

Installing and Utilizing DAPCy for Population Genetic Analysis

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

Compound of Interest

Compound Name: DAPCy

Cat. No.: B8745020

[Get Quote](#)

Application Notes and Protocols for Researchers, Scientists, and Drug Development Professionals

Introduction

DAPCy is a high-performance Python package for conducting Discriminant Analysis of Principal Components (DAPC), a multivariate statistical method used to identify and describe genetic population structure.^[1] Originally implemented in the R package adegenet, **DAPCy** offers a reimplementaion that leverages the scikit-learn library for enhanced scalability and efficiency, particularly with large genomic datasets.^{[2][3][4]} This document provides detailed instructions for installing **DAPCy** in a Python environment, along with protocols for its application in population genetic analyses, specifically using the Plasmodium falciparum Pf7 dataset as an illustrative example.

Installation Protocols

DAPCy can be installed via pip or conda/mamba, and it is highly recommended to perform the installation within a virtual environment to avoid conflicts with other packages.^[3]

Prerequisites:

- A Python version of 3.6 or higher is required.^[5]

- For Windows users intending to import VCF files, it is recommended to install **DAPCy** within a Windows Subsystem for Linux (WSL) environment due to a dependency on cyvcf2.[\[1\]](#)[\[3\]](#) Zarr files can be used as an input on Windows without this requirement.[\[1\]](#)[\[3\]](#)
- conda users should use Python version 3.12 or lower to avoid potential dependency conflicts.[\[1\]](#)[\[3\]](#)

Installation using pip:

- Create and activate a virtual environment:
- Install **DAPCy**:

Installation using conda/mamba:

- Create and activate a conda environment:
- Install **DAPCy** from the bioconda channel:

Experimental Protocols

DAPCy can be employed in two primary scenarios for analyzing population structure: with a priori knowledge of population groups or for de novo inference of genetic clusters using k-means clustering.[\[1\]](#) The following protocols are based on the official **DAPCy** tutorial using the Plasmodium falciparum Pf7 dataset.[\[1\]](#)

Protocol 1: Population Classification with a priori Population Labels

This protocol is applicable when the population groups of the samples are already known (e.g., country of origin).

Methodology:

- Data Loading and Preparation:
 - Load the genotype data (e.g., from VCF or BED files) and the corresponding sample metadata which includes the population labels.

- **DAPCy** includes a function (`geno2csr.py`) to extract genotype values and convert them into a compressed sparse row (csr) matrix, which is efficient for large datasets.[\[2\]](#)[\[4\]](#)
- Data Splitting:
 - Divide the dataset into training and testing sets to evaluate the performance of the DAPC classifier. This is a standard machine learning practice to assess how well the model generalizes to new data.
- Principal Component Analysis (PCA):
 - Perform PCA on the training data to reduce the dimensionality of the genetic data. **DAPCy** uses a truncated Singular Value Decomposition (SVD) for this step, which is computationally efficient for sparse matrices.[\[1\]](#)[\[4\]](#)
 - Determine the optimal number of principal components (PCs) to retain. A common approach is to use the k-1 criterion, where k is the number of known populations, to capture the essential variance for biological interpretation while maintaining computational efficiency.[\[1\]](#)
- Discriminant Analysis of Principal Components (DAPC):
 - Train the DAPC model using the retained principal components from the training set and their corresponding population labels.
 - The `dapc` class in **DAPCy** provides functions for model training and cross-validation.[\[2\]](#)[\[4\]](#)
- Model Evaluation:
 - Use the trained DAPC model to predict the population labels of the samples in the testing set.
 - Generate a classification report to assess the model's performance, including metrics like accuracy, precision, and recall for each population group.[\[1\]](#)[\[4\]](#)
- Visualization:

- Visualize the results by plotting the individuals on the discriminant axes. This allows for a visual inspection of the separation between the predefined population groups.

Protocol 2: De Novo Inference of Genetic Clusters using K-means Clustering

This protocol is used when there is no prior knowledge of the population structure.

Methodology:

- Data Loading:
 - Load the genotype data.
- Principal Component Analysis (PCA):
 - Perform PCA on the entire dataset to reduce dimensionality.
- K-means Clustering:
 - Apply the k-means clustering algorithm to the principal components to identify genetic clusters.
 - To determine the optimal number of clusters (k), the sum of squared errors (SSE) is calculated for a range of k values. The "elbow" point in the plot of SSE against k indicates the optimal number of clusters.[\[1\]](#)
 - The `kmeans_group()` function in **DAPCy** can be used for this purpose.[\[1\]](#)
- Discriminant Analysis of Principal Components (DAPC):
 - Once the optimal number of clusters is determined and individuals are assigned to these inferred clusters, proceed with the DAPC analysis as described in Protocol 1, using the inferred clusters as the population labels.
- Model Evaluation and Visualization:

- Evaluate the DAPC model and visualize the results to understand the genetic structure of the inferred populations.

Data Presentation

DAPCy provides functions to generate classification reports that summarize the performance of the DAPC model. These reports typically include key metrics for evaluating the accuracy of the classification.

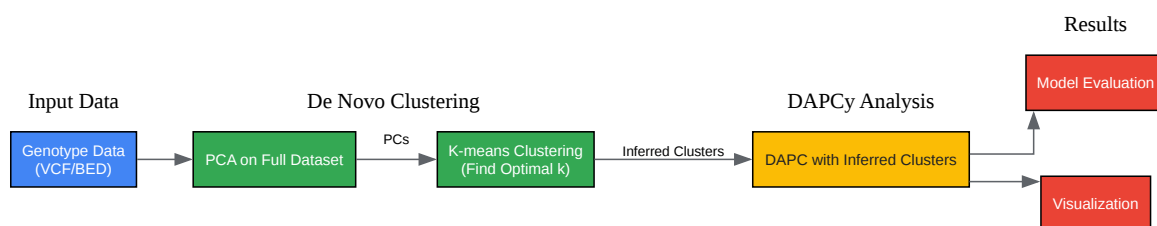
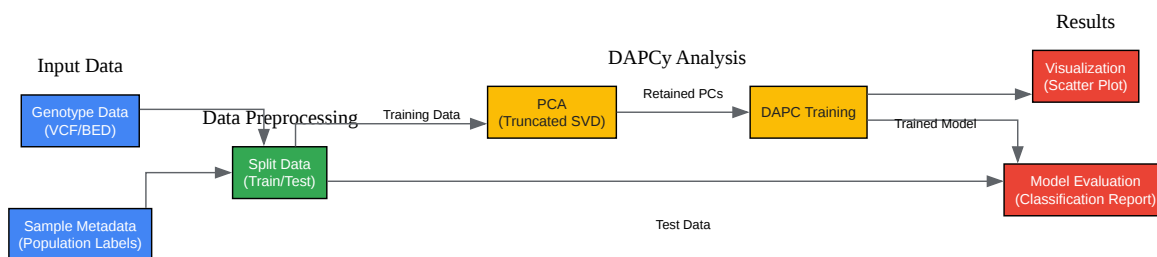
Table 1: Representative Classification Report for a DAPC Analysis

Metric	Population 1	Population 2	Population 3	...	Overall Mean
Precision	0.95	0.92	0.98	...	0.95
Recall	0.96	0.91	0.97	...	0.95
F1-Score	0.95	0.91	0.97	...	0.95
Support	100	120	95	...	315
Accuracy	-	-	-	-	0.95

- Precision: The ratio of correctly predicted positive observations to the total predicted positive observations.
- Recall: The ratio of correctly predicted positive observations to all observations in the actual class.
- F1-Score: The weighted average of Precision and Recall.
- Support: The number of actual occurrences of the class in the specified dataset.
- Accuracy: The ratio of correctly predicted observations to the total observations.

Visualizations

The following diagrams illustrate the key workflows in a **DAPCy** analysis.



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

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. academic.oup.com [academic.oup.com]
- 2. wellcomeopenresearch-files.f1000.com [wellcomeopenresearch-files.f1000.com]
- 3. An open dataset of Plasmodium falciparum genome variation in 7,000 worldwide samples - PMC [pmc.ncbi.nlm.nih.gov]

- 4. adegenet.r-forge.r-project.org [adegenet.r-forge.r-project.org]
- 5. researchgate.net [researchgate.net]
- To cite this document: BenchChem. [Installing and Utilizing DAPCy for Population Genetic Analysis]. BenchChem, [2025]. [Online PDF]. Available at: [<https://www.benchchem.com/product/b8745020#how-to-install-dapcy-in-a-python-environment>]

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