

2013

Utilities for Quantifying Separation in PCA/PLS-DA Scores Plots

Bradley Worley

University of Nebraska-Lincoln, bradley.worley@huskers.unl.edu

Steven M. Halouska

University of Nebraska-Lincoln, halouska@huskers.unl.edu

Robert Powers

University of Nebraska-Lincoln, rpowers3@unl.edu

Follow this and additional works at: <http://digitalcommons.unl.edu/chemistrypowers>

Worley, Bradley; Halouska, Steven M.; and Powers, Robert, "Utilities for Quantifying Separation in PCA/PLS-DA Scores Plots" (2013). *Robert Powers Publications*. 28.

<http://digitalcommons.unl.edu/chemistrypowers/28>

This Article is brought to you for free and open access by the Published Research - Department of Chemistry at DigitalCommons@University of Nebraska - Lincoln. It has been accepted for inclusion in Robert Powers Publications by an authorized administrator of DigitalCommons@University of Nebraska - Lincoln.

Published in final edited form as:

Anal Biochem. 2013 February 15; 433(2): 102–104. doi:10.1016/j.ab.2012.10.011.

Copyright © 2012 Elsevier Inc.

Utilities for Quantifying Separation in PCA/PLS-DA Scores Plots

Bradley Worley, Steven Halouska, and Robert Powers*

Department of Chemistry, University of Nebraska-Lincoln, Lincoln, NE 68588-0304

Abstract

Metabolic fingerprinting studies rely on interpretations drawn from low-dimensional representations of spectral data generated by methods of multivariate analysis such as PCA and PLS-DA. The growth of metabolic fingerprinting and chemometric analyses involving these low-dimensional scores plots necessitates the use of quantitative statistical measures to describe significant differences between experimental groups. Our updated version of the PCAToTree software provides methods to reliably visualize and quantify separations in scores plots through dendrograms employing both nonparametric and parametric hypothesis testing to assess node significance, as well as scores plots identifying 95% confidence ellipsoids for all experimental groups.

Keywords

PCA; PLS-DA; MVA; UPGMA; Hierarchical clustering; Bootstrapping; Statistical hypothesis testing; Mahalanobis distance; Hotelling T^2 distribution; Metabolomics

Introduction

A trademark of metabolomics experiments – more specifically metabolic fingerprinting and non-targeted metabolic profiling studies – is the use of multivariate analysis techniques, most commonly principal components analysis (PCA) and projection to latent structures discriminant analysis (PLS-DA) [1,2]. While these techniques provide low-dimensional representations of complex datasets through visually interpretable scores plots, the task of inferring biologically relevant conclusions from scores plots has been largely based on subjective examinations by expert users. Correspondingly, the continued growth in metabolomics and the associated application of chemometric analysis has created a strong need for a quantitative means to justify conclusions drawn from these scores plots. Towards this goal, we recently described the application of our PCAToTree software to generate metabolic tree diagrams from scores plots and the use of standard bootstrapping techniques to infer the statistical significance of each resulting tree node [3]. This note presents a new set of portable software tools that enhances and improves upon our original methodology. Our updated version of the PCAToTree software provides quantification of scores-space separation using both nonparametric bootstrapping and multivariate Hotelling's T^2 hypothesis testing to generate easily interpretable dendrograms of differences between

© 2012 Elsevier Inc. All rights reserved.

*To whom correspondence should be addressed, Robert Powers, University of Nebraska-Lincoln, Department of Chemistry, 722 Hamilton Hall, Lincoln, NE 68588-0304, rpowers3@unl.edu, Phone: (402) 472-3039, Fax: (402) 472-9402.

Publisher's Disclaimer: This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final citable form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

experimental groups. Notably, the new software is now stand-alone and no longer dependent on PHYLIP (<http://www.phylip.com/>) [4].

Scores plots generated from unsupervised PCA or supervised PLS-DA methods provide visualizable representations of information-rich spectral data by means of dimensionality reduction. In the case of PCA, orthogonal lines of maximum gross variation are found within the data, termed the ‘principal axes’, onto which the input data is transformed [5]. This operation preserves as much original gross variation as possible in the first few transformed dimensions, and reveals separations between experimental groups only when within-group variability is sufficiently less than between-group variability. Alternatively, PLS-DA is a supervised method that guides this transformation informed by between-group variability to better reveal group structure [6,7]. In any case, the resultant two- or three-dimensional scores plot is used to identify spectral features contributing to between-group variability based on separations observed between groups in the scores plot.

The importance placed on interpretation of PCA and PLS-DA scores plots necessitates the use of quantitative procedures to determine the significance of these group separations. However, no *de facto* protocol or metric exists to provide a means of reporting the degree or significance of cluster separation [3,8,9]. Anderson *et al.* used the J_2 criterion [10,11] to assess the quality of resulting scores clusters according to the average within-group and between-group scatters for all groups. However, the J_2 metric only provides an overall estimation of cluster separation without fine-grained information on each pair of groups [11]. A similar problem exists with the related Davies-Bouldin index [12], which chooses a worst-case estimate of cluster overlap as its figure of merit. Dixon *et al.* also comprehensively reported the performances of four cluster separation indices based on modifications of metrics used to validate separation for unsupervised clustering algorithms [13]. Alternatively, our PCAToTree protocol constructs dendrograms from distance matrices based on PCA scores for the PHYLIP software suite using a bootstrapping routine to determine node significance [3,4]. However, it was recently shown that hypothesis testing using a Mahalanobis distance metric and the T^2 and F distributions can provide a statistical means to infer cluster similarity [8], suggesting the possibility of returning p -values for full statistical quantitation of PCA group separations.

Methods

The methods described below were implemented in software using the C programming language with minimal external dependencies, so the programs may be compiled and executed on any modern GNU/Linux distribution.

Probability calculation

Under the assumption that each group in the scores space is distributed as a multivariate normal random variable, the distances between groups may be calculated using the squared Mahalanobis distance metric [14]:

$$D_M^2 = (\mathbf{u}_j - \mathbf{u}_i)^T \mathbf{S}_p^{-1} (\mathbf{u}_j - \mathbf{u}_i)$$

Here, \mathbf{u}_i and \mathbf{u}_j are the p -variate sample means of groups i and j , respectively, and \mathbf{S}_p is the pooled p -by- p variance-covariance matrix, a weighted average of the covariance matrices from groups i and j . The Mahalanobis distance may then be related to a Hotelling’s T^2 statistic by the following scaling [15]:

$$T^2 = \left(\frac{n_i n_j}{n_i + n_j} \right) D_M^2$$

where n_i and n_j are the number of data points in groups i and j , respectively. This T^2 statistic is an extension of the Student's t statistic to hypothesis tests in multiple dimensions, and can be related to an F -distribution by a final scaling [15]:

$$x_F = \frac{n_i + n_j - p - 1}{p(n_i + n_j - 2)} T^2 \sim F(p, n_i + n_j - p - 1)$$

It can be seen from this final relation that evaluation of the complement of the cumulative F -distribution function at x_F yields the p -value for accepting the null hypothesis: that the points in groups i and j are in fact drawn from the same multivariate normal distribution.

Tree generation

The implementation of the tree generation procedure is a classical UPGMA algorithm [16]. When p -values are reported at each branch point, a single tree is generated based on the matrix of Mahalanobis distances between groups. In the case of bootstrapped trees, the groups are randomly resampled with replacement while preserving group size. The desired number of trees is then generated using Euclidean distances between group means. The final tree used to report bootstrap probabilities is built using a Euclidean distance matrix calculated from the original (non-resampled) dataset.

Confidence ellipse calculation

When viewing PCA and PLS-DA scores plots, it is common practice to apply hand-drawn ellipses to inform group membership or to even omit such ellipses entirely. This may lead to inconsistent or erroneous interpretation of experimental results. Instead, the fact that the Mahalanobis distances of a set of p -variate points from their sample mean follow a chi-square distribution having p degrees of freedom [17] may be leveraged to estimate 95% confidence ellipsoids for scores in any number of dimensions. The sample mean \mathbf{u} and covariance matrix \mathbf{S} for each group must first be calculated from its scores space data. Then, the group covariance matrix is decomposed into its eigenvalues and eigenvectors:

$$\mathbf{S} = \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^{-1}$$

where \mathbf{Q} is a p -by- p matrix whose columns are the eigenvectors of \mathbf{S} , and $\mathbf{\Lambda}$ is a diagonal matrix of the corresponding eigenvalues of \mathbf{S} . For the case of two-dimensional scores data, the 95% confidence ellipse for the group follows:

$$\begin{bmatrix} x(t) \\ y(t) \end{bmatrix} = \mathbf{u} + \mathbf{Q} \sqrt{\mathbf{\Lambda} F_{0.95,2}^{-1}} \begin{bmatrix} \cos t \\ \sin t \end{bmatrix}$$

where $F_{0.95,2}^{-1}$ is the value of the inverse chi-square cumulative distribution function at $\alpha = 0.05$ and two degrees of freedom, and the square root is taken element-wise over $\mathbf{\Lambda}$. Similarly, a three-dimensional (3D) confidence ellipsoid may be obtained from the following parametric equation:

$$\begin{bmatrix} x(u, v) \\ y(u, v) \\ z(u, v) \end{bmatrix} = \mathbf{u} + \mathbf{Q} \sqrt{\lambda F_{0.95,3}^{-1}} \begin{bmatrix} \cos u \cos v \\ \cos u \sin v \\ \sin v \end{bmatrix}$$

where the parameters t , u and v are all evaluated on $(0, 2\pi)$. These methods allow for the inclusion of confidence regions onto two- and three-dimensional scores plots that reflect the 95% membership boundaries for each group. The approach assumes normally distributed data. Figure 1 illustrates the inclusion of these group confidence regions in representative PCA and OPLS-DA scores plots [18,19]. The ellipses and ellipsoids clearly define statistically significant class separation and also provide an example where multiple groups actually belong to the same biological classification.

Discussion

Our updated and enhanced PCAtree software package consists of a set of stand-alone C programs that generate dendrograms from PCA/PLS-DA scores, report p -values and bootstrap numbers, and incorporate confidence ellipse/ellipsoids into scores plots. The p -values reported for every pair of distinct groups in a PCA/PLS-DA scores plot provide a truly quantitative means to discuss group separations. We also included support for the generation of dendrograms which use these p -values at each branch point to address the question of tree uniqueness. This eliminated the prior dependency on PHYLIP [4]. The reporting of p -values is complementary to bootstrapping methods in cases of highly overlapped groups, where it provides a more direct, interpretable quantitation of group separation.

The PCAtree software package now uses Mahalanobis distances because this metric is more appropriate for multivariate data. De Maesschalck *et al.* provides an exceptional introduction to the use of Mahalanobis distances with PCA [20]. Specifically, Mahalanobis distances account for different variances in each direction (PC1, PC2, PC3) and is scale-invariant. Moreover, the use of a Mahalanobis distance metric for dendrogram generation includes cluster shape and orientation in the analysis of group separation. Also, Mahalanobis distances calculated between groups in PCA scores space will closely approximate those calculated on the original data while avoiding possibly collinearity of the original variables. This is not true of Mahalanobis distances in PLS-DA scores space, due to the underlying supervision of PLS. These features differ from the Euclidean metric, which is a special case of the Mahalanobis metric with the group covariance matrices equaling the identity. Figure 2 illustrates the differences in dendrogram structure based on the use of Euclidean and Mahalanobis distances determined from the same set of scores.

It is important to note that our software is not a means of inferring the reliability of PCA or PLS-DA models, but only a toolset for quantifying the scores that those models produce. In the case of PCA scores, significance of the principal components used must be inferred based on explained sum of squares or another cross-validation technique [21,22]. PLS-DA models require rigorous cross-validation to ensure model reliability, as they almost always yield perfect separations between the scores of different groups [23]. With that in mind, separations between groups not under discrimination may be due to true experimental differences in PLS-DA scores plots, as opposed to the forced separations between discriminated groups. Thus, interpretation of the results of our PCAtree software must be done with the knowledge of the underlying algorithm's mathematical intent, and only after the model has been validated. While we demonstrated our software using only 2D and 3D scores plot, our software places no restrictions on the number of components or on which

components are used during dendrogram generation and p -value calculation. Any dimensionality or choice of scores may be used with our PCAtoTree software provided all components are suitably validated.

Our updated and enhanced PCAtoTree software package provides novel means of quantifying and visualizing separation significance in PCA and PLS-DA scores plots. Importantly, our new software enables single-step methodologies for generating informative scores plots and dendrograms of experimental groups in *any* studies utilizing PCA or PLS-DA to elucidate group structure in complex datasets, including metabolic fingerprinting and non-targeted metabolic profiling. The tools are distributed under version 3.0 of the GNU General Public License and are freely available at <http://bionmr.unl.edu/pca-utils.php>.

Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

Acknowledgments

The authors would like to acknowledge Teklab Gebregiorgis, Bo Zhang and Shulei Lei for their generous contribution of representative PCA and OPLS-DA scores plots used to develop and test the updated PCAtoTree software. This work was supported in part by funds from the National Institute of Health to (RO1 AI087668, R21 AI087561), from the NIH National Center for Research Resources (P20 RR-17675), by the America Heart Association (0860033Z), and the Nebraska Research Council. The research was performed in facilities renovated with support from the National Institutes of Health (NIH, RR015468-01).

References

1. Gebregiorgis T, Powers R. Application of NMR Metabolomics to Search for Human Disease Biomarkers. *Comb. Chem. High Throughput Screening*. 2012; 15:595–610.
2. Zhang B, Powers R. Using NMR-based metabolomics to study the regulation of biofilm formation. *Future Med. Chem.* 2012; 4:1273–1306. [PubMed: 22800371]
3. Werth MT, Halouska S, Shortridge MD, Zhang B, Powers R. Analysis of metabolomic PCA data using tree diagrams. *Anal. Biochem.* 2010; 399:58–63. [PubMed: 20026297]
4. Retief JD. Phylogenetic analysis using. *Methods Mol. Biol.* (Totowa, N. J.). 2000; 132:243–258.
5. Jolliffe, IT. *Principal Component Analysis*. New York: Springer; 2002.
6. Barker M, Rayens W. Partial least squares for discrimination. *J. Chemom.* 2003; 17:166–173.
7. Wold S, Sjostrom M, Eriksson L. PLS-regression: a basic tool of chemometrics. *Chemometr. Intell. Lab.* 2001; 58:109–130.
8. Goodpaster AM, Kennedy MA. Quantification and statistical significance analysis of group separation in NMR-based metabolomics studies. *Chemometr. Intell. Lab.* 2011; 109:162–170.
9. Goodpaster AM, Romick-Rosendale LE, Kennedy MA. Statistical significance analysis of nuclear magnetic resonance-based metabolomics data. *Anal. Biochem.* 2010; 401:134–143. [PubMed: 20159006]
10. Anderson PE, Reo NV, DelRaso NJ, Doom TE, Raymer ML. Gaussian binning: a new kernel-based method for processing NMR spectroscopic data for metabolomics. *Metabolomics*. 2008; 4:261–272.
11. Koutroumbas, K.; Theodoridis, S. *Pattern Recognition*. Amsterdam, Boston: Elsevier/Academic Press; 2006.
12. Davies DL, Bouldin DW. A cluster separation measure. *IEEE Trans. Pattern Anal. Mach. Intell.* 1979; 1:224–227. [PubMed: 21868852]
13. Dixon SJ, Heinrich N, Holmboe M, Schaefer ML, Reed RR, Trevejo J, Brereton RG. Use of cluster separation indices and the influence of outliers: application of two new separation indices, the modified silhouette index and the overlap coefficient to simulated data and mouse urine metabolomic profiles. *J. Chemom.* 2009; 23:19–31.

14. Mahalanobis, PC. Proc. Natl. Inst. Sci. Vol. 2. India: 1936. On the generalized distance in statistics; p. 7
15. Mardia, KV.; Kent, JT.; Bibby, JM. Multivariate analysis. London; New York: Academic Press; 1979.
16. Sokal C, Michener C. A statistical method for evaluating systematic relationships. University of Kansas Science Bulletin. 1958; 38:30.
17. Hotelling H. The generalization of Student's ratio. Annals of Mathematical Statistics. 1931; 2:360–378.
18. Chaika NV, Gebregiorgis T, Lewallen ME, Purohit V, Radhakrishnan P, Liu X, Zhang B, Mehla K, Brown RB, Caffrey T, Yu F, Johnson KR, Powers R, Hollingsworth MA, Singh PK. MUC1 mucin stabilizes and activates hypoxia-inducible factor 1 alpha to regulate metabolism in pancreatic cancer. Proc. Natl. Acad. Sci. U S A. 2012; 109:13787–13792. [PubMed: 22869720]
19. Halouska S, Fenton RJ, Barletta RG, Powers R. Predicting the in Vivo Mechanism of Action for Drug Leads Using NMR Metabolomics. ACS Chem. Biol. 2012; 7:166–171. [PubMed: 22007661]
20. De Maesschalck R, Jouan-Rimbaud D, Massart DL. The Mahalanobis distance. Chemometr. Intell. Lab. 2000; 50:1–18.
21. Eastment HT, Krzanowski WJ. Cross-Validatory Choice of the Number of Components from a Principal Component Analysis. Technometrics. 1982; 24:73–77.
22. Krzanowski WJ. Cross-Validation in Principal Component Analysis. Biometrics. 1987; 43:575–584.
23. Kjeldahl K, Bro R. Some common misunderstandings in chemometrics. J. Chemom. 2010; 24:558–564.

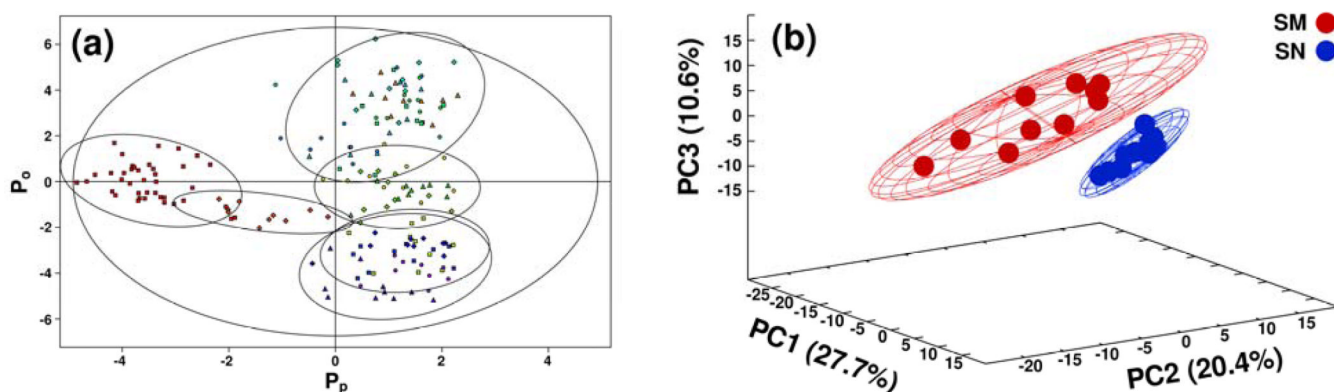


Figure 1.

(a) 2D OPLS-DA scores plot illustrating 95% confidence ellipses for data having one predictive and one orthogonal PLS component. The symbol shape and color of each point correspond to the groups in Figure 2. Discrimination in the first component is between wild-type and antibiotic-treated *Mycobacterium smegmatis*, and separations along the second component indicate metabolic differences between various antibiotic treatments. The antibiotics cluster together based on a shared biological target (cell wall synthesis, mycolic acid biosynthesis, or transcription, translation and DNA supercoiling). Three compounds of unknown *in vivo* activity were shown to cluster together with inhibitors of cell wall synthesis inferring a potential biological target. Interestingly, the *M. smegmatis* strain is resistant to ampicillin resulting in the ampicillin-treated cells clustering closer to untreated cells. (b) 3D PCA scores plot with superimposed 95% confidence ellipsoids drawn as meshes containing group points. The ellipses and ellipsoids define the statistical significance of class separation and provide an illustration where two groups actually belong to the same biological classification. Group ‘SN’ refers to mock-transfected pancreatic cancer cells grown as a control group, while ‘SM’ refers to MUC1-overexpressing pancreatic cancer cells. Separations in scores space relate to metabolic differences in pancreatic cancer due to MUC1 overexpression.

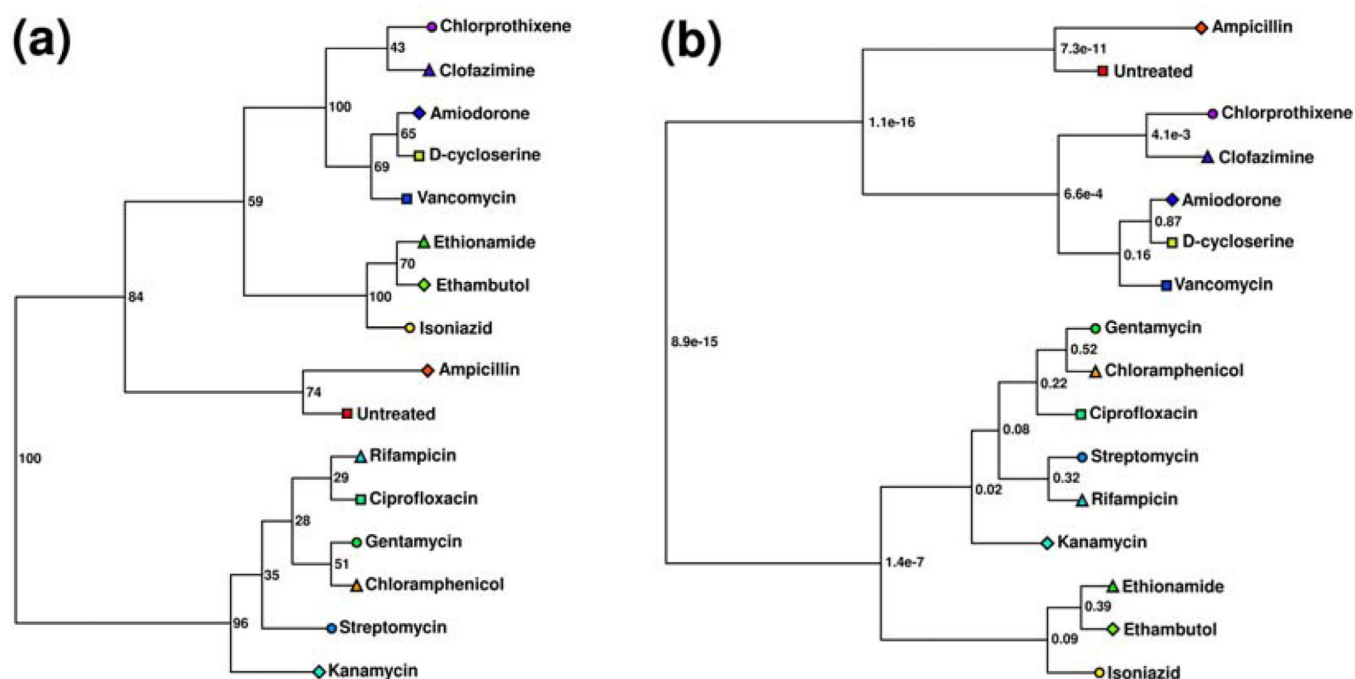
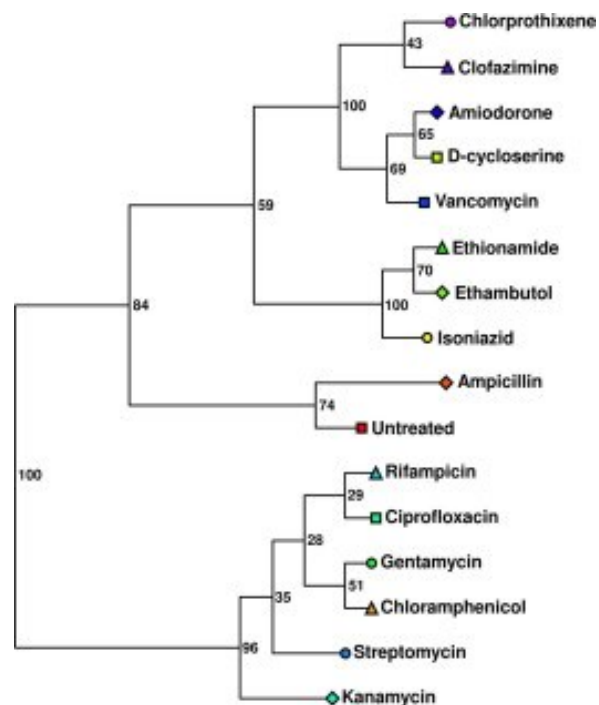


Figure 2.

(a) Dendrogram generated using Euclidean distances between group means from the OPLS-DA scores in Figure 1(a). Bootstrap statistics reported at each branch are for 5,000 bootstrap iterations. (b) Dendrogram generated from identical scores using Mahalanobis distances, with p -values for the null hypothesis reported at each branch.



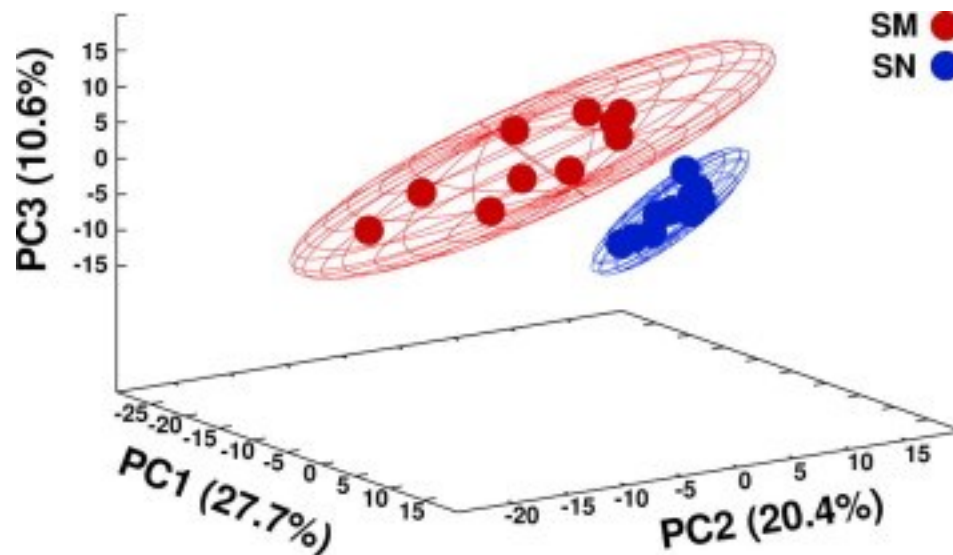
Supplementary Figure 1. Dendrogram generated using Euclidean distances between group means from the OPLS-DA scores in Figure 1(a). Bootstrap statistics reported at each branch are for 5,000 bootstrap iterations.

Bradley Worley, Steven Halouska, Robert Powers

Utilities for quantifying separation in PCA/PLS-DA scores plots

Analytical Biochemistry, Volume 433, Issue 2, 2013, 102–104

<http://dx.doi.org/10.1016/j.ab.2012.10.011>



Supplementary Figure 2. 3D PCA scores plot with superimposed 95% confidence ellipsoids drawn as meshes containing group points. The ellipsoids define the statistical significance of class separation and provide an illustration where two groups actually belong ...

Bradley Worley, Steven Halouska, Robert Powers

Utilities for quantifying separation in PCA/PLS-DA scores plots

Analytical Biochemistry, Volume 433, Issue 2, 2013, 102–104

<http://dx.doi.org/10.1016/j.ab.2012.10.011>