# Technical Support Center: Discriminant Analysis of Principal Components (DAPC)

**Author**: BenchChem Technical Support Team. **Date**: December 2025

| Compound of Interest | |
|---|---|
| Compound Name: | DAPCy |
| Cat. No.: | B8745020 |

Get Quote

Welcome to the technical support center for DAPC analysis. This guide provides troubleshooting information and frequently asked questions (FAQs) to help researchers, scientists, and drug development professionals optimize their experiments, with a specific focus on selecting the optimal number of principal components (PCs).

## Frequently Asked Questions (FAQs)

## Q1: What is Discriminant Analysis of Principal Components (DAPC) and what is its primary application?

Discriminant Analysis of Principal Components (DAPC) is a multivariate statistical method used to identify and describe clusters of genetically related individuals.[1] The method works in two stages. First, the data is transformed using Principal Component Analysis (PCA) to reduce dimensionality. Second, a Discriminant Analysis (DA) is performed on the retained principal components to maximize the separation between groups while minimizing variation within them.[2][3] This makes it an excellent tool for exploring the genetic structure of populations without assuming the data conforms to specific population genetics models, such as Hardy-Weinberg equilibrium.[2]

## Q2: Why is the selection of the number of Principal Components (PCs) a critical step in DAPC?

Choosing the number of PCs to retain is a critical decision that balances information retention against model overfitting.[4]

- Retaining too few PCs: This can lead to underfitting, where not enough genetic variation is captured, potentially obscuring the true structure within the data.

- Retaining too many PCs: This can lead to overfitting. The model may start to capture random noise rather than a true biological signal, leading to unstable and unreliable cluster assignments.[4] This is particularly problematic when the number of variables is much larger than the number of individuals.

The goal is to find the "sweet spot" that captures the meaningful biological variation and discards the noise, leading to a stable and reliable model.

## Q3: What are the primary methods for choosing the optimal number of PCs?

There are two widely accepted, objective methods for determining the optimal number of PCs to retain in a DAPC analysis, both available in the adegenet package in R:

- Cross-Validation: This method, implemented with the xvalDapc function, assesses the predictive power of the DAPC model with varying numbers of PCs.[2][5] It is often considered the most robust approach.

- A-score Optimization: This method, implemented with the optim.a.score function, evaluates the trade-off between discriminatory power and overfitting.[6] It helps identify a DAPC solution that is both stable and discriminative.[7]

## Q4: How does cross-validation work to find the best number of PCs?

Cross-validation objectively identifies the optimal number of PCs by assessing the stability and predictive accuracy of the DAPC.[4] The procedure, executed by the xvalDapc function, involves partitioning the data:

- The data is repeatedly split into two subsets: a training set (typically 90% of the data) and a validation set (the remaining 10%).[4][5]

- A DAPC is performed on the training set for a range of different numbers of retained PCs.

- The group assignments of individuals in the validation set are then predicted based on the DAPC model built from the training set.[2]

- The success of this prediction is measured across many replicates. The optimal number of PCs is the one that provides the highest mean prediction success and, more importantly, the lowest Root Mean Squared Error (RMSE).[5]

## Q5: What is the 'a-score' and how does it optimize PC selection?

The 'a-score' is a metric that measures how well a DAPC model can be successfully re-assigned to its original clusters compared to random clusters.[7][8] It is calculated as the difference between the probability of correct assignment of the true clusters and the probability of correct assignment of randomly permuted clusters.[7]

- An a-score close to 1 indicates a stable and highly discriminating model.

- An a-score close to 0 or lower suggests weak discrimination or an unstable model that is likely overfitted.[7]

The optim.a.score function calculates this score for different numbers of retained PCs. The optimal number of PCs is the one that maximizes the a-score, thus balancing discriminatory power with model stability.[8]

## Troubleshooting Guide

## Problem: The cross-validation plot of mean success is flat, or the RMSE does not show a clear minimum.

This situation can arise when the underlying population structure is very weak or non-existent. If there are no clear genetic clusters in your data, the ability to correctly predict group membership will not improve significantly with an increasing number of PCs, as no meaningful between-group variance can be maximized.

Solution:

- Re-evaluate Prior Clusters: If you are using pre-defined populations, consider whether this grouping is biologically meaningful. You can use the find.clusters function in adegenet to identify clusters based on the data itself and see if this reveals a more robust structure.[1]

- Check the RMSE: The Root Mean Squared Error (RMSE) is often more informative than the mean success rate. Look for the number of PCs that corresponds to the lowest RMSE value, even if the curve is relatively flat.[5]

- Consider Alternative Methods: If DAPC does not reveal a clear structure, it may be that the genetic differentiation is too low to be detected by this method. Consider other population structure analyses or methods that are more suited for detecting subtle differentiation.

## Problem: The xvalDapc function suggests retaining a very low number of PCs (e.g., 1 or 2).

This is not necessarily an error. If the vast majority of the discriminatory power is contained within the first few principal components, then retaining more may only add noise. This is common in datasets with a very strong and simple population structure (e.g., two very distinct species). Trust the cross-validation result, as its purpose is to objectively identify this point.

## Problem: The optim.a.score function runs very slowly.

The a-score optimization can be computationally intensive because it involves permutations and repeated analyses.

Solution:

- Use the smart parameter: The optim.a.score function in adegenet has a smart parameter which is TRUE by default.[8] This uses a faster algorithm that evaluates a subset of evenly distributed PC numbers and interpolates the results using splines to find the optimum.[8][9] Ensure you are using this default setting.

- Reduce the Range: Instead of testing a vast range of PCs (e.g., 1 to 300), first run a preliminary, coarse analysis (e.g., testing every 10th PC) to identify a promising region. Then, perform a second, finer-grained analysis within that smaller range.

# Experimental Protocol: PC Selection via Cross-Validation

This protocol outlines the standard procedure for using the xvalDapc function in the R package adegenet.

Objective: To identify the optimal number of principal components to retain in a DAPC for maximizing the predictive accuracy of group assignment.

Methodology:

- Load Data: Load your genetic data (e.g., from a VCF, genepop, or other file) into R as a genind or genlight object. Ensure your object contains the a priori group or population assignments for each individual.

- Execute Cross-Validation: Use the xvalDapc() function. Specify your data object and the grouping factor. It is recommended to run a sufficient number of replicates for a stable result.

- Interpret the Results: The xvalDapc function returns a list of results and a plot. The key outputs are:

  - Mean Successful Assignment: The proportion of validation individuals correctly assigned to their group for each number of PCs tested.

  - Root Mean Squared Error (RMSE): A measure of the error in predicted assignments.

- Identify Optimal PC Number: The optimal number of PCs is the one associated with the lowest RMSE.[5] The function output will explicitly state this number.

- Perform Final DAPC: Run the final DAPC analysis using the optimal number of PCs identified in the previous step.
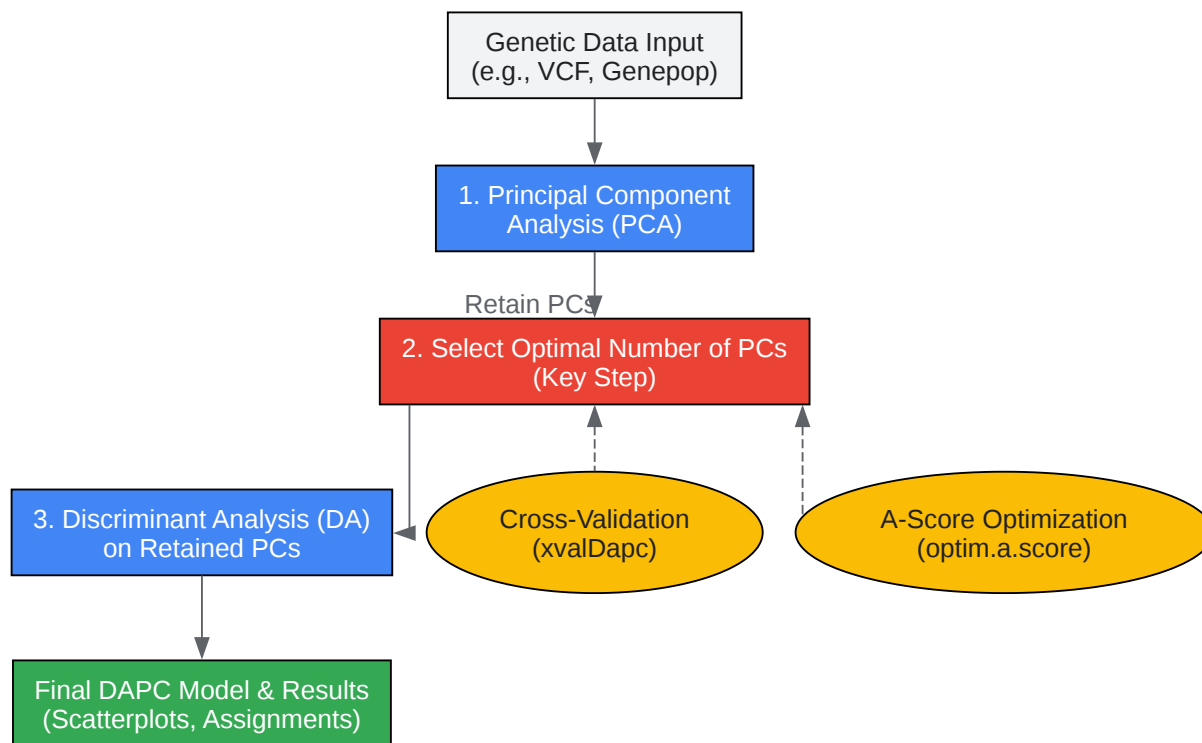
## Data Presentation

The results from the cross-validation can be summarized in a table for clear interpretation and reporting.

| Number of PCs Retained | Mean Success (%) | Standard Deviation | RMSE |
|---|---|---|---|
| 10 | 85.2 | 3.1 | 0.247 |
| 20 | 92.5 | 2.5 | 0.187 |
| 30 | 96.1 | 1.9 | 0.125 |
| 40 | 97.3 | 1.5 | 0.098 |
| 50 | 97.4 | 1.6 | 0.101 |
| 60 | 97.2 | 1.8 | 0.115 |
| 70 | 96.8 | 2.0 | 0.132 |

Table 1: Example output from a DAPC cross-validation analysis. The lowest Root Mean Squared Error (RMSE) is achieved when retaining 40 PCs, which is therefore the optimal number for this analysis.

# Visualizations

## DAPC Workflow Diagram

```
┌──────────────────────────┐
│     Genetic Data Input   │
│   (e.g., VCF, Genepop)   │
└──────────────────────────┘
             │
             ▼
┌──────────────────────────┐
│  1. Principal Component  │
│      Analysis (PCA)      │
└──────────────────────────┘
             │ Retain PCs
             ▼
┌──────────────────────────┐
│ 2. Select Optimal Number │
│   of PCs (Key Step)      │
└──────────────────────────┘
```

Genetic Data Input (e.g., VCF, Genepop)

1. Principal Component Analysis (PCA)

Retain PCs

2. Select Optimal Number of PCs (Key Step)

3. Discriminant Analysis (DA) on Retained PCs

Cross-Validation (xvalDapc)

A-Score Optimization (optim.a.score)

Final DAPC Model & Results (Scatterplots, Assignments)

Model Performance ← ─── ─── Number of Principal Components Retained →

Underfitting (Too few PCs, Information Loss)

Optimal Model (Balanced Fit)

Overfitting (Too many PCs, Noise Captured)

opt_line_start

opt_line_end

Model Fit (e.g., RMSE)

Model Fit (e.g., RMSE)

Model Fit (e.g., RMSE)

Model Fit (e.g., RMSE)

Model Fit (e.g., RMSE)

Click to download full resolution via product page

> **Need Custom Synthesis?**
>
> *BenchChem offers custom synthesis for rare earth carbides and specific isotopiclabeling.*
>
> *Email: info@benchchem.com or Request Quote Online.*

# References

- 1. Discriminant analysis of principal components: a new method for the analysis of genetically structured populations - PMC [pmc.ncbi.nlm.nih.gov]
- 2. Discriminant analysis of principal components (DAPC) [grunwaldlab.github.io]
- 3. RPubs - DAPC [rpubs.com]

- 4. adegenet.r-forge.r-project.org [adegenet.r-forge.r-project.org]

- 5. HTTP redirect [search.r-project.org]

- 6. GitHub - laurabenestan/DAPC: Discriminant Analysis in Principal Components (DAPC) [github.com]

- 7. R: Compute and optimize a-score for Discriminant Analysis of... [search.r-project.org]

- 8. ms.mcmaster.ca [ms.mcmaster.ca]

- 9. raw.githubusercontent.com [raw.githubusercontent.com]

- To cite this document: BenchChem. [Technical Support Center: Discriminant Analysis of Principal Components (DAPC)]. BenchChem, [2025]. [Online PDF]. Available at: [https://www.benchchem.com/product/b8745020#choosing-the-optimal-number-of-principal-components-in-dapcy]

---

**Disclaimer & Data Validity:**

**Technical Support:**The protocols provided are for reference purposes. Unsure if this reagent suits your experiment? [Contact our Ph.D. Support Team for a compatibility check]

**Need Industrial/Bulk Grade?**   Request Custom Synthesis Quote

# BenchChem

Our mission is to be the trusted global source of essential and advanced chemicals, empowering scientists and researchers to drive progress in science and industry.

Contact

Address: 3281 E Guasti Rd

Ontario, CA 91761, United States

Phone: (601) 213-4426

Email: info@benchchem.com