

# ***Barycentric Discriminant Analysis (BADIA)***

Hervé Abdi\* & Lynne J. Williams\*\*

\**The University of Texas at Dallas*    \*\**University of Western Ontario*

Barycentric discriminant analysis (BADIA) generalizes discriminant analysis, and like discriminant analysis, it is performed when measurements made on some observations are combined to assign these observations, or new observations, to *a priori* defined categories. For example, BADIA can be used (a) to assign people to a given diagnostic group (e.g., patients with Alzheimer’s disease, patients with other dementia, or people aging without dementia) on the basis of brain imaging data or psychological tests (here the *a priori* categories are the clinical groups), (b) to assign wines to a region of production on the basis of several physical and chemical measurements (here the *a priori* categories are the regions of production), (c) to use brain scans taken on a given participant to determine what type of object (e.g., a face, a cat, a chair) was watched by the participant when the scans were taken (here the *a priori* categories are the types of object), or (d) to use DNA measurements to predict whether a person is at risk for a given health problem (here the *a priori* categories are the types of health problem).

BADIA is more general than standard discriminant analysis because it can be used in cases for which discriminant analysis cannot be used. This is the case, for example, when there are more variables than observations or when the measurements are categorical.

---

<sup>0</sup>Cite as Abdi, H. & Williams, L.J. (2010). Barycentric discriminant analysis (BADIA). In N.J. Salkind, D.M., Dougherty, & B. Frey (Eds.): *Encyclopedia of Research Design*. Thousand Oaks (CA): Sage. pp. 64–75.

Address correspondence to  
Hervé Abdi  
Program in Cognition and Neurosciences, MS: Gr.4.1,  
The University of Texas at Dallas,  
Richardson, TX 75080-, USA  
**E-mail:** herve@utdallas.edu    <http://www.utd.edu/~herve>

BADIA is a class of methods that all rely on the same principle: Each category of interest is represented by the *barycenter* of its observations (i.e., the weighted average; the barycenter is also called the *center of gravity* of the observations of a given category), and a generalized principal components analysis (GPCA) is performed on the category by variable matrix. This analysis gives a set of discriminant factor scores for the categories and another set of factor scores for the variables. The original observations are then projected onto the category factor space, providing a set of factor scores for the observations. The distance of each observation to the set of categories is computed from the factor scores, and each observation is assigned to the closest category. The comparison between the *a priori* and *a posteriori* category assignments is used to assess the quality of the discriminant procedure. The prediction for the observations that were used to compute the barycenters is called the *fixed-effect* prediction. Fixed-effect performance is evaluated by counting the number of correct and incorrect assignments and storing these numbers in a confusion matrix. Another index of the performance of the fixed-effect model—equivalent to a squared coefficient of correlation—is the ratio of category variance to the sum of the category variance plus the variance of the observations within each category.

This coefficient is denoted  $R^2$  and is interpreted as the proportion of variance of the observations explained by the categories or as the proportion of the variance explained by the discriminant model. The performance of the fixed-effect model can also be represented graphically as a *tolerance* ellipsoid that encompasses a given proportion (say 95%) of the observations. The overlap between the tolerance ellipsoids of two categories is proportional to the number of mis-classifications between these two categories.

New observations can also be projected onto the discriminant factor space, and they can be assigned to the closest category. When the actual assignment of these observations is not known, the model can be used to *predict* category membership. The model is then called a *random* model (as opposed to the fixed model). An obvious problem, then, is to evaluate the quality of the prediction for new observations. Ideally, the performance of the random-effect

model is evaluated by counting the number of correct and incorrect classifications for new observations and computing a confusion matrix on these new observations. However, it is not always practical or even feasible to obtain new observations, and therefore the random-effect performance is, in general, evaluated using computational cross-validation techniques such as the *jackknife* or the *bootstrap*. For example, a jackknife approach (also called *leave one out*) can be used by which each observation is taken out of the set, in turn, and predicted from the model built on the other observations. The predicted observations are then projected in the space of the fixed-effect discriminant scores. This can also be represented graphically as a *prediction* ellipsoid. A prediction ellipsoid encompasses a given proportion (say 95%) of the new observations. The overlap between the prediction ellipsoids of two categories is proportional to the number of mis-classifications of new observations between these two categories.

The stability of the discriminant model can be assessed by a cross-validation model such as the bootstrap. In this procedure, multiple sets of observations are generated by sampling with replacement from the original set of observations, and the category barycenters are computed from each of these sets. These barycenters are then projected onto the discriminant factor scores. The variability of the barycenters can be represented graphically as a confidence ellipsoid that encompasses a given proportion (say 95%) of the barycenters. When the confidence intervals of two categories do not overlap, these two categories are significantly different.

In summary, BADIA is a Generalized Principal Component Analysis performed on the category barycenters. GPCA encompasses various methods, such as correspondence analysis, biplot, Hellinger distance analysis, discriminant analysis, and canonical variate analysis. For each specific type of GPCA, there is a corresponding version of BADIA. For example, when the GPCA is correspondence analysis, this is best handled with the most well-known version of BADIA: discriminant correspondence analysis (DiCA). Because BADIA is based on GPCA, it can also analyze data tables obtained by the concatenation of blocks (i.e.,

subtables). In this case, the importance (often called the contribution) of each block to the overall discrimination can also be evaluated and represented as a graph.

See also

- Bootstrapping
- Canonical Correlation Analysis
- Correspondence Analysis
- Discriminant Analysis
- Jackknife
- Matrix Algebra
- Principal Components Analysis

## Further readings

- [1] Abdi, H. (2007). Discriminant correspondence analysis (DCA). In N.J. Salkind (Ed.): *Encyclopedia of Measurement and Statistics*. Thousand Oaks (CA): Sage. pp. 270–275.
- [2] Abdi, H, Dunlop, J.P., & Williams, L.J. (2009). How to compute reliability estimates and display confidence and tolerance intervals for pattern classifiers using the Bootstrap and 3-way multidimensional scaling (DISTATIS). *NeuroImage*, **45**, 89–95.
- [3] Efron, B., & Tibshirani, R.J. (1993). *An Introduction to the Bootstrap*. New York: Chapman & Hall.
- [4] Saporta, G., & Niang, N. (2006). Correspondence Analysis and Classification. In M. Greenacre & J. Blasius (Eds.): *Multiple Correspondence Analysis and Related Methods*. Boca raton (FL): Chapman & Hall/CRC. pp. 371–392.