

slapnap: Super LeArner Prediction of NAb Panels

David Benkeser, Brian D. Williamson, Craig A. Magaret, Bhavesh R. Borate, Peter B. Gilbert

May 24, 2020

Contents

Welcome	1
1 Docker	2
2 CATNAP Database	3
3 Running the slapnap container	3
3.1 Mounting a local directory	3
3.2 slapnap options	4
3.3 Interactive sessions	5
4 Examples	5
4.1 Basic calls to slapnap	5
4.2 Super learning	6
4.3 Train an algorithm	6
4.4 Pull and clean data	6
5 Report details	6
6 Data details	6
7 Method details	6
7.1 Outcome definitions	6
7.2 Super learner details	6
7.3 Variable importance details	6
8 References	6

Welcome

The **slapnap** container is a tool for using the Compile, Analyze and Tally NAb Panels (CATNAP) database to develop predictive models of HIV-1 neutralization sensitivity to one or several broadly neutralizing antibodies (bNAbs).



Crystal structure of HIV-1 gp120 glycoprotein. Highlighted residues indicating sites most-predictive of VRC01 neutralization resistance. [magaret2019prediction]

In its simplest form, **slapnap** can be used simply to access and format data from CATNAP in a way that is usable for machine learning analysis. However, the tool also offers fully automated and customizable machine learning analyses based on up to five different neutralization endpoints, complete with automated report generation to summarize results and identify the most predictive features.

This document serves as the user manual for the **slapnap** container. Here, we describe everything needed to utilize the **slapnap** container and understand its output. The documentation is organized into the following sections:

- Section 1 provides a brief overview of Docker, including information on installing Docker and downloading the **slapnap** container.
- Section 2 provides a brief overview of the CATNAP database and the specifics of how the data were accessed to build the **slapnap** container.
- Section 3 provides a detailed description of how to make calls to the **slapnap** repository, including descriptions of all options that are available.
- Section 4 includes several example calls to the **slapnap** container and descriptions of their output.
- Section 6 provides a description of the data set created in the **slapnap** container.
- Section 7 provides an overview of the methodology that is used in within the **slapnap** analysis.

If you have any issues or questions about using **slapnap**, please file an issue on GitHub.

1 Docker

Docker is a free platform for building containers. Containers are standard units of software that package code and all its dependencies, so that the code can be executed reliably irrespective of computing environment. The **slapnap** tool relies on machine learning implemented in the R language and relies on several packages. Achieving full reproducibility for such analyses is challenging in that it requires synchronization across the specific version of R and dependent packages. In other words, two users running two versions of R or two versions of the same R package may arrive at different output when running the same code. Containerization ensures that this does not happen. Any two runs of **slapnap** with the same input options will yield the same output every time.

Installing Docker is necessary for running the **slapnap** tool. While it is not necessary for execution of the **slapnap** container, readers interested in learning more about Docker should consult the Docker documentation for information about getting started using Docker.

Once Docker has been installed on your local computer, you can download **slapnap** using the following command.

```
docker pull slapnap/slapnap
```

This command pulls the image from DockerHub. Once the image has been downloaded, we are ready to learn about how to execute **slapnap** jobs. The next section contains information on the source data used by **slapnap**. Users familiar with the CATNAP data may wish to skip directly to Section 3.2.

2 CATNAP Database

The CATNAP database is a web server hosted by Los Alamos National Laboratory [Yoon et al., 2015]. The database integrates antibody neutralization and HIV-1 sequence data from published studies. Neutralization is measured in terms of half maximal inhibitory concentration (IC_{50}) and 80% inhibitory concentration (IC_{80}). These measures of neutralization against HIV envelope pseudoviruses are available for many broadly neutralizing antibodies (bNAbs) and for some combination bNAbs. Also available on each pseudovirus are amino acid (AA) sequence features for the gp160 protein. These are detailed in Section 6.

During each build of the **slapnap** container, all raw data are downloaded from CATNAP. At run time, the relevant data are selected and processed into a format that is amenable for predictive machine learning analyses. The CATNAP data are updated periodically. To check the date the raw data were pulled from CATNAP to **slapnap**, you can check the date of the **latest** build [here](#).

3 Running the slapnap container

To run the **slapnap** container, we make use of the **docker run** command. Note that administrator (**sudo**) privileges are needed to execute this command.

There are several options that are necessary to include in this command to control the behavior of **slapnap**. These are discussed in separate subsections below.

3.1 Mounting a local directory

At the end of a **slapnap** run, user-specified output will be saved (see option **return** in Section 3.2). To retrieve these files from inside the container, we *mount* a local directory to an output directory (**/home/out/**) in the container using the **-v** option. That is, all files in the mounted local directory will be visible to programs running inside the container and any items saved to the output directory in the container (file path in the container **/home/out/**) will be available in the mounted directory.

Suppose **/path/to/local/dir** is the file path on a local computer in which we wish to save the output files from a **slapnap** run. A **docker run** of **slapnap** would include the option **-v /path/to/local/dir:/home/out**. After a run completes, the requested output should be viewable in **/path/to/local/dir**. See Section 4 for full syntax.

3.2 slapnap options

The user has control over many aspects of **slapnap**'s behavior. These options are passed in using the **-e** option¹. Semi-colon separated strings are used to set options. For example, to provide input for the option **option_name**, we would use **-e option_name="a;semi-colon;separated;string"**. Note that there are no spaces between the option name and its value and no spaces after semi-colons in the separated list. See Section 4 for full syntax.

Each description below lists the default value that is assumed if the option is not specified. Note that many of the default values are chosen simply so that `na{ }i` calls to **slapnap** compile quickly. Proper values should be determined based on scientific context.

-e options for slapnap

- **nab**: A semicolon-separated list of bNAbs (default = "VRC01"). A list of possible bNAbs can be found here. If multiple bNAbs are listed, it is assumed that the analysis should be of estimated **outcomes** for a combination of bNAbs (see Section 7.1 for details on how estimated outcomes for multiple bNAbs are computed).
- **outcomes**: A semicolon-separated string of outcomes to include in the analysis. Possible values are "ic50" (included in default), "ic80", "iip", "sens" (included in default), "estsens", "multsens". If only a single **nab** is specified, use **sens** to include a dichotomous endpoint. If multiple **nabs** are specified, use **estsens** and/or **multsens**. For detailed definitions of outcomes see Section 7.1).
- **sens_thresh**: A numeric value defining the IC₅₀ threshold for defining a sensitive versus resistant pseudovirus (default = 1). The dichotomous sensitivity/resistant **outcomes** are defined as the indicator that (estimated) IC₅₀ is greater than or equal to **sens_thresh**.
- **multsens_nab**: A numeric value used for defining whether a pseudovirus is resistant to a multi-nAb cocktail. The dichotomous outcome **multsens** is defined as the indicator that a virus has (estimated) IC₅₀ greater than **sens_thresh** for at least **multsens_nab** nAbs.
- **learners**: A semicolon-separated string of machine learning algorithms to include in the analysis. Possible values include "rf" (random forest, default), "xgboost" (eXtreme gradient boosting), and "lasso" (elastic net). If more than one algorithm is included, then it is assumed that a cross-validated-based ensemble (i.e., a super learner) is desired (see Section 7.2).
- **cvtune**: A boolean string (i.e., either "TRUE" or "FALSE" [default]) indicating whether the **learners** should be tuned using cross-validation and a small grid search. Defaults to "FALSE". If multiple **learners** are specified, then the super learner ensemble includes three versions of each of the requested **learners** with different tuning parameters.
- **cvperf**: A boolean string (i.e., either "TRUE" or "FALSE" [default]) indicating whether the **learners** performance should be evaluated using cross-validation. If **cvtune**="TRUE" or **learners** includes multiple algorithms, then nested cross-validation is used to evaluate the performance of the cross-validation-selected best value of tuning parameters for the specified algorithm or the super learner, respectively.
- **nfolds**: A numeric string indicating the number of folds to use in cross-validation procedures (default = "2").
- **importance_grp**: A semicolon-separated string indicating which group-level variable importance measures should be computed. Possible values are none "" (default), marginal "marg", conditional "cond". See Section 7.3.1 for details on these measures.
- **importance_ind**: A semicolon-separated string indicating which individual-level variable importance measures should be computed. Possible values are none "" (default), learner-level "pred", marginal "marg" and conditional "cond". The latter two take significant computation time to compute.
- **report_name**: A string indicating the desired name of the output report (default = `report_[_-separated list of nabs]_[date].html`).
- **return**: A semicolon-separated string of the desired output. Possible values are "report" (default), "learner" for the trained algorithm, "data" for the analysis dataset, "figures" for all figures from the report, and "vimp" for variable importance objects.

¹This sets an environment variable in the container environment. These variables are accessed by the various R and bash scripts in the container to dictate how the container executes code.

3.3 Interactive sessions

To simply enter the container and poke around, use an interactive session by including `-it` and overriding the container's entry-point.

```
docker run -it slapnap/slapnap /bin/bash
```

This will enter you into the container in a bash terminal. This may be useful for exploring the file structure, examining versions of R packages that are included in the container, etc.

4 Examples

4.1 Basic calls to `slapnap`

A call to `slapnap` with all default options can be run using the following command.

```
docker run -v /path/to/local/dir:/home/output slapnap/slapnap
```

Note that this call mounts the local directory `path/to/local/dir` to receive output from the container (Section 3.1).

When this command is executed,

4.2 Super learning

4.3 Train an algorithm

4.4 Pull and clean data

5 Report details

6 Data details

7 Method details

7.1 Outcome definitions

7.2 Super learner details

7.3 Variable importance details

7.3.1 Biological importance

7.3.2 Predictive importance

8 References

References

Hyejin Yoon, Jennifer Macke, Anthony P West Jr, Brian Foley, Pamela J Bjorkman, Bette Korber, and Karina Yusim. CATNAP: a tool to compile, analyze and tally neutralizing antibody panels. *Nucleic Acids Research*, 43(W1):W213–W219, 2015. doi: 10.1093/nar/gkv404.