

A Roadmap to Robust Discriminant Analysis of Principal 1 Components (DAPC)

Catherine Cullingham, Rhiannon M. Peery, Joshua M. Miller

NOTICE: This is the peer reviewed version of the following article: Cullingham, C., Peery, R. M., & Miller, J. M. (2023). A roadmap to robust discriminant analysis of principal components. *Molecular Ecology Resources*, 23(3), 519-522, which has been published in final form at <http://dx.doi.org/10.1111/1755-0998.13724>. This article may be used for non-commercial purposes in accordance with Wiley Terms and Conditions for self-archiving.

Permanent link to this version <https://hdl.handle.net/20.500.14078/3431>

License All Rights Reserved

1 **A roadmap to robust Discriminant Analysis of Principal Components (DAPC)**

2 Cullingham, CI¹, Peery RM¹ Miller JMM²

3 ¹Department of Biology, Carleton University, Ottawa, Ontario, Canada

4 ²Biological Sciences, MacEwan University, Edmonton, Alberta, Canada

5 Cullingham ORCID: <https://orcid.org/0000-0002-6715-0674>

6 Peery ORCID: <https://orcid.org/0000-0003-4711-4702>

7 Miller ORCID: <https://orcid.org/0000-0002-4019-7675>

8 Corresponding author: Catherine Cullingham, catherine.cullingham@carleton.ca

9 **Identification of population structure is a common goal for a variety of applications,**
10 **including conservation, wildlife management, and medical genetics. The outcome of these**
11 **analyses can have far reaching implications; therefore, it is important to ensure an**
12 **understanding of the strengths and weaknesses of the methodologies used. Increasing in**
13 **popularity, the Discriminant Analysis of Principal Components (DAPC) method**
14 **incorporates combinations of genetic variables (alleles) into a model that differentiates**
15 **individuals into genetic clusters. However, users may not have a full understanding of how**
16 **to best parameterize the model. In this issue of Molecular Ecology Resources, Thia (2022)**
17 **looks under the hood of the DAPC. Using simulated data, he demonstrates the importance**
18 **of careful parameter selection in developing a DAPC model, what the implications are for**
19 **over-fitting the model, and finally, how best to evaluate the results of DAPC models. This**
20 **work highlights the issues that can arise when over-parameterizing the DAPC model when**
21 **gene flow is high among clusters and provides important guidelines to ensure researchers**
22 **are making conclusions that are biologically relevant.**

23 The identification of clusters and the estimation of their divergence using genetic data has been a
24 pillar of population genetics for over 70 years and the implications of these analyses can have far
25 reaching impacts. For example, in the area of conservation biology, where resources are often
26 limited, incorrect identification of significantly differentiated groups can result in either wasted
27 efforts if groups are unnecessarily split, or inflated population size estimates if true population
28 structure is not identified (Fallon, 2007). Thus, rigorous and accurate assessment of population
29 structure is necessary (Haig et al., 2016). Similar to the maturing of the use of other clustering
30 software (e.g. STRUCTURE, Pritchard, Stephens, & Donnelly, 2000)), as a community we need to
31 establish a set of best practices for more recent software designed to find genetic structure.

32 The program STRUCTURE (Pritchard et al., 2000) has become the standard method for assessing
33 population structure. Concurrent with its popularity, several papers have explored the limits and
34 developed best practices for STRUCTURE (Gilbert et al., 2012; Janes et al., 2017; Lawson, van
35 Dorp, & Falush, 2018). However, as the number of loci included in analyses increases,
36 researchers often note the extensive time requirements as a barrier to its use (Wang, 2022).
37 Given studies using datasets with thousands, to hundreds of thousand loci are now common,
38 more expedient methods are needed. Discriminant analysis of principal components was
39 introduced by Jombart, Devillard, & Balloux (2010) to address not only issues with analysis of
40 large datasets, but also as a method that does not rely on the population genetic model
41 assumptions associated with STRUCTURE. Discriminant analysis of principal components can be
42 used to identify groups when they are unknown, visualize complex population structure, identify
43 genomic regions driving population differences, and test assignment of individuals to clusters.
44 Given this broad range of utility, together with the ease of use, it is now the third most cited
45 method for assessing population structure (Figure 1).

46 Discriminant analysis of principal components involves a two-step process. The first is to use a
47 Principal Component Analysis (PCA) to explain the population level variation among
48 uncorrelated combinations of alleles by generating eigenvectors that summarize the covariance
49 matrix generated from the data. The second is to select components that describe the between
50 cluster variation and use those in a Discriminant Analysis (DA), which builds a model that can
51 predict the population for each individual. The number of components selected from the PCA is
52 critical as selecting too many axes will result in overfitting the model and create an inflated
53 estimate of the amount of differentiation among clusters (Thia 2022). While the developers of

54 the DAPC method indicate the importance of this step, they do not provide a defined set of rules
55 for selecting the components (Jombart & Collins, 2022). To address this issue, Thia (2022)
56 examined simulated data for a five-population, finite-island model under three different levels of
57 migration, low, medium, and high (panmixia). He analyzed the data using DAPC, while selecting
58 different numbers of PC axes to use in the discriminant model. Based on these analyses, he
59 highlights several important issues, and makes clear recommendations to ensure the proper use
60 and interpretation of DAPC.

61 The selection of the number of PC axes has primarily relied on three approaches (Miller,
62 Cullingham, & Peery, 2020): 1) selecting the number of PC axes which explain an arbitrary
63 amount of variance in the dataset (often $\geq 80\%$), 2) the *xvalDapc* function, which uses a training
64 and testing set to find the optimal trade-off between selecting too few, and too many PC axes, or
65 3) the *optim.a.score* function, which looks at the assignment success of individuals based on
66 different numbers of PC axes retained, and selects the number of PCs that maximizes assignment
67 success. The work by Thia (2022) highlights caveats to consider for each of these methods. The
68 first method often leads to inclusion of components which explain individual differences, rather
69 than just including those which capture between population differences. When there is no
70 migration ($F_{ST} = 0.99$ in simulated data), there is little impact on the conclusions made because
71 the majority of variance is between populations. However, when migration is medium to high
72 ($F_{ST} = 0.09$ & 0.0009 , respectively) there is an inflation of how differentiated populations appear
73 (Figure 8 in Thia 2022). The second method suffered from similar issues to the first, where the
74 *xvalDapc* function indicates more PC axes be selected than those that just explained population
75 structure in the simulated datasets, again resulting in an over-fit model. The third method was the
76 most conservative, but again returned inflated values of PC axes to obtain when migration was
77 medium to high. This is highly relevant given a review in 2020 found that the majority of studies
78 using DAPC do not report the number of PC axes retained (64%) and for those that do, the range
79 of reported values (2-600, mean = 57) suggests many studies are selecting far too many axes
80 (Miller et al., 2020).

81 As an alternative to these approaches, Thia (2022) suggests selecting the number of PC axes
82 using the $K-1$ rule from Patterson, Price, & Reich, (2006), where K is the number of genetic
83 clusters. This initially seems like a very conservative approach; however, Thia (2022) clearly
84 demonstrates these components capture the between population differences, and when migration
85 is moderate, provide the best model for predicting population membership. Including additional
86 PC axes can inflate or obscure the true population structure.

87 Of course, many studies do not necessarily know the number of populations represented by their
88 samples, and therefore may be interested in using DAPC to determine the most likely value of K .
89 Thia (2022) recommends using alternative programs for finding K , however, as noted in Miller et
90 al. (2020) most researchers are using more than one method to corroborate a choice of K and
91 several authors use DAPC along with STRUCTURE or other admixture programs. To facilitate
92 decision making when employing DAPC we have generated a decision tree that includes some
93 recommendations made in Thia (2022) together with best practices in population genomic
94 analyses that are tailored to the overall goal of the study (Figure 2). In some instances,
95 conducting a PCA without a DA may be the most appropriate approach to use. As noted by Thia
96 (2022) and in Miller et al., (2020) reporting of decisions and parameters are critical for
97 evaluation of the DAPC analyses itself and for reproducibility.

98 In recognizing the important utility of a PCA based method, we must also mention the Achilles
99 heel of the PCA – no method is robust to poorly curated datasets. The matrix of alleles that is
100 used for the PCA needs to be properly filtered and pruned before analysis. Errors and bias in the
101 data can strongly influence the overall structure of the PC plot (e.g., Patterson et al., 2006),
102 especially when genetic differentiation is low. Missing data, a common feature of the
103 genotyping-by-sequencing method, will also influence the PC plot, sometimes suggesting
104 population structure where none exists (e.g., Yi & Latch, 2022). No matter the study objective,
105 best practices for assessing data integrity in genomic datasets need to be followed prior to
106 examining population structure (Laurie et al., 2010; Shafer et al., 2017).

107 Identifying and characterizing population structure will continue to be an important goal in
108 molecular ecology. To this end, methods that are rapid, consistent, and can define clusters when
109 genetic differentiation is limited are desirable. Discriminant analysis of principal components fits
110 these criteria, yet the parameterization of the model has not been thoroughly addressed in the
111 literature. Thia (2022) address some of these outstanding concerns, and demonstrates it is a
112 powerful tool, but can lead to false conclusions when too many PC axes are included in the
113 model, especially when genetic differentiation is moderate to low. Thia (2022) then goes on to
114 present a list of recommendations (included in Figure 2), that if followed, will allow researchers
115 to arrive at conclusions that are relevant for the biology of their species. We expect with these
116 clear recommendations, the use of this method will continue to increase, and users will benefit
117 from the utility of the tool.

118 **Data accessibility**

119 Data generated from the Web of Science search and the R code for generating Figure 1 can be
120 found at <https://doi.org/10.5683/SP3/UPONUS>.

121

122 **References**

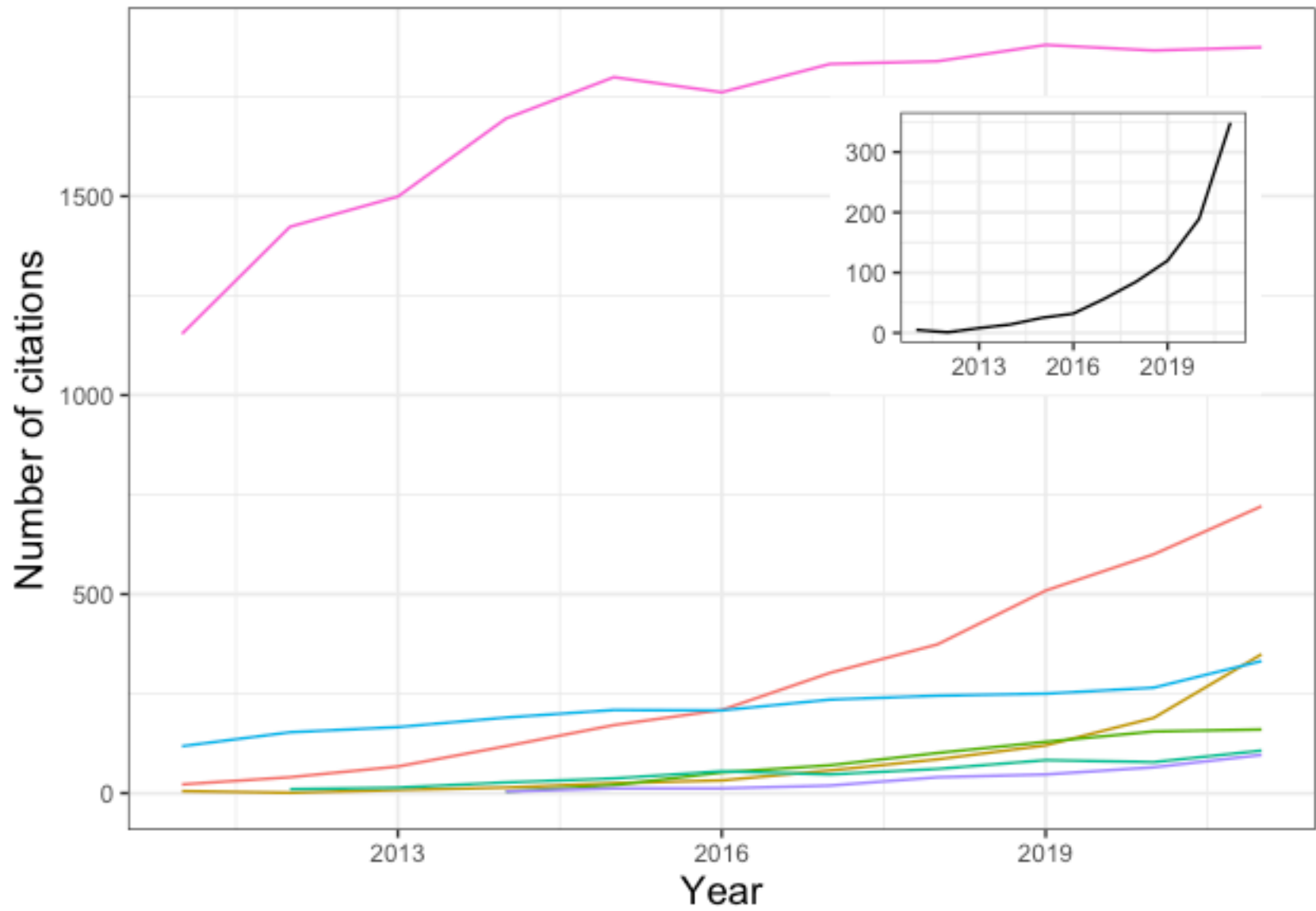
- 123 Fallon, S. M. (2007). Genetic Data and the Listing of Species Under the U.S. Endangered
124 Species Act. *Conservation Biology*, 21(5), 1186–1195.
- 125 Gilbert, K. J., Andrew, R. L., Bock, D. G., Franklin, M. T., Kane, N. C., Moore, J.-S., ... Vines,
126 T. H. (2012). Recommendations for utilizing and reporting population genetic analyses: the
127 reproducibility of genetic clustering using the program structure. *Molecular Ecology*,
128 21(20), 4925–4930.
- 129 Haig, S. M., Miller, M. P., Bellinger, R., Draheim, H. M., Mercer, D. M., & Mullins, T. D.
130 (2016). The conservation genetics juggling act: integrating genetics and ecology, science
131 and policy. *Evolutionary Applications*, 9(1), 181–195.
- 132 Janes, J. K., Miller, J. M., Dupuis, J. R., Malenfant, R. M., Gorrell, J. C., Cullingham, C. I., &
133 Andrew, R. L. (2017). The $K = 2$ conundrum. *Molecular Ecology*.
- 134 Jombart, T., & Collins, C. (2022). A tutorial for Discriminant Analysis of Principal Components
135 (DAPC) using *adegenet 2.1.6* (p. 43). p. 43.
- 136 Jombart, T., Devillard, S., & Balloux, F. (2010). Discriminant analysis of principal components:
137 a new method for the analysis of genetically structured populations. *BMC Genetics*, 11(1),

- 138 94.
- 139 Laurie, C. C., Doheny, K. F., Mirel, D. B., Pugh, E. W., Bierut, L. J., Bhangale, T., ... Weir, B.
140 S. (2010). Quality control and quality assurance in genotypic data for genome-wide
141 association studies. *Genetic Epidemiology*, *34*(6), 591–602.
- 142 Lawson, D. J., van Dorp, L., & Falush, D. (2018). A tutorial on how not to over-interpret
143 STRUCTURE and ADMIXTURE bar plots. *Nature Communications*, *9*(1).
- 144 Miller, J. M., Cullingham, C. I., & Peery, R. M. (2020). The influence of a priori grouping on
145 inference of genetic clusters: simulation study and literature review of the DAPC method.
146 *Heredity* 2020 125:5, *125*(5), 269–280.
- 147 Patterson, N., Price, A. L., & Reich, D. (2006). Population Structure and Eigenanalysis. *PLOS*
148 *Genetics*, *2*(12), e190.
- 149 Pritchard, J. K., Stephens, M., & Donnelly, P. (2000). Inference of population structure using
150 multilocus genotype data. *Genetics*, *155*(2), 945–959.
- 151 Shafer, A. B. A., Peart, C. R., Tusso, S., Maayan, I., Brelsford, A., Wheat, C. W., & Wolf, J. B.
152 W. (2017). Bioinformatic processing of RAD-seq data dramatically impacts downstream
153 population genetic inference. *Methods in Ecology and Evolution*, *8*(8), 907–917.
- 154 Thia, J. (2022). Guidelines for standardising the application of discriminant analysis of principal
155 components to genotype data. *Molecular Ecology Resources*, *XX*, XXX–XXX.
- 156 Wang, J. (2022). Fast and accurate population admixture inference from genotype data from a
157 few microsatellites to millions of SNPs. *Heredity* 2022 129:2, *129*(2), 79–92.
- 158 Yi, X., & Latch, E. K. (2022). Nonrandom missing data can bias Principal Component Analysis
159 inference of population genetic structure. *Molecular Ecology Resources*, *22*(2), 602–611.
- 160

161 Figure captions

162 Figure 1. Total number of citations, calculated via reverse citation search of each program's
163 initial publication in Web-of-Science on Sept. 9, 2022, for popular clustering algorithms. The
164 popularity of DAPC has steadily increased since its publication in 2011 (inset).

165 Figure 2. Decision tree for DAPC analyses designed in draw.io. The first phase is goal setting
166 (shapes with stippled outlines), the second phase determines the model parameters and user
167 inputs (bold outlines), and the final phase tests the model fit (thin, solid outline). The data
168 pruning and reporting phases recommended by (Laurie et al., 2010; Miller et al., 2020; Shafer et
169 al., 2017) are not shown.



— ADMIXTURE — fastSTRUCTURE — smartPCA — STRUCTURE
— DAPC — fineSTRUCTURE — sNMF

