

Application Notes and Protocols for Cross-Validation Techniques in DAPCy

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

Compound of Interest

Compound Name: DAPCy

Cat. No.: B8745020

[Get Quote](#)

Audience: Researchers, scientists, and drug development professionals.

Introduction to DAPCy and the Role of Cross-Validation

Discriminant Analysis of Principal Components (DAPC) is a multivariate statistical method used to identify and describe clusters of genetically related individuals.[1][2][3] It is a powerful tool for inferring population structure from genetic markers like single nucleotide polymorphisms (SNPs).[4][5] The methodology first employs Principal Component Analysis (PCA) to reduce the dimensionality of the genetic data, transforming the correlated variables into a set of uncorrelated principal components (PCs).[1][2] Subsequently, it applies Discriminant Analysis (DA) to these PCs to maximize the separation between predefined or inferred groups.[2]

A critical step in DAPC is determining the optimal number of PCs to retain. Retaining too few PCs may result in the loss of important genetic information, while retaining too many can lead to overfitting, where the model captures random noise rather than the true population structure.[2][5] This can result in a model that performs well on the sampled data but poorly on new, unseen data.[2] Cross-validation is an essential technique to objectively select the optimal number of PCs, thereby ensuring the robustness and predictive accuracy of the DAPC model.[1][6][7]

Cross-Validation Methodologies in DAPC

Cross-validation assesses a model's ability to generalize to an independent dataset.[8][9] In the context of DAPC, this involves partitioning the data into a training set and a testing (or validation) set.[7][8] The DAPC model is built on the training set using a specific number of PCs, and its ability to correctly classify individuals in the testing set is evaluated.[1][7] This process is repeated for a range of different numbers of retained PCs, and the number that provides the best predictive performance is selected for the final analysis.[6][7]

Two primary software packages are used for DAPC, each with its own approach to cross-validation: the R package *adegenet* and the Python package **DAPCy**.

Repeated Random Sub-sampling Cross-Validation (in *adegenet*)

The *adegenet* package in R utilizes a repeated random sub-sampling or bootstrapping approach for cross-validation, implemented in the `xvalDapc` function.[6] This method involves repeatedly splitting the data, typically using 90% for the training set and 10% for the validation set.[6][7] To ensure that all groups are represented in both sets, stratified sampling is used.[7] The performance of the DAPC model for a given number of PCs is then averaged over many replicates to provide a robust estimate of the prediction success.[6] The optimal number of PCs is the one that minimizes the Mean Squared Error (MSE) or maximizes the proportion of successful predictions.[6]

k-Fold Cross-Validation (in **DAPCy**)

The **DAPCy** Python package, built on the scikit-learn library, offers several k-fold cross-validation schemes, which are more computationally efficient for large datasets.[4][10][11] In k-fold cross-validation, the dataset is divided into 'k' equal-sized folds.[12] The model is then trained on k-1 folds and tested on the remaining fold. This process is repeated k times, with each fold serving as the test set once.[12] The performance is then averaged across all k trials. **DAPCy** supports:

- Standard k-fold cross-validation: Randomly partitions the data into k folds.[10]
- Stratified k-fold cross-validation: Ensures that each fold has the same proportion of individuals from each group as the original dataset, which is crucial for imbalanced datasets. [10]

- Leave-one-out cross-validation (LOOCV): An extreme case of k-fold cross-validation where k is equal to the number of individuals. Each individual is used once as the test set.[\[10\]](#)[\[13\]](#)

DAPCy employs a grid-search approach to automatically test a range of PC numbers and identify the one with the highest classification accuracy.[\[10\]](#)[\[14\]](#)

Experimental Protocols

Protocol 1: Cross-Validation using xvalDapc in R (adegenet)

This protocol outlines the steps to determine the optimal number of PCs to retain in a DAPC analysis using the adegenet package in R.

Methodology:

- Load the necessary library and data:
- Perform the cross-validation: Run the xvalDapc function, specifying the genetic data, the group assignments, and the range of PCs to test.
- Interpret the results: The output of xvalDapc includes a plot showing the mean successful assignment per number of PCs retained. The number of PCs with the highest success rate (or lowest root mean squared error) is considered optimal.[\[7\]](#)
- Run the final DAPC with the optimal number of PCs:

Protocol 2: Grid Search Cross-Validation in Python (DAPCy)

This protocol describes how to find the optimal number of PCs using the grid search and cross-validation functionality in the **DAPCy** Python package.

Methodology:

- Install and import the necessary libraries:
- Initialize the **DAPCy** object:

- Retrieve and interpret the results: The best number of components and the corresponding accuracy are stored in the **DAPCy** object.
- Fit the final DAPC model: The `grid_search` function automatically fits the final model with the optimal number of PCs. You can access it for further analysis.

Data Presentation

The quantitative results from the cross-validation procedures can be summarized in tables for easy comparison.

Table 1: Example Output from `xvalDapc` in `ade4`

Number of PCs Retained	Mean Successful Assignment (%)	Standard Deviation	Mean Squared Error (MSE)
5	75.2	3.1	0.248
10	88.9	2.5	0.111
15	94.1	1.8	0.059
20	95.3	1.5	0.047
25	95.1	1.6	0.049
30	94.8	1.7	0.052

The optimal number of PCs is 20, as it corresponds to the lowest Mean Squared Error.

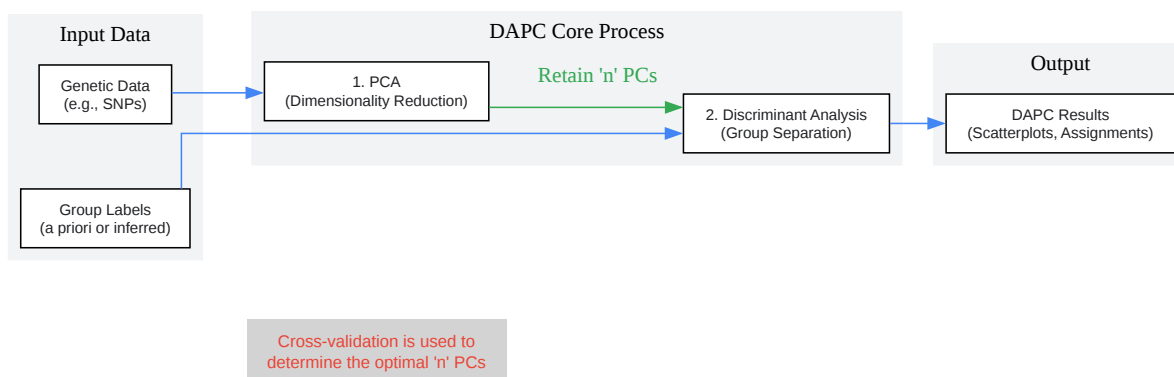
Table 2: Example Output from Grid Search in **DAPCy**

Number of PCs	Mean CV Accuracy	Standard Deviation of CV Accuracy
5	0.761	0.032
10	0.893	0.028
15	0.945	0.021
20	0.958	0.019
25	0.956	0.020
30	0.952	0.022

The optimal number of PCs is 20, achieving the highest mean cross-validation accuracy.

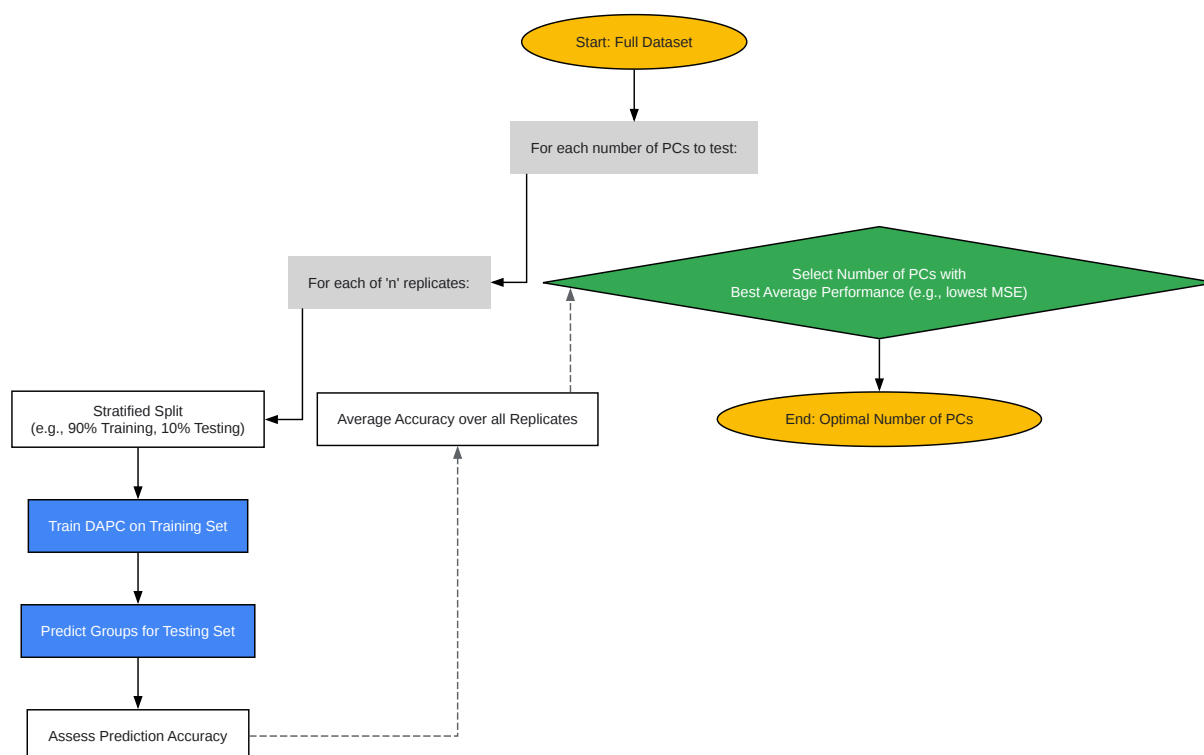
Visualizations

The following diagrams illustrate the workflows and logical relationships of the cross-validation techniques in DAPC.



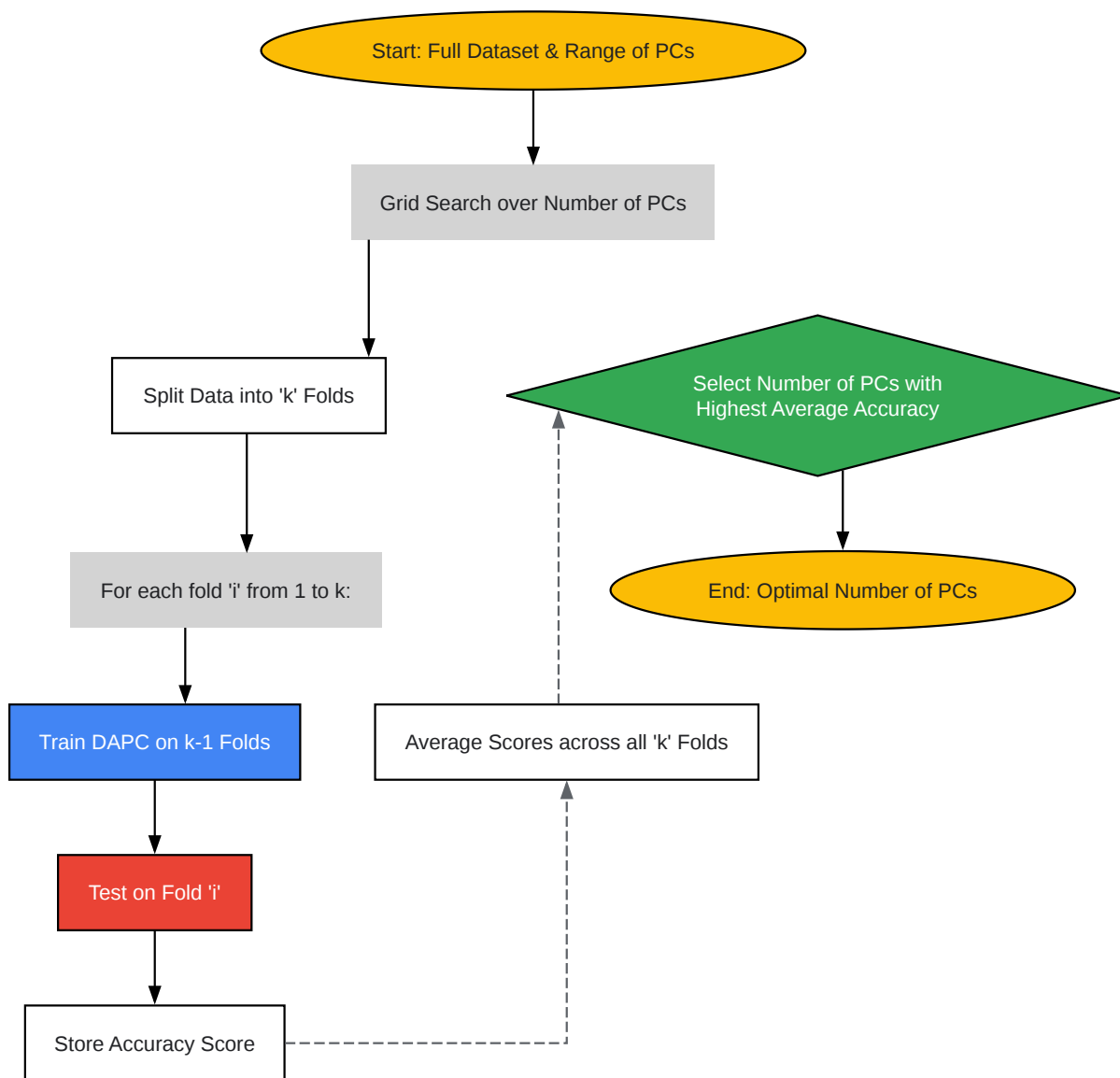
[Click to download full resolution via product page](#)

Caption: Overall workflow of the Discriminant Analysis of Principal Components (DAPC).



[Click to download full resolution via product page](#)

Caption: Workflow for xvalDapc in the R adegenet package.



[Click to download full resolution via product page](#)

Caption: Workflow for k-fold cross-validation grid search in the **DAPCy** Python package.

Need Custom Synthesis?

BenchChem offers custom synthesis for rare earth carbides and specific isotopic labeling.

Email: info@benchchem.com or [Request Quote Online](#).

References

- 1. Discriminant analysis of principal components (DAPC) [grunwaldlab.github.io]
- 2. Discriminant analysis of principal components: a new method for the analysis of genetically structured populations - PMC [pmc.ncbi.nlm.nih.gov]
- 3. The influence of a priori grouping on inference of genetic clusters: simulation study and literature review of the DAPC method - PMC [pmc.ncbi.nlm.nih.gov]
- 4. DAPCy [uhasselt-bioinfo.gitlab.io]
- 5. Genetic diversity, linkage disequilibrium, population structure and construction of a core collection of *Prunus avium* L. landraces and bred cultivars - PMC [pmc.ncbi.nlm.nih.gov]
- 6. HTTP redirect [search.r-project.org]
- 7. Discriminant analysis of principal components and pedigree assessment of genetic diversity and population structure in a tetraploid potato panel using SNPs - PMC [pmc.ncbi.nlm.nih.gov]
- 8. Cross-validation (statistics) - Wikipedia [en.wikipedia.org]
- 9. CrossValidation Techniques for Classification Models|Keylabs [keylabs.ai]
- 10. academic.oup.com [academic.oup.com]
- 11. DAPCy: a Python package for the discriminant analysis of principal components method for population genetic analyses - PubMed [pubmed.ncbi.nlm.nih.gov]
- 12. neptune.ai [neptune.ai]
- 13. dambe.bio.uottawa.ca [dambe.bio.uottawa.ca]
- 14. DAPCy Tutorial: MalariaGEN *Plasmodium falciparum* - DAPCy [uhasselt-bioinfo.gitlab.io]
- To cite this document: BenchChem. [Application Notes and Protocols for Cross-Validation Techniques in DAPCy]. BenchChem, [2025]. [Online PDF]. Available at: [https://www.benchchem.com/product/b8745020#cross-validation-techniques-in-dapcy]

Disclaimer & Data Validity:

The information provided in this document is for Research Use Only (RUO) and is strictly not intended for diagnostic or therapeutic procedures. While BenchChem strives to provide accurate protocols, we make no warranties, express or implied, regarding the fitness of this product for every specific experimental setup.

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