivariate analyses (Figure 1). There were less spectral overlapping in X (*F*₂) projections of the 2D *J*-resolved spectra than conventional ¹H single pulse spectra because there were no splitting from *J*_{HH} [1]. Three sample replicates were performed.

Bucket integrations were performed on the series of spectra (0.02 ppm width each, between 0.5-9.0 ppm). At the time, variables integrals of which are equivalent to those of noise were removed before multivariate analyses. As a result, the matrix (multivariate for multivariate analyses) consisted of 100 variables * 18 samples were obtained. PCA and HCA were performed with free statistical functions in “R”.

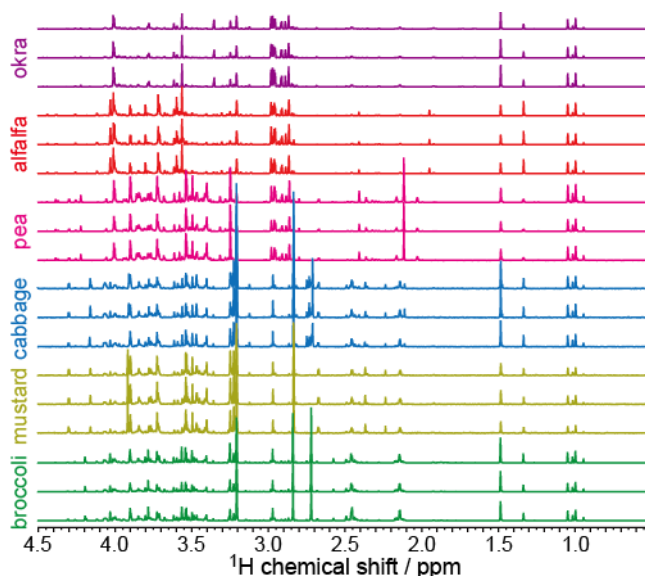


Figure 1. Series of ¹H spectra of polar metabolites in sprouts. X (*F*₂) projection of 2D *J*-resolved spectra were used. NMR spectra were recorded continuously on a JNM-ECZ400S spectrometer equipped with a ROYAL probe and ASC30 auto sample changer.

Applying PCA

Principal components (PCs) are linear combinations of each variable in a data set that result in a maximizing of the variances. The result of this process is that almost all data can be explained using far less dimensions of data. Formula 1 shows a relation between scores, loadings, and original variables. Normally, a centering and scaling of each variable are performed before PCA.

original variables (centering and scaling)

$$\underset{\text{score}}{t_{ij}} = \underset{\text{loadings}}{\bar{x}_{i1}w_{j1}} + \bar{x}_{i2}w_{j2} + \cdots + \bar{x}_{in}w_{jn} \quad (1)$$

Here, subscripts *i* represents *i*th sample, and *j* represents *j*th variable or PC. We can interpret interesting information (ex. characteristic metabolites in particular sample group) from scores and loadings. In this case, eigenvectors were used as loadings.

The score and loading plots in PCA are shown in **Figure 2**. Brassicaceous sprouts including radish, broccoli, and mustard had negative scores in PC1, and others had positive scores (**Figure 2(a)**). In loading plot, asparagine had positive value, on the other hand, glutamine had negative (**Figure 2(b)**).

It is known that amino acids which have an amino group on their side chain such as asparagine and glutamine have an important role in nitrogen metabolism in young plants [2]. The ratio of these amino acids are different by plant species and growing environment. For example, it was reported that brassicaceous and leguminous plants used glutamine and asparagine as the major amino acid has amino group on side chain, respectively [3]. It was also reported that plants synthesize asparagine more than glutamine when they cannot photosynthesize because C/N ratio is smaller in the former than the latter [4].

Applying HCA

HCA is able to summarize relations between each sample and also variable as a dendrogram. Dendrograms are generated according to distance between data and clusters. HCA with heat map visually shows relations between samples and variables (*i.e.* metabolites). While Euclidean distance is one of the most common definition of distance (formula 2), correlation distance (formula 3) which is based on Pearson's correlation coefficient is used in the analysis in this note.

$$dist_{\text{euclidean}}(A, B) = \left[\sum_{i=1}^n (a_i - b_i)^2 \right]^{1/2} \quad (2)$$

$$dist_{\text{corr}}(A, B) = (1 - \text{corr}(A, B)) / 2 \quad (3)$$

Here, a_i and b_i are i th variable in A and B, respectively, and $\text{corr}(A, B)$ is Pearson's correlation coefficient.

HCA with the heat map is shown in **Figure 3**. Columns and rows in the heat map represent variables and samples, respectively. Dendrograms on left- and upper-side show distance between variables and samples hierarchically. In the dendrogram of sample, brassicaceous sprouts were discriminated from the others. In the dendrogram of variable, variables from the same metabolite and variables from different metabolites which had similar behavior built clusters. The heat map was replaced according to dendrograms; Thus, we can visualize which metabolites were present in a particular sample group. For example, sulforaphane, which is an isothiocyanate, is concentrated more heavily in broccoli, and sinigrin, which is a glucosinolate, is found in higher concentrations in mustard [5].

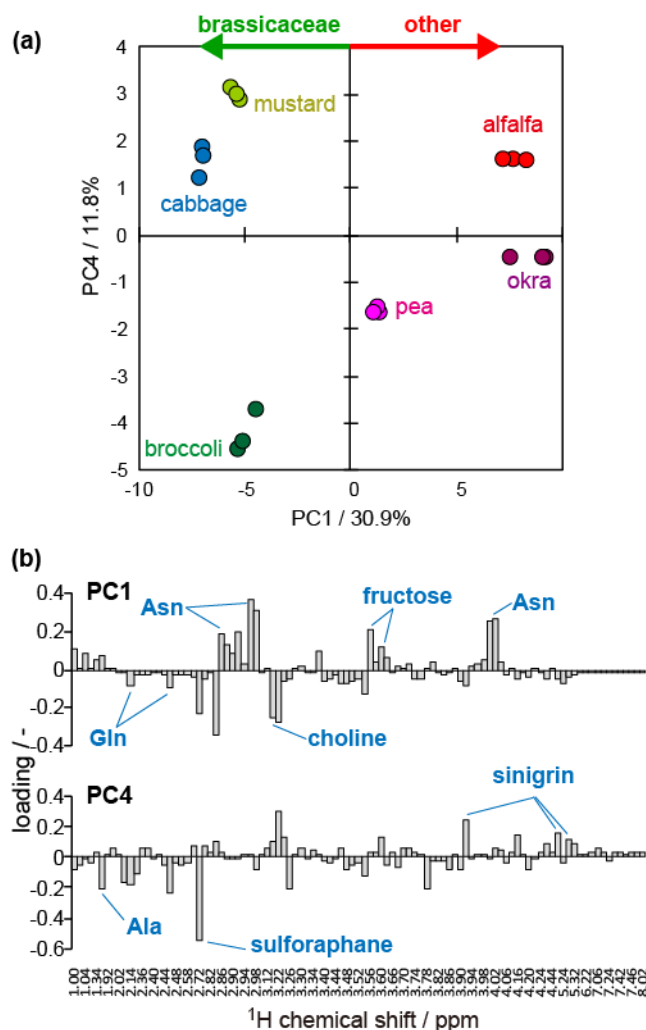


Figure 2. PCA of polar metabolites in sprouts. (a) score plot. (b) loading plot. PCA was performed after the matrix was normalized by the constant sum method and centering was performed.

