



SAnDReS User Guide

STATISTICAL ANALYSIS OF DOCKING RESULTS AND SCORING
FUNCTIONS

Gabriela Bitencourt-Ferreira and Walter F. de Azevedo, Jr. | Version 2.0 |
February 25, 2022

Contents

| | |
|--|----|
| 1. Conventions and Availability..... | 02 |
| 2. Introduction..... | 03 |
| 3. Installing SAnDReS (Linux)..... | 04 |
| 4. Installing SAnDReS (Windows)..... | 06 |
| 5. Tutorial 1. Cycling-Dependent Kinase 2 with IC ₅₀ Data..... | 07 |
| 5.1. Tutorial 1. Setup..... | 08 |
| 5.2. Tutorial 1. Dataset..... | 11 |
| 5.3. Tutorial 1. Docking Hub..... | 17 |
| 5.4. Tutorial 1. Scoring Functions..... | 22 |
| 5.5. Tutorial 1. Virtual Screening..... | 26 |
| 5.6. Tutorial 1. Machine Learning Box (For Modeling)..... | 30 |
| 5.7. Tutorial 1. Machine Learning Box (For Docking Results)..... | 37 |
| 5.8. Tutorial 1. Machine Learning Box (For Virtual Screening Results)..... | 39 |
| 5.9. Tutorial 1. Statistical Analysis..... | 41 |
| 6. Tutorial 2. Molegro Virtual Docker Data..... | 46 |
| 6.1. Tutorial 2. Setup..... | 47 |
| 6.2. Tutorial 2. Machine Learning Box..... | 49 |
| 6.3. Tutorial 2. Statistical Analysis..... | 52 |
| 7. Tutorial 3. Docking of One Structure..... | 54 |
| 7.1. Tutorial 3. Setup..... | 55 |
| 7.2. Tutorial 3. Dataset..... | 56 |
| 7.3. Tutorial 3. Docking Hub..... | 58 |
| 8. Tutorial 4. Cycling-Dependent Kinase 2 with K _i Data..... | 60 |
| 8.1. Tutorial 4. Setup..... | 61 |
| 8.2. Tutorial 4. Dataset..... | 63 |
| 8.3. Tutorial 4. Docking Hub..... | 65 |
| 8.4. Tutorial 4. Scoring Functions..... | 66 |
| 8.5. Tutorial 4. Machine Learning Box..... | 67 |
| 9. Scoring Function Space..... | 69 |
| 10. FAQ. Frequently Asked Questions..... | 70 |
| 11. GitHub..... | 76 |
| 12. Authors..... | 77 |
| 13. Acknowledgments..... | 78 |
| 14. Colophon..... | 79 |
| 15. References..... | 80 |
| Appendix. Regression Methods..... | 83 |

1. Conventions and Availability

This User Guide shows how to install and use the various attributes of the SAnDReS program (Version 2.0.0). This guide includes the capabilities of the program, how to apply these capabilities, and how to install SAnDReS on Linux and Windows OS.

Here, we have the following typographical conventions:

Arial font with italic

Indicates filenames, folders (directories) in the main text and commands in the SAnDReS menu.

Courier New font with Italic

Used for Linux/Windows commands, PDB listings, and data to be typed by the user.

SAnDReS is open-source software and freely distributed under GNU General Public License v3.0 (GPL-3.0 License). Its code is available to download on GitHub (<https://github.com/azevedolab/sandres>).

In the following sections, we describe the installation guidelines and four tutorials showing how to use SAnDReS for machine learning and docking simulations.

2. Introduction

SAnDReS (Statistical Analysis of Docking Results and Scoring functions) draws inspiration from several protein systems we have been working on in the last decades. These projects began in the 1990s with pioneering studies focused on intermolecular interactions between cyclin-dependent kinase (CDK) (EC 2.7.11.22) and inhibitors (de Azevedo et al., 1996; de Azevedo et al., 1997). SAnDReS is a free and open-source (GPL-3.0 License) computational environment for the development of machine-learning models (Bitencourt-Ferreira & de Azevedo, 2019; Bitencourt-Ferreira et al., 2021; Bitencourt-Ferreira, Rizzotto et al., 2021) for the prediction of ligand-binding affinity (Xavier et al., 2016; Bitencourt-Ferreira & de Azevedo, 2019; Veit-Acosta & de Azevedo, 2021). We developed SAnDReS using Python 3 programming language and SciPy, NumPy, Scikit-Learn (Pedregosa et al., 2011), XGBoost, and Matplotlib libraries as a computational tool to explore the Scoring Function Space (Heck et al., 2017; Bitencourt-Ferreira & de Azevedo, 2019).

This version is a new software with the same goal. SAnDReS explores the Scoring Function Space using the machine learning methods available in Scikit-Learn (version 1.0.2). SAnDReS 1.0 generated machine learning models based on nine regression methods only. It had a fixed polynomial scoring function set that used three features to generate up to 512 polynomial equations taking squared and mixed independent variables (x^2 and $x.y$) (Xavier et al., 2016). SAnDReS 2.0 has 64 regression methods, with no limitation in the number of features used in the machine learning models. SAnDReS 1.0 used AutoDock Vina 1.1.2 (Trott & Olson, 2010), now we have the most recent version, AutoDock Vina 1.2.3 (Eberhardt et al., 2021). There are a lot of minor modifications, but these previously highlighted are the most important ones.

SAnDReS 2.0.0 brings together the most advanced tools for protein-ligand docking simulation and machine-learning modeling. We have the newest version of AutoDock Vina (Trott & Olson, 2010; Eberhardt et al., 2021), available in February 2022 (version 1.2.3), as a docking engine. Also, SAnDReS 2.0.0 uses the latest version of Scikit-Learn, available in February 2022 (version 1.0.2). SAnDReS predicts binding affinity for a specific protein system with superior performance compared against classical scoring functions. In summary, SAnDReS makes it possible for you to design a scoring function adequate to the protein system of your interest.

You need Python 3 installed on your computer to run SAnDReS 2.0.0. In addition, you need Matplotlib, NumPy, Scikit-Learn, SciPy, and XGBoost. It is also necessary to have MGLTools 1.5.7 (Morris et al., 2009). You can make the installation of Python packages faster by installing Anaconda.

3. Installing SAnDReS (Linux)

You should type all commands shown here in a Linux terminal. The easiest way to open a Linux terminal is to use the Ctrl+Alt+T key combination.

Step 1. Download MGLTools 1.5.7 (<https://ccsb.scripps.edu/mgltools/downloads/>).

Type the following commands:

```
cd ~
cp Downloads/mgltools_Linux-x86_64_1.5.7_install .
chmod u+x mgltools_Linux-x86_64_1.5.7_install
./mgltools_Linux-x86_64_1.5.7_install
```

Step 2. Download Anaconda Installer for Linux

(https://repo.anaconda.com/archive/Anaconda3-2021.11-Linux-x86_64.sh). Go to the directory where you have the installer file and type the following commands:

```
chmod u+x Anaconda3-2021.11-Linux-x86_64.sh
./Anaconda3-2021.11-Linux-x86_64.sh
```

Follow the instructions of the installer.

Step 3. To run SAnDReS properly, you need Scikit-Learn 1.0.2. To be sure you have version 1.0.2, open a terminal, and type the following commands:

```
python3 -m pip uninstall scikit-learn
python3 -m pip install scikit-learn==1.0.2
```

Step 4. To install XGBoost

(<https://xgboost.readthedocs.io/en/latest/install.html#python>), type the following command in a terminal:

```
python3 -m pip install xgboost
```

Step 5. Download SAnDReS

(<https://github.com/azevedolab/sandres/raw/master/sandres2.zip>). Copy the sandres2 zipped directory (sandres2.zip) to wherever you want it and unzip the zipped directory.

Type the following command:

```
unzip sandres2.zip
```

change to *sandres2* directory then, type:

```
python3 sandres2.py
```

Now you have the GUI window for SAnDReS

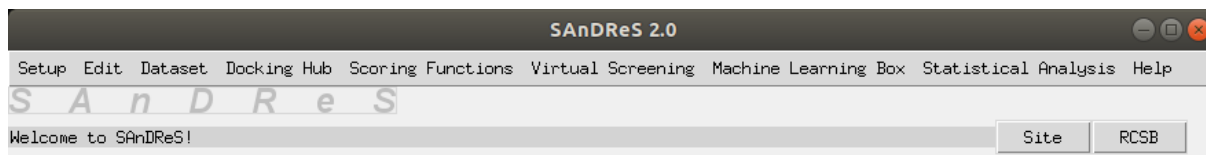


Figure 1. SAnDReS Main Menu (Linux Version).

That's it, good SAnDReS session!

4. Installing SAnDReS (Windows)

Step 1. Install MGLTools 1.5.7 (<https://ccsb.scripps.edu/mgltools/downloads/>).

Step 2. Install Anaconda (<https://www.anaconda.com/download/>). Click on the Windows Start Menu icon and select the Anaconda prompt. From now on, insert all commands in the Anaconda prompt.

Step 3. To run SAnDReS properly, you need Scikit-Learn 1.0.2. To be sure you have version 1.0.2, open an Anaconda prompt, and type the following commands:

```
python -m pip uninstall scikit-learn
```

```
python -m pip install scikit-learn==1.0.2
```

Step 4. To install XGBoost

(<https://xgboost.readthedocs.io/en/latest/install.html#python>), type the following command:

```
python -m pip install xgboost
```

Step 5. Download SAnDReS

(https://github.com/azevedolab/sandres/blob/master/sandres2_win.zip). Copy the *sandres2_win* zipped directory (*sandres2_win.zip*) to wherever you want it and unzip the zipped directory. Open the Anaconda Prompt, and cd to *sandres2_win* directory then, type:

```
python sandres2.py
```

Now you have the GUI window for SAnDReS. That's it, good SAnDReS session!

5. Tutorial 1. Cyclin-Dependent Kinase 2 with IC₅₀ Data

In this tutorial, we will generate a machine learning model to predict binding affinity for cyclin-dependent kinase 2 (CDK2). The idea is to use the experimental information (crystallographic and IC₅₀) and build a machine-learning model.

We will carry out docking simulations applying SAnDReS to the CDK2 structures. We will perform statistical analysis of these docked structures evaluating docking accuracy (DA) and the root-mean-squared deviations (RMSD). Also, once SAnDReS created a machine-learning model based on the crystallographic structures, we will evaluate the ranking of poses using this new scoring function targeted to CDK2.

We consider that you have successfully installed SAnDReS as described in the previous sections ([Section 3](#) and [Section 4](#)). We captured all images used to illustrate this tutorial running the Windows version. The overall sequence is the same for the Linux version.

In this tutorial, I used the SAnDReS 2.0.0 installed on *C:\Users\Walter\sandres2_win*

Open an Anaconda3 Prompt and cd to *sandres2_win* directory. For Linux, cd to the *sandres2* directory. Type the following commands:

```
cd c:\Users\Walter\sandres2_win
```

```
python sandres2.py
```

It is necessary to use the directory you have SAnDReS on your computer.

Upon starting the program SAnDReS, you have the following GUI.

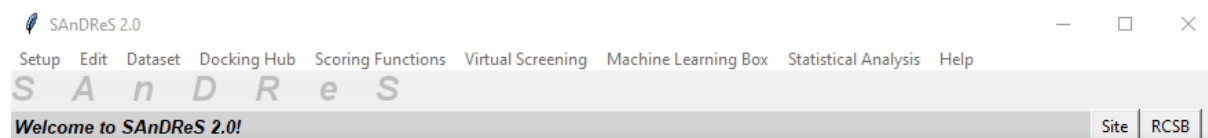


Figure 2. SAnDReS Graphical User Interface (Windows Version).

In the rest of this tutorial, we will show the running of all necessary tasks through the above GUI.

Note. After SAnDReS installation, we have the following folders in the *sandres2_win* directory (*sandres2* for Linux): *datasets*, *Ligands*, *misc*, *MLRegMPy*, and *tools*. The last two folders have Python codes used by SAnDReS. We find auxiliary files necessary to run SAnDReS in the *misc* folder. We keep the project directories in the *dataset* folder, one project directory for each protein system (e.g., *CDK2_IC50*). In Tutorials [1](#) and [2](#), we need files found in the *Ligands* folder.

5.1. Tutorial 1. Setup

Here, we will download the ligand databases and define the project directory. SAnDReS has predefined ligand databases with experimental binding affinities (K_i , K_d , and IC_{50}) and generated PDBQT files for ligand structures. We used the affinity data available in the PDBbind (Wang et al., 2004) to generate these databases. The project directory is where SAnDReS keeps all files generated during its execution. We expect one for each protein system.

Click on *Setup->Check Ligand Datasets*, as shown below.

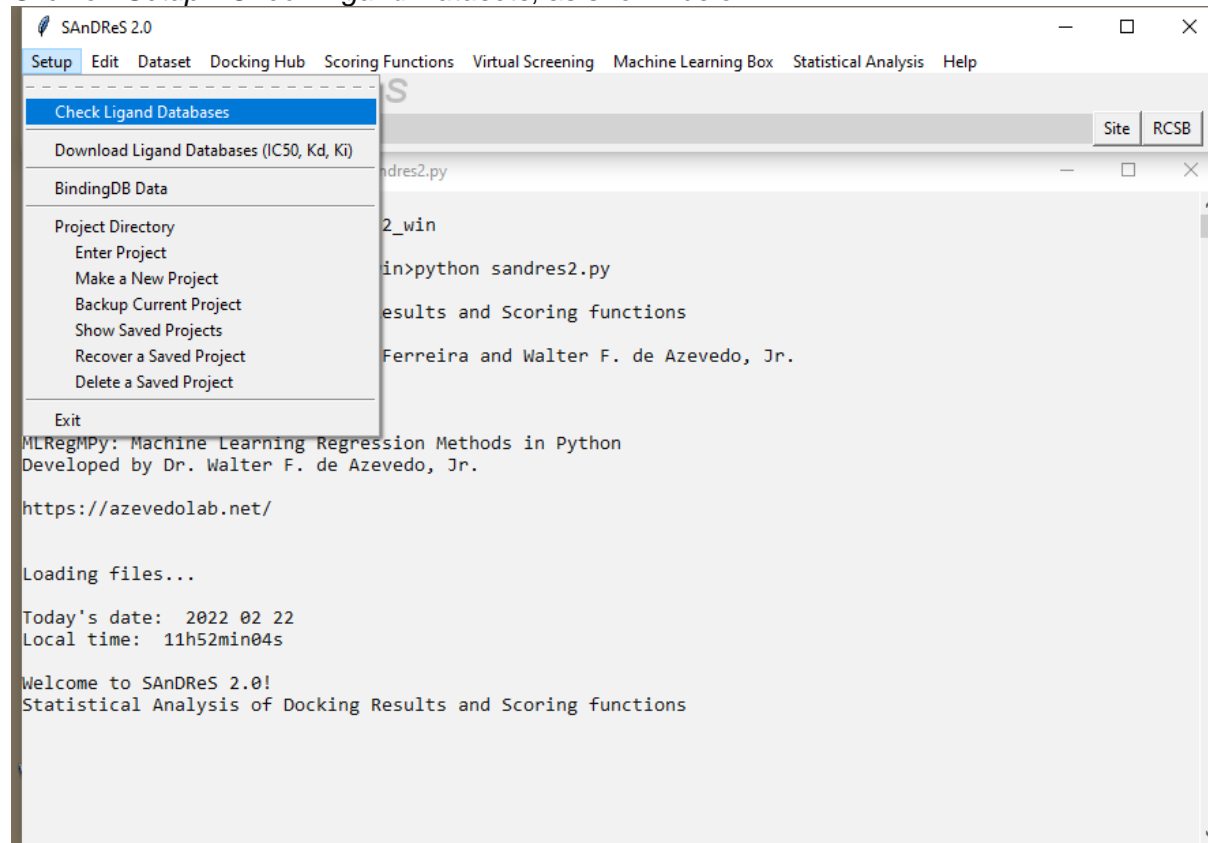


Figure 3. Dataset Menu (Windows Version).

SAnDReS will show the number of ligands in the databases. As we can see, we have zero ligands in the databases. SAnDReS communicates with you through the message bar, right below the SAnDReS logo.

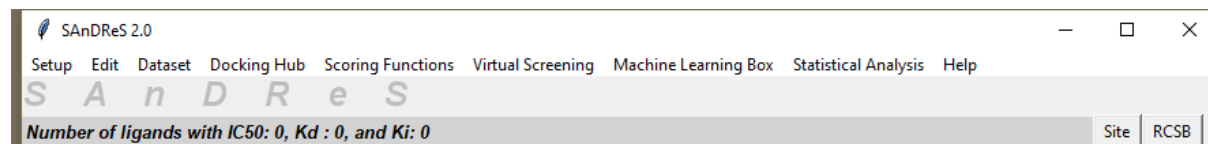


Figure 4. Main Menu with information on the message bar (Windows Version).

To download the ligand databases, click on *Setup->Download Ligand Databases (IC50, Kd, Ki)*. Click on the Yes option. After finishing downloading, you get the following message:

SAnDReS updated Ligand Databases!

Next, we will enter the project directory.

Click on *Setup->Project Directory->Enter Project*, as shown below.

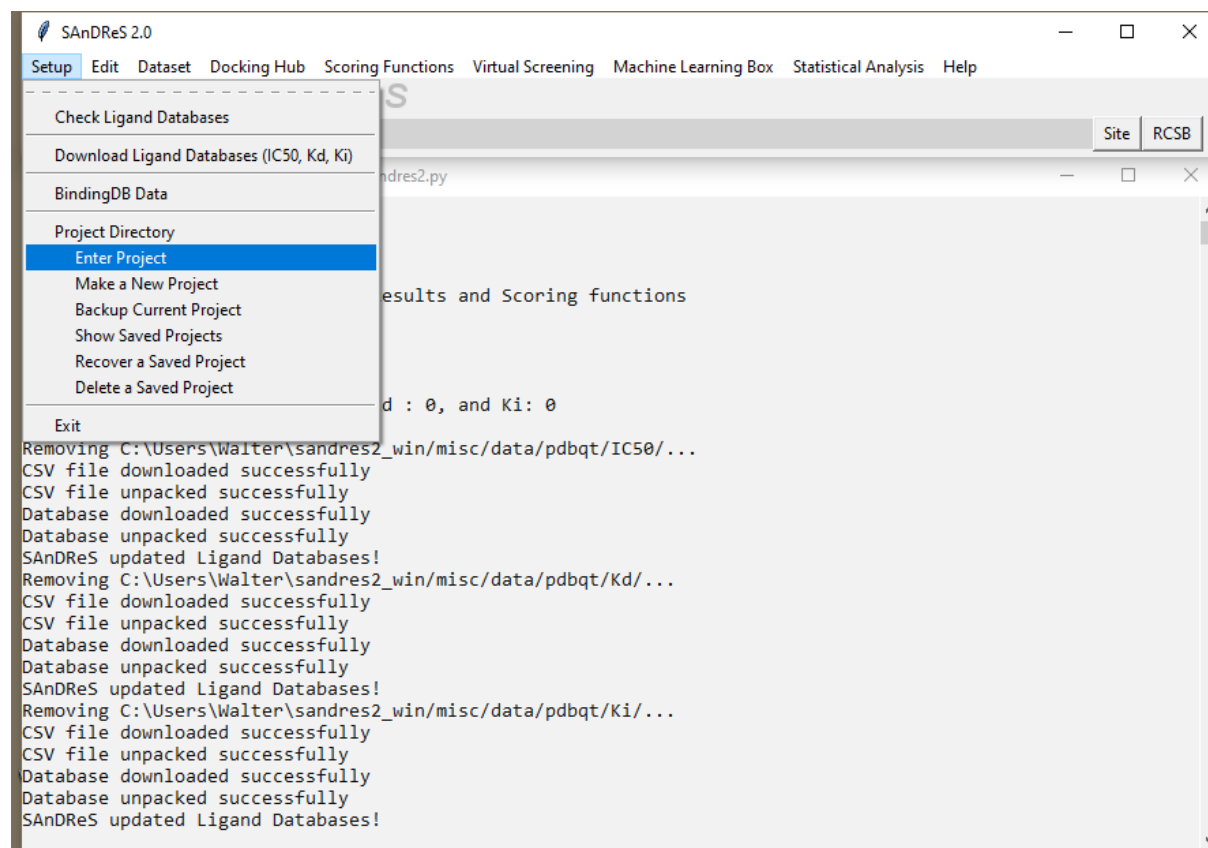


Figure 5. Dataset Menu (Windows Version).

SAnDReS will open an editor (Fast Editor), and you should insert the directory where we will have all data related to this modeling, as shown in Figure 6. With this approach, we have a project directory for each protein system.

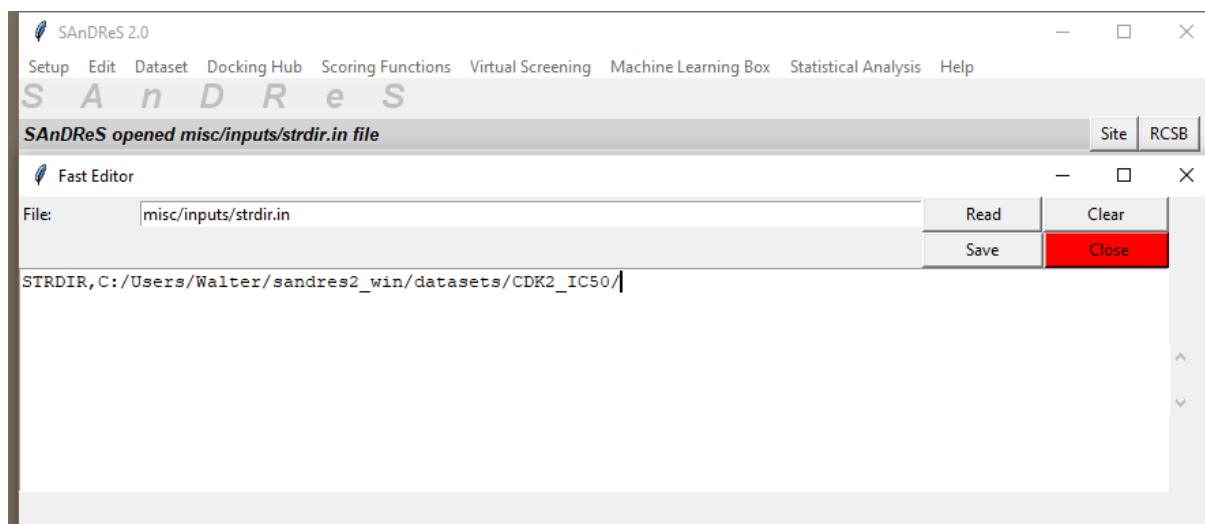


Figure 6. Fast Editor (Windows Version).

After writing the project directory in the Fast Editor, click on the *Save* button and *Close* button. I recommend creating your project directory in the *datasets* folder. For instance, you want to create a project directory named *CDK2_IC50*. You type the following string in the Fast Editor:

```
C:/Users/Walter/sandres2_win/datasets/CDK2_IC50/
```

Do not forget the final slash */*.

Then, you click on *Setup->Project Directory->Make a New Project*. Click on the *Yes* option. SAnDReS will generate a new directory named *CDK2_IC50*. You should get the following message:

Successfully created the directory

```
C:/Users/Walter/sandres2_win/datasets/CDK2_IC50/
```

Now, we finished the Setup. For this tutorial, you do not need to interact with the other options of the Setup Menu.

Note. During a SAnDReS session, we keep all graphic files in the *plots* folder of the project directory. SAnDReS saves all downloaded PDB files in the *pdb* folder of the project directory. We find generated PDBQT files in the *pdbqt* folder.

5.2. Tutorial 1. Dataset

In this part, we will download crystal structures from the Protein Data Bank (PDB) and select the binding affinity data. We will generate PDBQT files for all entries in the dataset.

Click on *Edit->Parameters->Project Summary*. This part is not mandatory is only to add some description about this dataset. With this task, SAnDReS edits a file named *summary.txt*. You may add some explanation, for instance: CDK2 with IC50 data. After adding this sentence, click on *Save* button and *Close* button.

Now, click on *Dataset->Enter PDB Access Codes*. You will have the following screen.

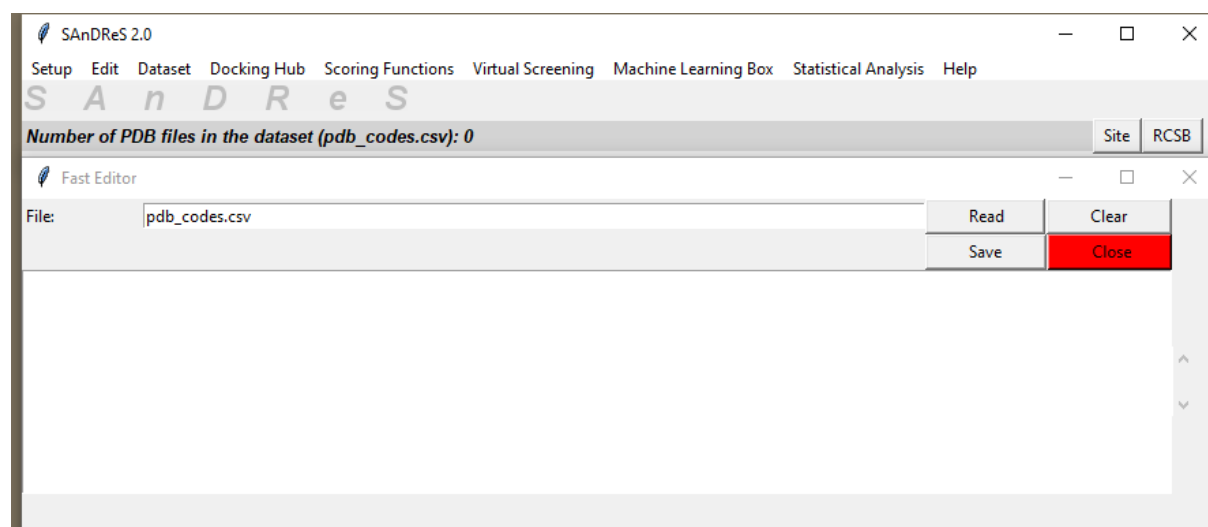


Figure 7. Fast Editor (Windows Version).

Now, you Ctrl C the following PDB access codes and Ctrl V to the SAnDReS editor. You may download the *pdb_codes.csv* file by clicking [here](#). Save it in the project directory.

```
3IG7,3WBL,2VTH,3IGG,2VTA,2VTP,2VTO,2VTN,2VTM,3R8V,2VTL,3R8U,2V
TI,4LYN,3UNJ,3QTZ,3QTX,3QTW,4EZ3,3TIY,3QTU,1PXL,3QTS,3QTR,3QTQ
,3QU0,2W17,3TI1,1Y91,1G5S,1P XK,2W1H,5D1J,3EZV,2W06,2W05,3EZR,2
A4L,1V1K,3LFN,3QOK,2R64,3RAH,2B52,1W0X,2B53,2B55,1R78,2C6I,2C6
K,3RAL,2C6L,2C6M,3RPY,3RPV,2BTR,1P2A,2BTS,3FZ1,2UZO,2UZN,2R3M,
2R3N,2R3O,3S2P,2R3G,2C5Y,2R3H,2R3I,2A0C,3S1H,1DI8,4ERW,3SQQ,3S
00,3PJ8,2DUV,3RNI,4RJ3,1KE5,1KE6,1KE7,1KE8,3RMF,1VYZ,3PY1,3PXZ
,4GCJ,2VV9,3NS9,3RK9,3RK7,3RK5,3RKB,2VTT,2VTS,2VTQ,3R8Z,1OIQ,1
OIR,3R9D,3RJC,3R9N,2DS1
```

For Linux, you may have to edit the PDB access codes to have the following screen.

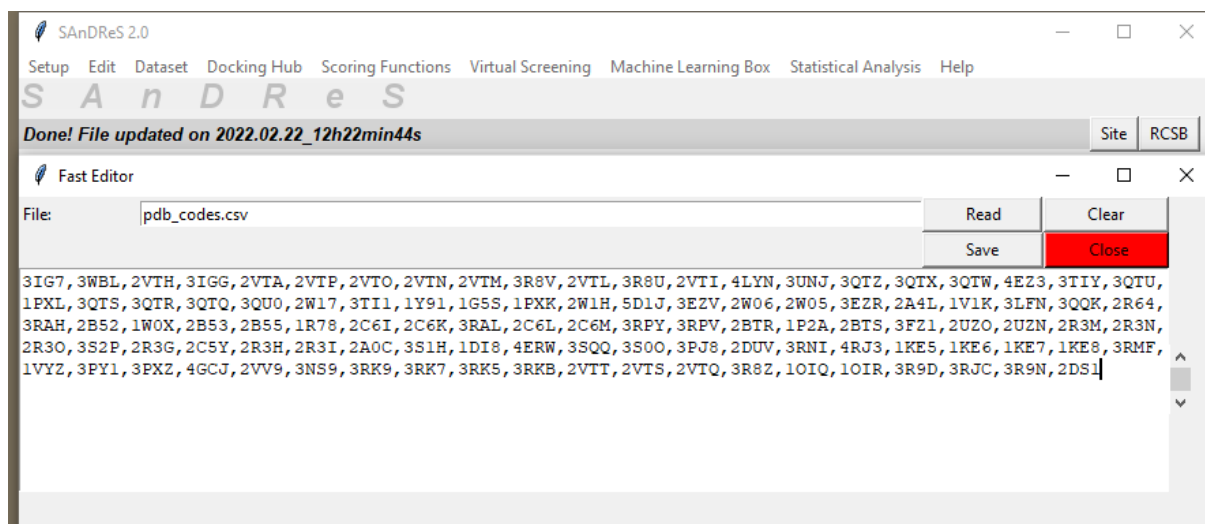


Figure 8. Fast Editor (Windows Version).

Click on the Save button and the Close button. SAnDReS saved the PDB access codes in a file named *pdb_codes.csv*. If you reopen the *pdb_codes.csv* file, SAnDReS should show that you have 104 structures. To avoid repeated PDB codes, click on *Dataset->Unify PDB Access*. You get the following message:

Done! Number of repeated PDB files: 0.

Click on *Dataset->Add->Binding Affinity Data*. In the pop-up window, click on the *Half-maximal Inhibitory Concentration (IC50)* button. Then, click on the Start button. Once finished, you will get the following message:

SAnDReS finished "Add Binding Affinity Data" request!

SAnDReS added experimental binding affinity and generated a file named *bind_IC50.csv*. Click on the Done button and the Close button.

In the following, we have the first lines of the *bind_IC50.csv* file.

| PDB | Ligand | Chain | Number | Resolution(A) | Ligand Occupation Fe | IC50(M) | log(IC50) | pIC50 |
|------|--------|-------|--------|---------------|----------------------|----------|-----------|---------|
| 3IG7 | EFP | A | 999 | 1.8 | 1 | 6.3e-08 | -7.20066 | 7.20066 |
| 3WBL | PDY | A | 302 | 2 | 1 | 2.3e-05 | -4.63827 | 4.63827 |
| 2VTH | LZ2 | A | 1300 | 1.9 | 1 | 0.00012 | -3.92082 | 3.92082 |
| 3IGG | EFQ | A | 999 | 1.8 | 1 | 6.65e-08 | -7.17718 | 7.17718 |
| 2VTA | LZ1 | A | 1301 | 2 | 1 | 0.000185 | -3.73283 | 3.73283 |
| 2VTP | LZ9 | A | 1299 | 2.1 | 1 | 3e-09 | -8.52288 | 8.52288 |
| 2VTO | LZ8 | A | 1299 | 2.1 | 1 | 1.4e-07 | -6.85387 | 6.85387 |
| 2VTN | LZ7 | A | 1299 | 2.2 | 1 | 8.5e-07 | -6.07058 | 6.07058 |
| 2VTM | LZM | A | 1299 | 2.2 | 1 | 0.001 | -3 | 3 |
| 3R8V | Z62 | A | 473 | 1.9 | 1 | 2.9e-06 | -5.5376 | 5.5376 |
| 2VTL | LZ5 | A | 1299 | 2 | 1 | 9.7e-05 | -4.01323 | 4.01323 |
| 3R8U | Z31 | A | 465 | 2 | 1 | 5e-06 | -5.30103 | 5.30103 |
| 2VTI | LZ3 | A | 1299 | 2 | 1 | 6.6e-07 | -6.18046 | 6.18046 |
| 4LYN | 1YG | A | 301 | 2 | 1 | 6e-08 | -7.22185 | 7.22185 |
| 3UNJ | OBX | A | 299 | 1.9 | 1 | 1.1e-05 | -4.95861 | 4.95861 |

Figure 9. Partial view of the *bind_IC50.csv* file.

Click on *Dataset->Add->Structures PDB*. In the pop-up window, click on the *with IC50 data* button. Then, click on the *Start* button. SAnDReS will start the downloading of the PDB files. It may take a while since we have more than 100 structures in this dataset. Figure 10 shows the progress bar during the downloading.

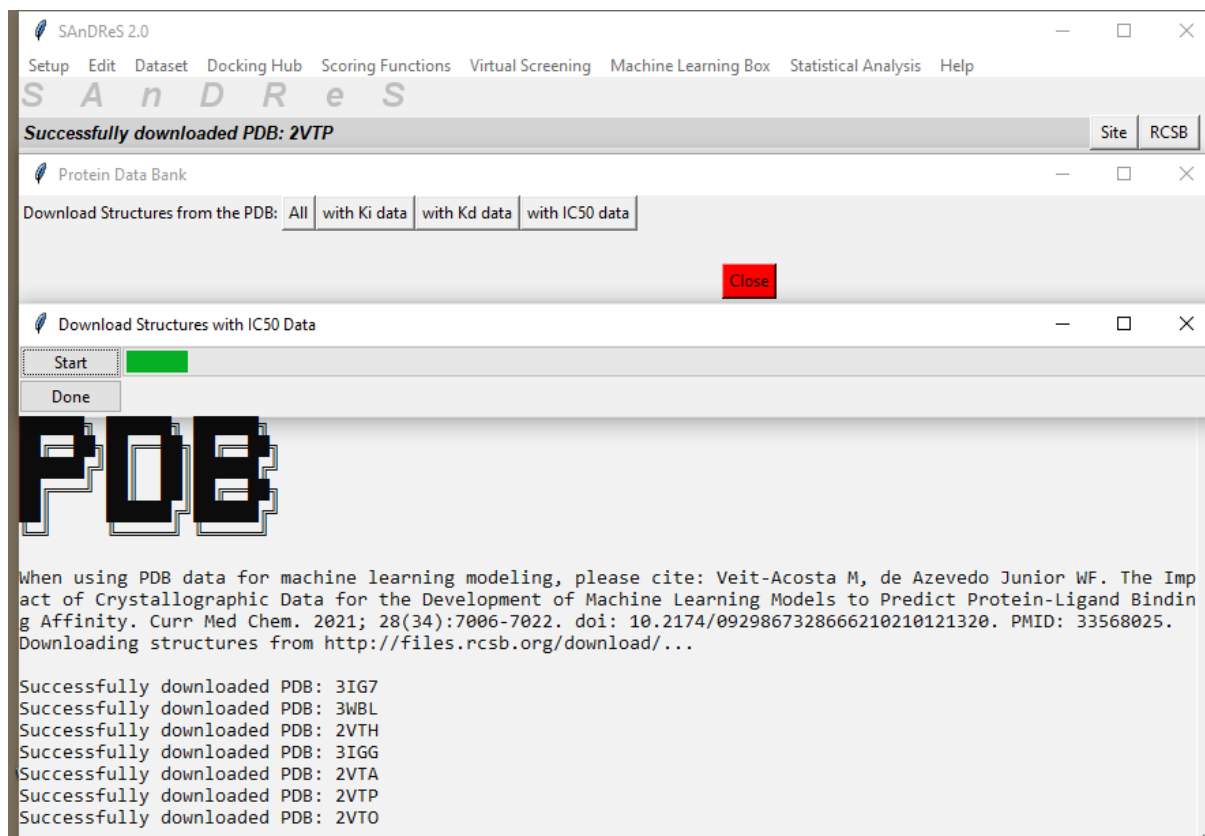


Figure 10. Progress bar of the downloading process (Windows Version).

After finishing the downloading, you get the following message:

SANdReS finished the “Add Structures (PDB)” request!

Click on the *Done* button and then on the *Close* button.

SANdReS created a *pdb* folder in the project directory with all downloaded structures and generated a plot named *resol_binding.pdf*. This plot is in the *plots* folder.

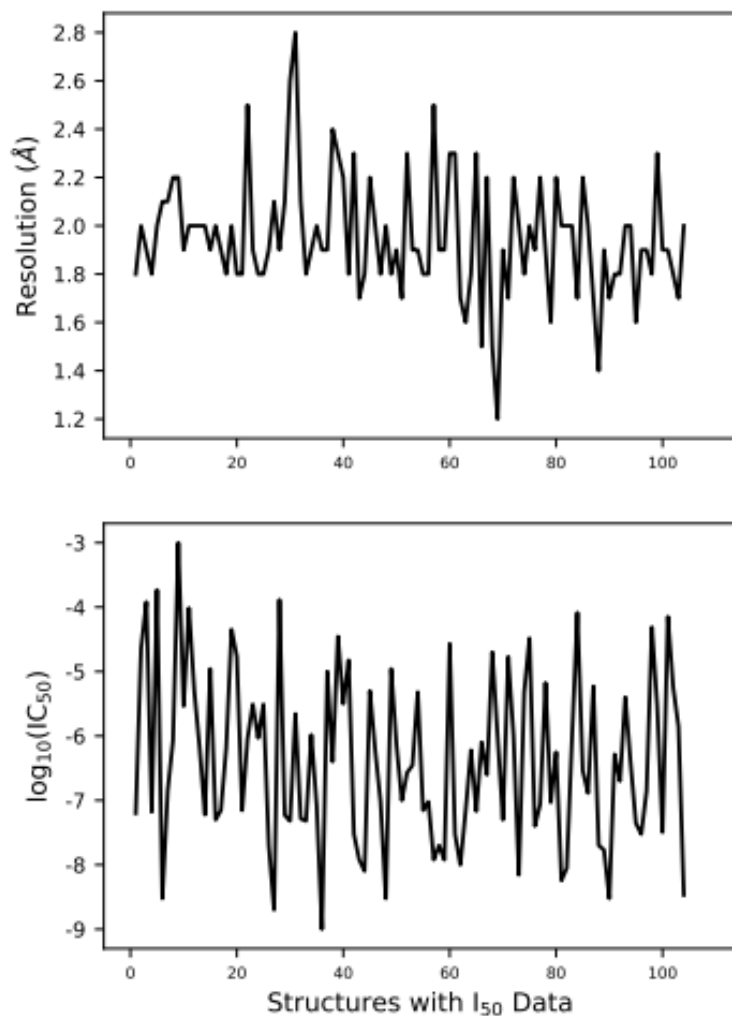


Figure 11. Crystallographic resolution and $\log(\text{IC}_{50})$ for the structures of the CDK2_IC50 dataset.

Click on *Dataset->Add->Structures PDBQT*. In the warning window, click on the *Yes* button.

SANdReS will show the following pop-up window.

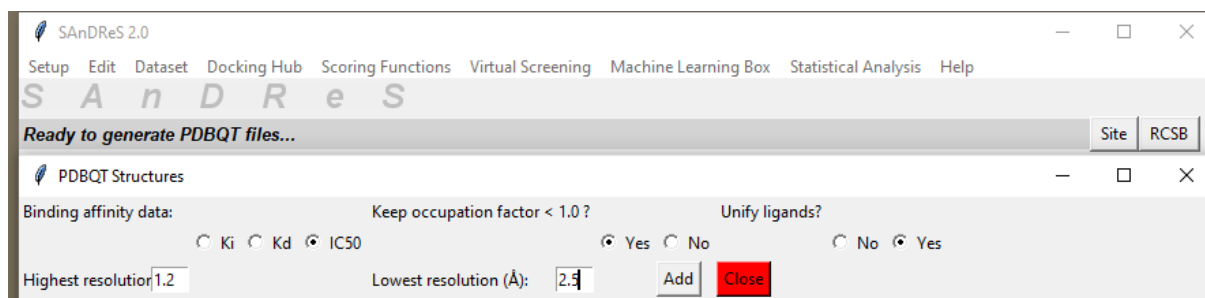


Figure 12. PDBQT Structures Menu (Windows Version).

Set the lowest resolution to 2.5 Å and the binding affinity to IC50. Choose Yes for the *Unify ligands* option. *Unify ligands* option will eliminate repeated ligands. You may have different PDB codes for CDK2 with the same ligand. In this case, SAnDReS chooses the best crystallographic resolution structure.

Click on the *Add* button, then on the *Start* button.

You can follow the progress as shown below.

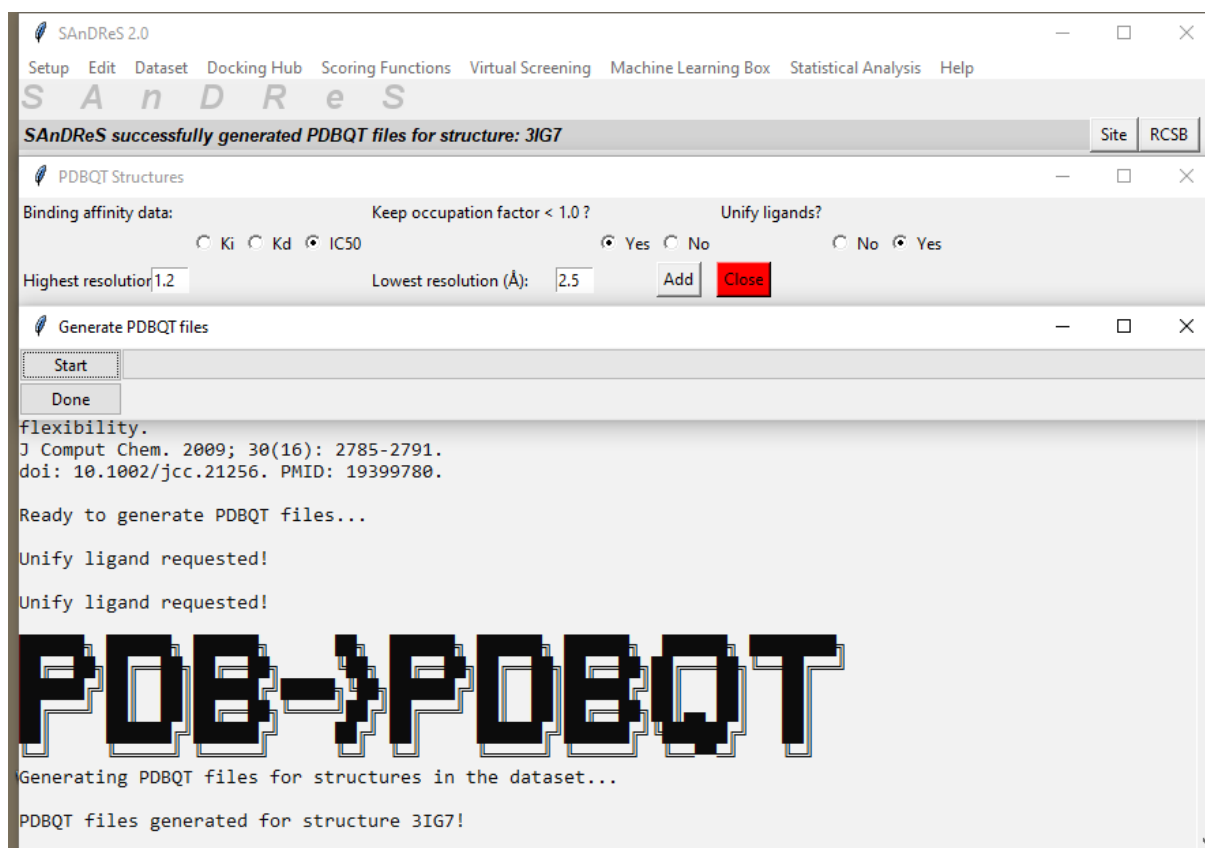


Figure 13. PDBQT Structures Menu with progress bar (Windows Version).

After this task, you will get the following message:

SAnDReS finished the “Add Structures (PDBQT)” request!

Click on the *Done* button and the *Close* button.

Note. SAnDReS could face some problems in generating PDBQT files. In this case, you may delete the problematic PDB access codes from the dataset or use AutoDockTools4 (Morris et al., 2009) to generate the missing files.

If you decide to do so, you should generate PDBQT files in the same directory created by SAnDReS. For instance, in the structure *XXXX*, we should keep all files in the *C:\Users\Walter\sandres2_win\datasets\CDK2_IC50\pdbqt\XXXX* directory. Since we have over 100 structures, we believe that is enough for all purposes (machine learning modeling) and we will delete any PDB code with missing PDBQT files. We have more than five observations (structures) for each feature used in the model (Gramatica, 2013).

Click on *Dataset->Check Directories for Current Dataset*. Click on the *Yes* option. SAnDReS will check any missing directory in the *pdbqt* folder.

You get the following message:

Done! SAnDReS updated dataset!

Click on *Dataset->Check Missing PDBQT Files*. SAnDReS will check any missing PDBQT files.

You get the following message:

Done! No missing PDBQT files in the dataset!

Since we do not have missing PDBQT files, we do not have to update this dataset. Now, our dataset has now 102 structures. SAnDReS deleted the following structures: 1G5S and 1PXK.

We finished the Dataset part of this tutorial.

5.3. Tutorial 1. Docking Hub

Here, we will carry out redocking for all crystallographic structures in the dataset. The goal here is to validate a docking protocol using AutoDock Vina 1.2.3. We will apply the best docking protocol to perform virtual screening (VS). Also, we will employ a SAnDReS-generated scoring function to sort VS results and predict binding for poses generated during the docking simulations.

Click on *Docking Hub*->*Enter Vina Parameters*. SAnDReS will open the file *vina_par.csv*, as shown in the figure below.

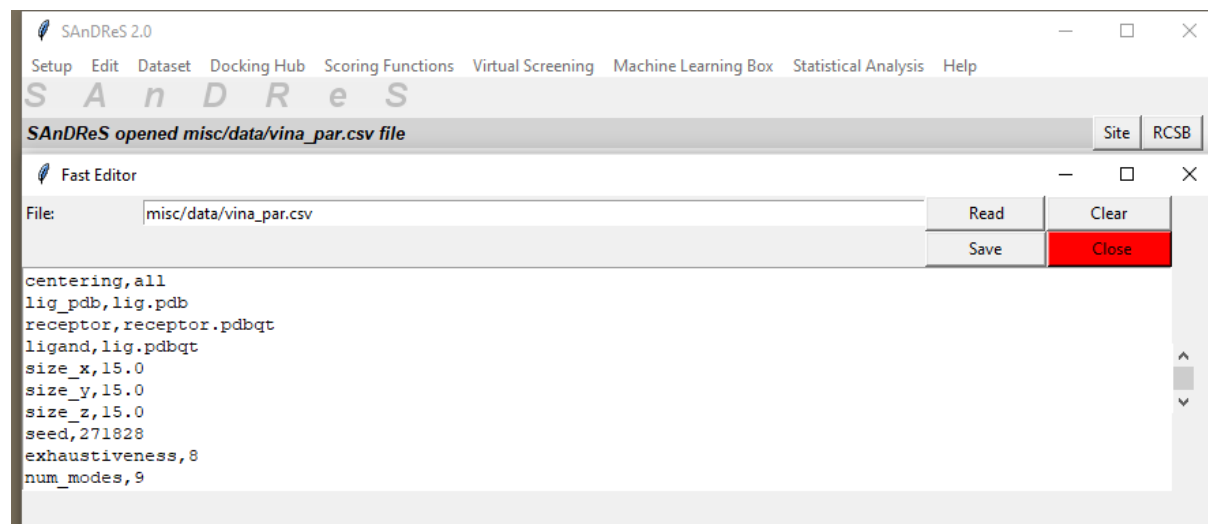


Figure 14. Fast Editor (Windows Version).

If you are familiar with AutoDock Vina 1.2.3, you may play with the parameters. For this tutorial, do not change the following parameters: *centering*, *lig_pdb*, *receptor*, *ligand*, *out*, *docking_engine*, and *score_only_vina*.

Click on the *Save* and *Close* buttons.

SAnDReS 2.0 has two main approaches to running docking simulations with AutoDock Vina. The first, named “One-Click Docking with Vina (Centering: CM, EC, GC)” performs docking simulation using three different centering schemes. The first centering is the center of mass (CM) which is the standard procedure used by AutoDock Vina, the center of mass of ligand atoms. The electric center (EC) carries out an average of the charges found in the PDBQT files and takes it as the center for the ligand. SAnDReS takes the absolute values ($|q|$). It intends to locate the center close to the concentration of electric charges in the ligand structure. We define the coordinates of electric centers (x_{EC} , y_{EC} , z_{EC}) as follows:

$$x_{EC} = \frac{1}{N} \sum_{i=1}^N (|q_i| x_i)$$

$$y_{EC} = \frac{1}{N} \sum_{i=1}^N (|q_i| y_i)$$

$$z_{EC} = \frac{1}{N} \sum_{i=1}^N (|q_i| z_i)$$

Where x_i , y_i , z_i are the atomic coordinates of each non-hydrogen atom in the ligand structure, N is the number of non-hydrogen atoms in the ligand, and q_i is the partial charge of each non-hydrogen atom. The geometric center (GC) is the center where we consider that all non-hydrogen atoms have masses equal to 1.0.

SAnDReS performs three docking simulations for each structure in the dataset and considers the best docking protocol for each entry, the one with the lowest RMSD.

The option named “One-Click Docking with Vina (Centering: CM)” performs docking simulation using CM only.

In this tutorial, we will carry out docking simulations with “One-Click Docking with Vina (Centering: CM, EC, GC)” to explore this new approach implemented in SAnDReS 2.0.

Click on *Docking Hub*->*One-Click Docking with Vina (Centering: CM, EC, GC)*->*Run Simulation*. Click on the *Yes* option, then click on the *Run* button. You will have the following menu.

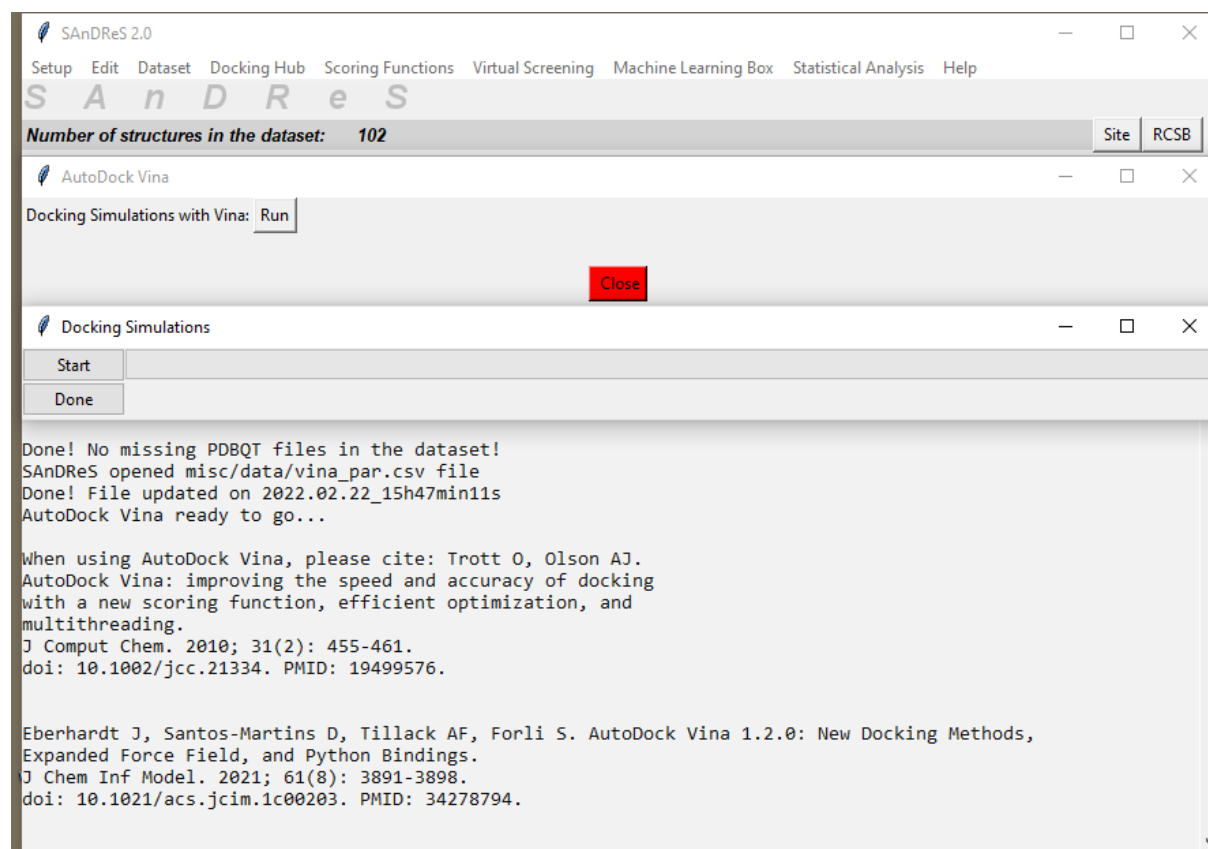


Figure 15. Docking Simulations with Vina (Windows Version).

Click on the *Start* button to initiate the docking simulations.

This process may take a few hours, depending on your CPU. You may follow the evolution of the process by observing the progress bar, as shown in Figure 16.

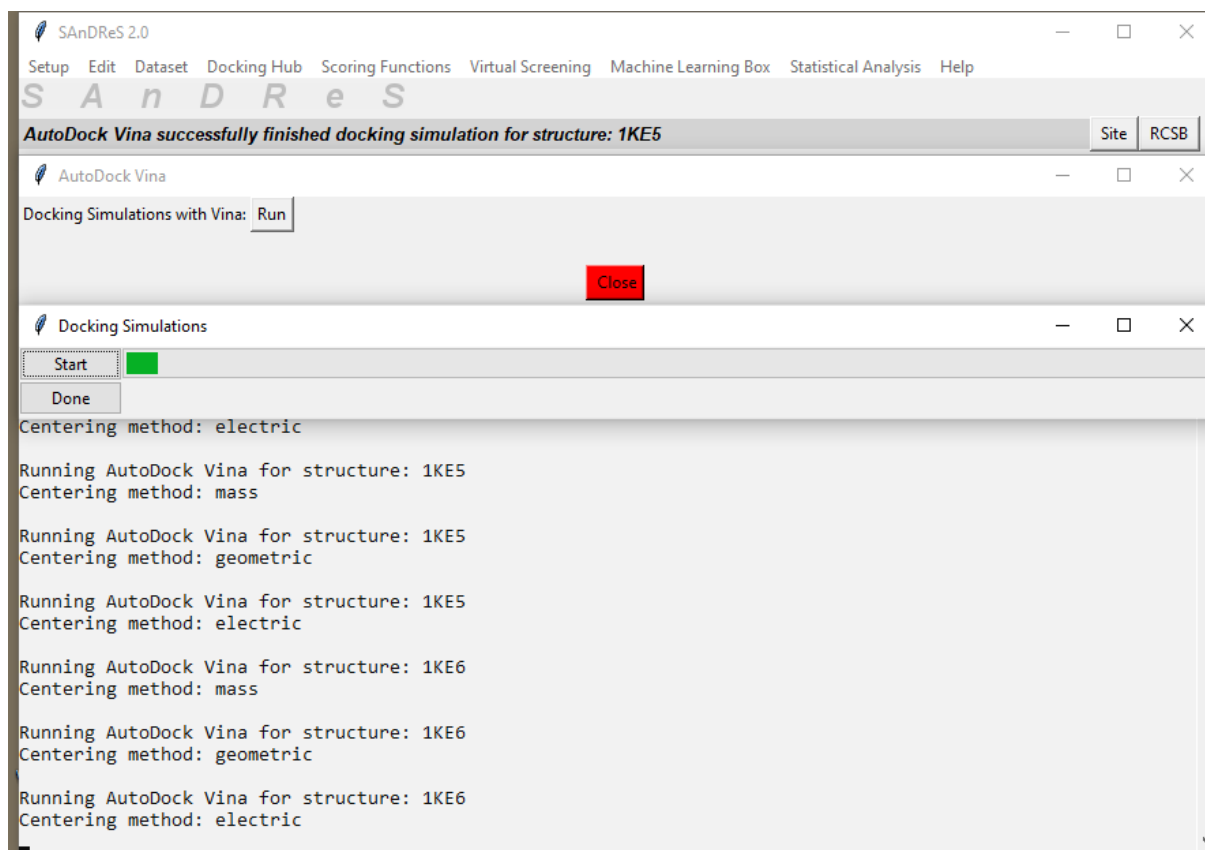


Figure 16. Docking Simulations with Vina (Windows Version).

For any unsuccessful docking simulation for a structure XXXX, you get a warning message as follows:

Warning! AutoDock Vina couldn't run docking simulation for structure: XXXX!

If the unsuccessfully docked structure is meaningful, you may try to run AutoDock Vina 1.2.3 outside SAnDReS. If you do that, keep the same file names and directories. SAnDReS will use them for further statistical analysis of the docking results.

Once finished the docking simulations, you have the following message:

SAnDReS finished the “One-Click Docking with Vina” request using AutoDock Vina!

Click on the *Done* button and the *Close* button.

All files generated during docking simulations are on the related folders in the *pdbqt* directory of the project directory. For instance, for the structure 1DI8, we have the following files shown in Figure 17.



















| | |
|---|------------------|
|  1DI8 | 22/02/2022 14:32 |
|  biomatrix | 22/02/2022 14:32 |
|  config_electric | 22/02/2022 15:57 |
|  config_geometric | 22/02/2022 15:56 |
|  config_mass | 22/02/2022 15:55 |
|  lig | 22/02/2022 14:32 |
|  lig.pdbqt | 22/02/2022 14:32 |
|  lig_out_electric.pdbqt | 22/02/2022 15:57 |
|  lig_out_geometric.pdbqt | 22/02/2022 15:57 |
|  lig_out_mass.pdbqt | 22/02/2022 15:56 |
|  receptor | 22/02/2022 14:32 |
|  receptor.pdbqt | 22/02/2022 14:32 |
|  vina_results_electric | 22/02/2022 15:57 |
|  vina_results_geometric | 22/02/2022 15:57 |
|  vina_results_mass | 22/02/2022 15:56 |
|  vina_simulation_electric | 22/02/2022 15:57 |
|  vina_simulation_geometric | 22/02/2022 15:57 |
|  vina_simulation_mass | 22/02/2022 15:56 |

Figure 17. Files found in the directory

C:\Users\Walter\sandres2_win\datasets\CDK2_IC50\pdbqt\1DI8 (Windows Version).

SAnDReS identifies the docking results (e.g., *vina_simulation_electric.txt*) and input files (e.g., *config_electric.txt*) for each centering with strings: electric, geometric, and mass. The file *biomatrix.csv* has the rotation matrix and translation vector used to generate the biological assembly if needed. SAnDReS creates biological units because you may have a protein for which the active site lays between units of the biological assembly e.g., human purine nucleoside phosphorylase (PNP) (de Azevedo et al., 2003). The asymmetric unit of human PNP is a monomer, and the biological assembly is a trimer. Also, the binding pocket of human PNP sits between the monomers (see structure 1PF7) (de Azevedo et al., 2003).

Now, we will carry out a statistical analysis of docking results.

Click on *Docking Hub*->*One-Click Docking with Vina (Centering: CM, EC, GC)*->*Statistical Analysis of Docking Results*. Click on the Yes button.

Once finished the statistical analysis, you have the following message:

SAnDReS finished the “Statistical Analysis of Docking Results” request using AutoDock Vina!

SAnDReS generated a file with the best docking results for each structure (*vina_rmsd_results.csv*). Figure 18 shows part of these results.

| PDB | Centering Method | RMSD(A) |
|------|------------------|---------|
| 2DS1 | electric | 0.238 |
| 2BTS | geometric | 0.514 |
| 2W17 | electric | 0.548 |
| 2W06 | electric | 0.556 |
| 3RAH | electric | 0.566 |
| 2VTS | mass | 0.655 |
| 3R9N | electric | 0.668 |
| 3LFN | mass | 0.671 |
| 3QTU | mass | 0.956 |
| 3RMF | mass | 0.963 |
| 2R3I | electric | 0.992 |
| 4GCJ | mass | 0.999 |
| 1KE6 | electric | 1.175 |
| 2VTH | electric | 1.215 |
| 2VTI | electric | 1.283 |
| 2C5Y | electric | 1.3 |
| 3QU0 | electric | 1.42 |
| 3RKB | mass | 1.444 |
| 3QTS | mass | 1.466 |
| 3WBL | geometric | 1.468 |
| 3UNJ | geometric | 1.48 |
| 1KE5 | electric | 1.606 |
| 1Y91 | geometric | 1.639 |

Figure 18. Partial view of the *vina_rmsd_results.csv* file.

As you can see, the inclusion of two alternative centering schemes for docking simulations generated better results. We obtained most of the lowest docking RMSD for electric and geometric centers. The center of mass was the best center for 40 out of 102 structures. This simple variation of the center of docking simulations improves docking accuracy (DA).

Next, we check the docking results. This part identifies unsuccessful docking simulations. We may delete these structures from the dataset or run docking simulations outside SAnDReS.

Click on *Docking Hub->Check Unsuccessful Docking Simulations*. If you have successful docking simulations, you click on the Yes option.

After checking docking simulations, you have the following message:

Number of structures with unsuccessful docking simulations: 0

We finished the docking hub part of this tutorial.

5.4. Tutorial 1. Scoring Functions

Now, we will calculate the energy terms available in AutoDock Vina 1.2.3 for the crystallographic positions and poses generated during the docking simulations. To do that, we employ the following energy terms: Gauss 1, Gauss 2, Repulsion, Hydrophobic, Hydrogen, and Torsional. We will use these energy terms as features of our machine learning models. Also, SAnDReS can calculate some descriptors such as the average charge for ligands, charge for ligands, number of atoms found in the structure of the ligands (e.g., C, N, O). It is possible to employ these descriptors to generate a hybrid scoring function involving AutoDock Vina energy terms and selected descriptors.

Click on *Scoring Functions->Enter Vina Parameters*.

We have the editor with the Vina parameters. As highlighted in the previous section, we do not change the following parameters: *centering*, *lig_pdb*, *receptor*, *ligand*, *out*, *docking_engine*, and *score_only_vina*.

Click on the *Close* button.

Click on *Scoring Functions->Add->Descriptors*.

We will have the following menu.

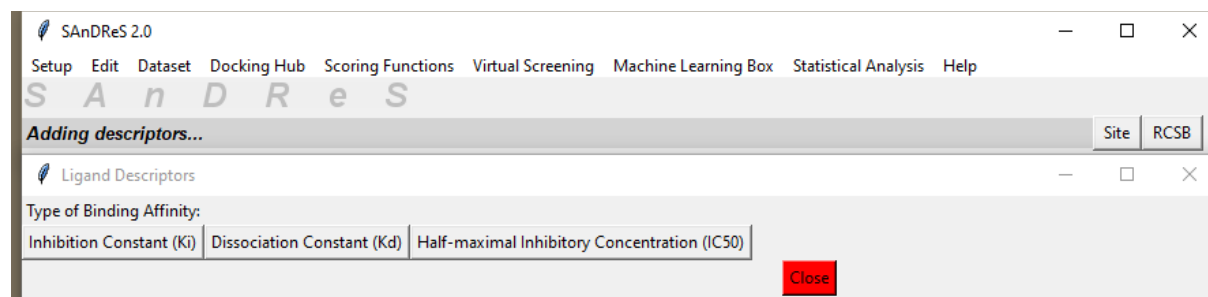


Figure 19. Ligand Descriptors Menu (Windows Version).

Click on the *IC50* button. Click on the *Yes* option.

After adding descriptors, you get the following message:

Ligand descriptors added to bind_IC50.csv file!

Click on the *Close* button. SAnDReS added new descriptors to the *bind_IC50.csv* file, as shown in Figure 20.

| Torsions | Q | Average Q | Ligand Bfactor(A2) | Receptor Bfactor(A2) | Bfactor ratio (Ligand/ C | N | O | P | S | F | Cl | Br |
|----------|--------------|--------------|--------------------|----------------------|--------------------------|----|---|---|---|---|----|----|
| 6 | -2.77556e-17 | -9.91271e-19 | 28.1364 | 33.2951 | 0.845062 | 18 | 4 | 3 | 0 | 0 | 0 | 0 |
| 5 | 0.001 | 3.125e-05 | 31.0403 | 22.077 | 1.406 | 20 | 6 | 1 | 0 | 0 | 0 | 0 |
| 3 | 2.77556e-17 | 1.54198e-18 | 35.1056 | 33.2169 | 1.05686 | 10 | 1 | 3 | 0 | 1 | 0 | 0 |
| 6 | 0.001 | 4e-05 | 20.9452 | 32.5998 | 0.642496 | 16 | 3 | 3 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 32.9291 | 45.8042 | 0.718909 | 7 | 2 | 0 | 0 | 0 | 0 | 0 |
| 4 | -0.001 | -3.33333e-05 | 39.7963 | 34.7462 | 1.14534 | 17 | 4 | 2 | 0 | 0 | 3 | 0 |
| 4 | -0.001 | -3.57143e-05 | 32.6564 | 27.7886 | 1.17518 | 17 | 4 | 2 | 0 | 0 | 1 | 0 |
| 3 | -0.002 | -8.69565e-05 | 35.8948 | 37.8154 | 0.94921 | 12 | 4 | 2 | 0 | 0 | 1 | 0 |
| 0 | -0.001 | -7.69231e-05 | 44.3138 | 36.0121 | 1.23052 | 7 | 4 | 0 | 0 | 0 | 0 | 0 |
| 8 | -1.002 | -0.0385385 | 17.9881 | 24.0363 | 0.748371 | 13 | 4 | 3 | 0 | 1 | 0 | 0 |
| 2 | -0.001 | -5.88235e-05 | 32.8071 | 32.0243 | 1.02444 | 10 | 3 | 1 | 0 | 0 | 0 | 0 |
| 7 | -0.001 | -3.7037e-05 | 19.8011 | 24.3172 | 0.814283 | 17 | 3 | 1 | 0 | 1 | 0 | 0 |
| 4 | 0.001 | 3.7037e-05 | 16.8948 | 26.7415 | 0.631784 | 14 | 4 | 3 | 0 | 1 | 0 | 0 |
| 7 | 0.001 | 3.33333e-05 | 34.016 | 27.2581 | 1.24792 | 20 | 3 | 2 | 0 | 2 | 0 | 0 |
| 5 | -0.002 | -7.40741e-05 | 17.4889 | 27.4703 | 0.636647 | 17 | 4 | 2 | 0 | 0 | 0 | 0 |
| 7 | -0.002 | -6.25e-05 | 16.2934 | 31.6888 | 0.51417 | 16 | 4 | 3 | 0 | 2 | 1 | 0 |
| 9 | -1.001 | -0.0286 | 17.6369 | 31.6028 | 0.558079 | 16 | 5 | 5 | 0 | 2 | 0 | 0 |
| 6 | -0.002 | -7.40741e-05 | 21.1659 | 33.5003 | 0.631812 | 15 | 4 | 1 | 0 | 1 | 0 | 0 |
| 5 | 2.77556e-17 | 1.06752e-18 | 23.5077 | 29.4049 | 0.799449 | 11 | 4 | 4 | 0 | 1 | 0 | 0 |
| 4 | -0.001 | -5e-05 | 31.9165 | 28.6488 | 1.11406 | 11 | 0 | 5 | 0 | 0 | 0 | 0 |
| 8 | -0.001 | -2.7027e-05 | 19.447 | 32.9704 | 0.589834 | 16 | 5 | 5 | 0 | 3 | 0 | 0 |
| 4 | 1.004 | 0.0358571 | 49.3332 | 44.9275 | 1.09806 | 16 | 4 | 0 | 0 | 1 | 3 | 0 |

Figure 20. Partial view of the bind_IC50.csv file with recent added descriptors.

Click on *Scoring Functions->Add->Energy Terms (Vina) for Crystal Structures*.

Click on the Yes option. Click on the *Run* button.

We will have the familiar menu used to run AutoDock Vina.

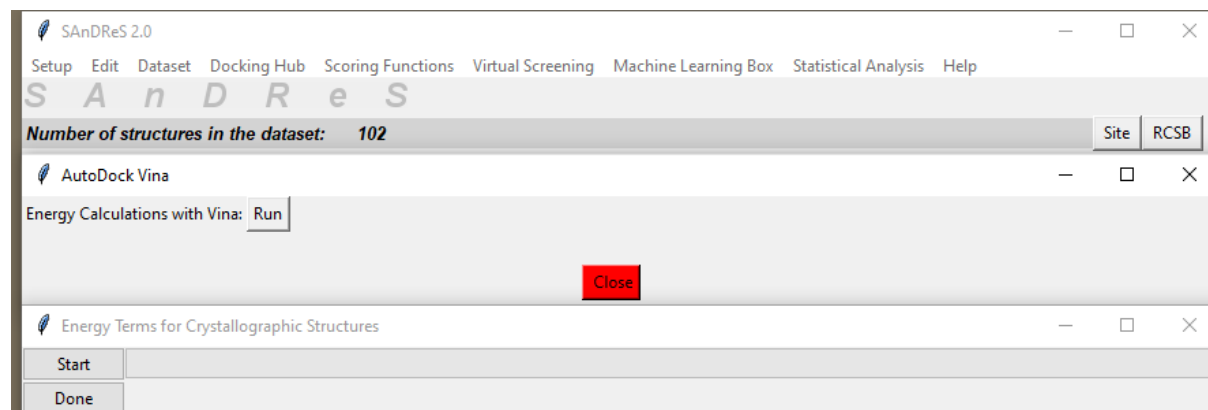


Figure 21. Energy Terms for Crystallographic Structures Menu (Windows Version).

Click on the *Start* button.

You may follow the progression.

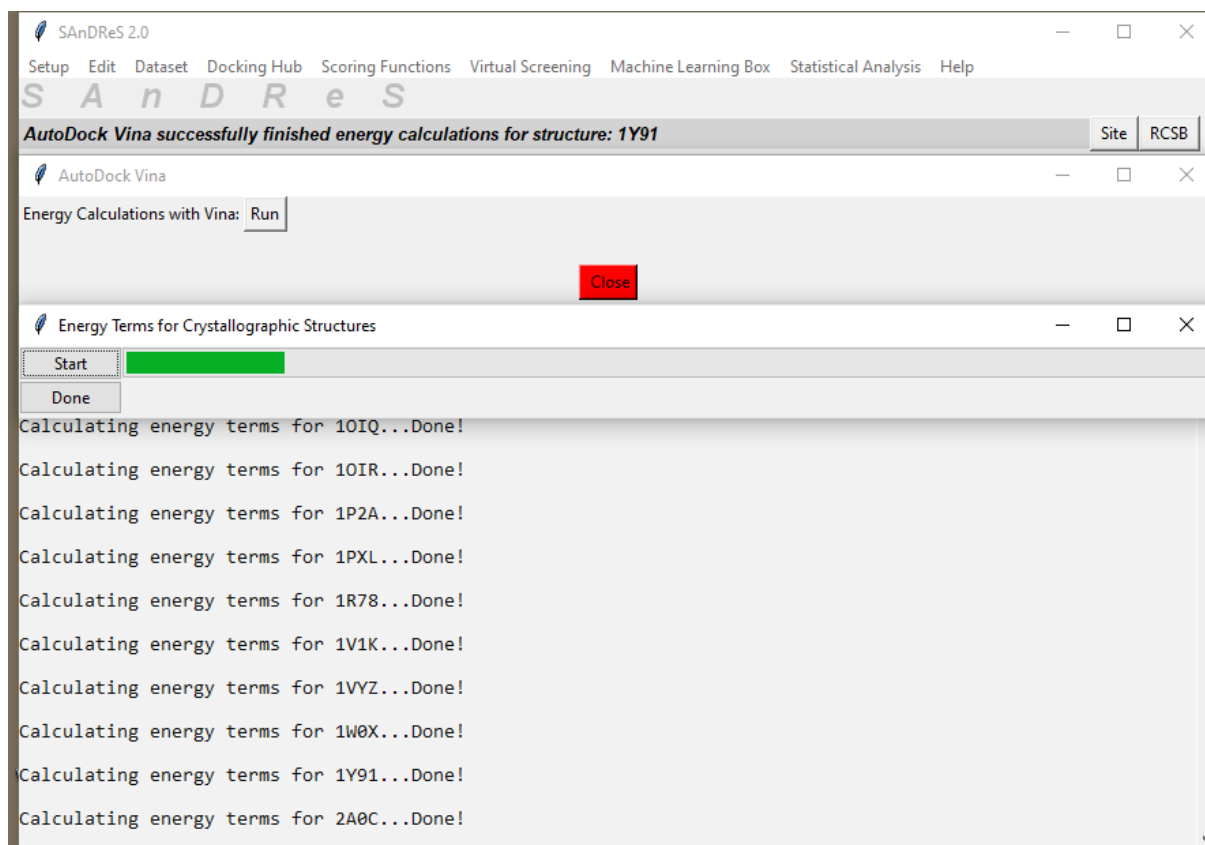


Figure 22. Energy Terms for Crystallographic Structures Menu with Anaconda prompt command (Windows Version).

After finishing all calculations, you get the following message:

SAnDReS finished the "Energy Terms" request using AutoDock Vina!

Click on the *Done* and *Close* buttons.

Click on *Scoring Functions->Add->Energy Terms (Vina) for Poses (Centering: CM, EC, GC)*.

Click on the *Yes* option. Click on the *Run* button. Click on the *Start* button.

Be careful to use the same centering option used for docking simulations. Here, we employed the centering: CM, EC, GC.

After ending the calculations, SAnDReS shows the following message:

SAnDReS finished the "Add Energy Terms (Vina) for Poses" request!

Click on the *Done* and *Close* buttons. Now, we check the energy calculations.

Click on *Scoring Functions->Check Unsuccessful Scoring Function Calculations for Poses*.

Click on the Yes option.

Then, SAnDReS shows the following message:

Done! SAnDReS updated scores4poses.csv file!

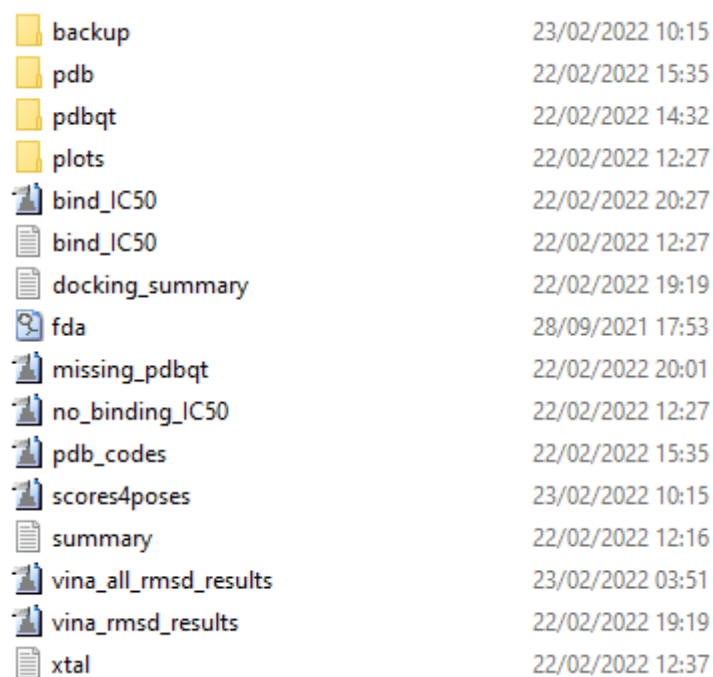
We finished this part of Tutorial 1.

5.5. Tutorial 1. Virtual Screening

Here, we will use AutoDock Vina 1.2.3 to perform a docking screen of a set of potential ligands against our protein target, the CDK2. SAnDReS checks the docking results for all structures in the dataset and selects the *config.txt* file for the best result. We consider as the best docking result the one with the lowest RMSD. SAnDReS uses the coordinates of the receptor for which we have the lowest RMSD to run the virtual screening.

To illustrate the virtual screening (VS) with SAnDReS, we will perform the simulation using only a small dataset of small molecules from the FDA-approved drugs. We may adopt this approach using the complete FDA dataset for drug repurposing purposes. Our goal is to show you how to use SAnDReS for virtual screenings.

We have the file *fda.mol2* with five molecules in the *Ligands* folder of the *sandres2_win* (or *sandres2* for Linux) directory. We will employ it for training purposes only. We may use the complete version (*fda_full.mol2*) also available in the *Ligands* folder or any other dataset of small molecules in the mol2 format. Copy the *fda.mol2* file to the project directory. Figure 23 shows the files and folders in your project directory.



| | |
|-----------------------|------------------|
| backup | 23/02/2022 10:15 |
| pdb | 22/02/2022 15:35 |
| pdbqt | 22/02/2022 14:32 |
| plots | 22/02/2022 12:27 |
| bind_IC50 | 22/02/2022 20:27 |
| bind_IC50 | 22/02/2022 12:27 |
| docking_summary | 22/02/2022 19:19 |
| fda | 28/09/2021 17:53 |
| missing_pdbqt | 22/02/2022 20:01 |
| no_binding_IC50 | 22/02/2022 12:27 |
| pdb_codes | 22/02/2022 15:35 |
| scores4poses | 23/02/2022 10:15 |
| summary | 22/02/2022 12:16 |
| vina_all_rmsd_results | 23/02/2022 03:51 |
| vina_rmsd_results | 22/02/2022 19:19 |
| xtal | 22/02/2022 12:37 |

Figure 23. Files found in the project directory

C:\Users\Walter\sandres2_win\datasets\CDK2_IC50\pdbqt\1DI8 (Windows Version).

Click on *Virtual Screening*->*Enter Virtual Screening Parameters*.

SAnDReS opens the *vs_par.csv* file, as shown in Figure 24.

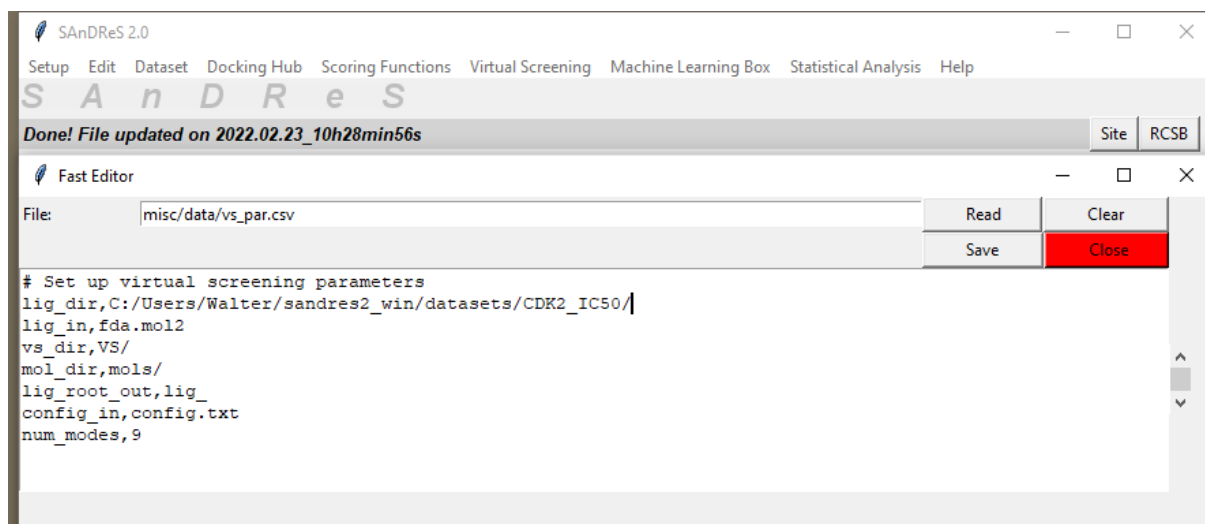


Figure 24. Fast Editor (Windows Version).

Make sure that the *lig_dir* is the project directory, and *lig_in* has *fda.mol2*, as shown above. Click on the **Save** button and the **Close** button.

Click on *Virtual Screening->Simulation->Import Ligands*.

Click on the **Yes** option.

When SAnDReS finished the conversion, we have the following message:

SAnDReS finished the "Virtual Screening->Import Ligands" request!

SAnDReS converted the molecules in the *fda.mol2* file to individual PDBQT files, one for each molecule found in the *fda.mol2* file. All PDBQT files are in the *VS/mols* folder in the project directory, as shown in Figure 25.


| | |
|---|------------------|
|  lig_1.pdbqt | 23/02/2022 10:30 |
|  lig_2.pdbqt | 23/02/2022 10:30 |
|  lig_3.pdbqt | 23/02/2022 10:30 |
|  lig_4.pdbqt | 23/02/2022 10:30 |
|  lig_5.pdbqt | 23/02/2022 10:31 |

Figure 25. Files found in the directory

C:\Users\Walter\sandres2_win\datasets\CDK2_IC50\pdbqt\1DI8\VS\mols (Windows Version).

Now, SAnDReS tests all docking results and imports the *config.txt* and *receptor.pdbqt* files of the lowest docking RMSD structure. SAnDReS will copy these files to the *VS* folder of the project directory.

Click on *Virtual Screening->Simulation->Import Receptor*.

Click on the Yes option. When SAnDReS finished the copying, we have the following message:

SAnDReS finished the "Virtual Screening->Import Receptor" request!

Now, in the VS folder, we have two additional files: *config.txt* and *receptor.pdbqt*.

You may edit the *config.txt* file, clicking on *Virtual Screening->Simulation->Edit config.txt*.

It is not necessary for this tutorial. Now, we are going to run the VS simulation.

Click on *Virtual Screening->Simulation->Run*.

Click on the Yes option. Click on the *Run* button.

As you expect, we have a familiar docking interface, as shown in Figure 26.

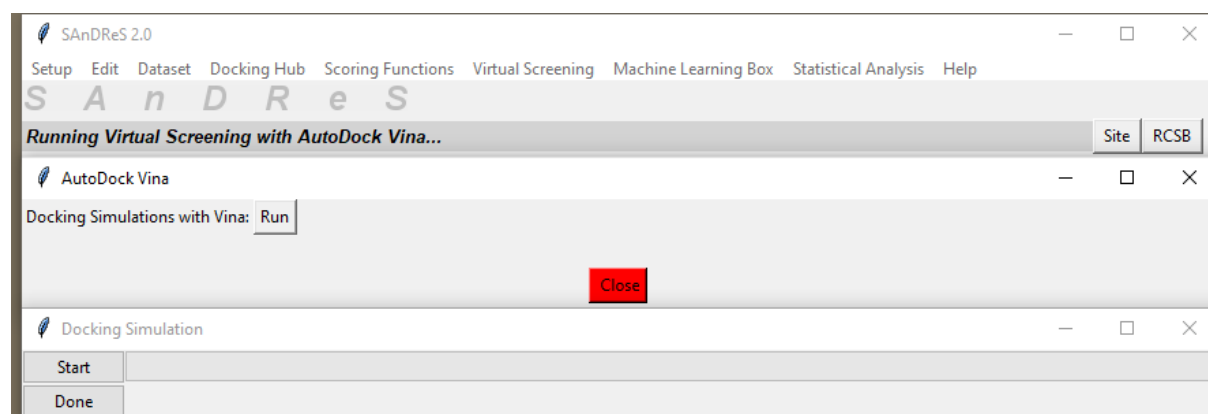


Figure 26. Docking Simulation Menu (Windows Version).

Click on the *Start* button.

After finishing the simulation, we have the following message:

SAnDReS finished the "Virtual Screening->Run" request!

Click on the *Done* button. Click on the *Close* button.

Click on *Virtual Screening->Simulation->Merge VS Results*. Click on the Yes option.

After finishing the merging, we have the following message:

SAnDReS finished the "Virtual Screening->Merge VS Results" request!

Click on *Virtual Screening->Simulation->Sort VS Results*.

Click on the Yes option.

You will have the following pop-up window.

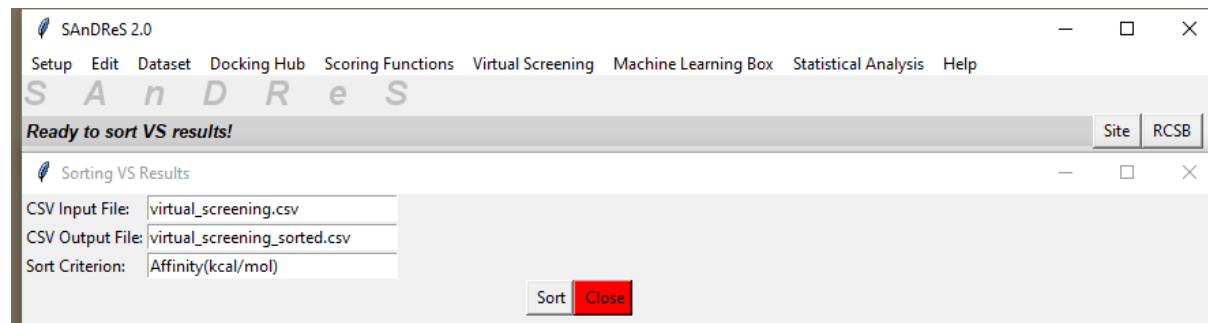


Figure 27. Sorting VS Results Menu (Windows Version).

Click on the *Sort* button.

After finishing the sorting, we have the following message:

SAnDReS finished the "Virtual Screening->Sort VS Results" request!

Click on the *Close* button.

SAnDReS generated the *virtual_screening_sorted_unified.csv* file with the poses sorted by AutoDock Vina 1.2.3 scoring function (Affinity).

Click on *Virtual Screening->Check Unsuccessful Virtual Screening Simulation*.

SAnDReS checks the VS results. If everything goes fine, you get the following message:

Number of structures with unsuccessful virtual screening simulations: 0

In this tutorial, we did not use the *Add BindingDB* option. With this option, we can carry out VS simulation for ligands downloaded from the BindingDB (Gilson et al., 2016) and add the experimental data to the VS results. Then, it is possible to verify the correlation between predicted (e.g., Affinity) and experimental binding affinities. It is also possible to employ these results to generate machine learning models based on the VS results. It is not the goal of this tutorial.

We finished out our VS simulations.

5.6. Tutorial 1. Machine Learning Box (For Modeling)

SAnDReS 2.0 relies on Scikit-Learn 1.0.2 to generate machine learning models to predict binding affinity. Here, we focus on crystallographic structures to develop our scoring functions, not docked poses. We have two sources of experimental data, the crystallographic structures and the affinity data ($\log(\text{IC}_{50})$). The $\log(\text{IC}_{50})$ is our target function. SAnDReS employs the crystal structures to calculate the energy terms. These terms are the features used in machine learning modeling.

We have an arsenal of 64 regression methods (see [Appendix](#) for details) available in SAnDReS to generate new scoring functions. This freedom to play with the features and regression methods makes it possible to explore a wider region of the Scoring Function Space (SFS) (Bitencourt-Ferreira et al., 2021), increasing the chances of finding an adequate model for our protein system. Although in computational modeling, some argue that the data is more relevant than the machine learning algorithms for complex problems (Halevy et al., 2009). Several studies showed that variation of the machine learning algorithms generates scoring functions with superior predictive power compared with classical scoring functions (Xavier et al., 2016; de Ávila et al., 2017; Pinto et al., 2017; Bitencourt-Ferreira & de Azevedo, 2018; Levin et al., 2018; Wójcikowski et al., 2019; Ballester, 2019; da Silva et al., 2020; Bitencourt-Ferreira et al., 2021; de Azevedo et al., 2021; de Azevedo, 2021). These results highlight the importance of trying different methodologies when studying complex systems. In SAnDReS, we adopt this idea. The idea of freedom to explore unexplored regions of the SFS, finding a model just right for your protein system.

One key aspect of SAnDReS is the freedom of choice. You have the liberty to test several machine learning models, and based on the predictive performance furnished by SAnDReS, choose the model you find is adequate to the protein you are studying.

To start exploring the SFS, click on *Machine Learning Box*->*Enter Machine Learning Parameters*.

SAnDReS opens the *ml_par.csv* file. In this file, we have the definition of the parameters necessary to apply the regression methods available in Scikit-Learn 1.0.2. We show part of the *ml_par.csv* file below.

```
preprocessing,StandardScaler
ml_parameters,ml.in
scoring_function_file,scores.csv
# Set up input parameters
n_features,6
features_in,Gauss 1,Gauss 2,Repulsion,Hydrophobic,Hydrogen,Torsional
test_size_in,0.3
seed_in,271828
# Set up regression methods
mlr_method,AdaBoostRegressor                                # AdaBoost Regression
mlr_method,AdaBoostRegressorCV                               # AdaBoost Regression
```

SAnDReS considers any line starting with # as a comment line. If it appears in the middle of a line, from that the part on is a comment.

The line `preprocessing,StandardScaler` defines the type of preprocessing. SAnDReS will scale the data. In the following, we have the definition of the input file with all specific parameters necessary to run machine learning methods available in Scikit-Learn 1.0.2. Please, see Appendix for details.

The line `scoring_function_file,scores.csv` shows a temporary file used during regression analysis. You do not need to change it. The next line starts with `#`. It is a comment line.

The line `n_features,6` establishes the number of features for machine learning methods. Once started the machine learning, you cannot change it. We will leave as it is.

The line `features_in,Gauss 1,Gauss 2,Repulsion,Hydrophobic,Hydrogen,Torsional` defines the features. In this case, we use the energy terms available in the AutoDock Vina scoring function. If you want to modify it, you must do all steps from this point again.

SAnDReS splits the data in training and test sets. The following line `test_size_in,0.3` defines the fraction used as test set. In this tutorial, we have 30 % of the dataset as test set.

The line `seed_in,271828` defines an integer used as a seed to generate pseudorandom numbers. SAnDReS employs these pseudorandom numbers to split the dataset into training and test sets. This definition allows us to reproduce the same results for the same operating system. For this tutorial, leave it as it is. The next line is a comment line. In the following, we have a sequence of 64 regression methods specified in each line, for instance, `mlr_method,AdaBoostRegressor` indicates that SAnDReS will use the AdaBoostRegressor as a method for regression. In the above list, we have the first two lines out of 64 representing all regression methods available in SAnDReS. If you want to omit any of the methods, it is necessary to add `#` as the first character in the line. You may also just delete the undesired method.

Click on the *Save* button and the *Close* button.

Click on *Machine Learning Box->For Modeling->Preprocess Data*.

Click on the *Yes* option.

After finishing preprocessing the data, we get the following message:

SAnDReS finished the "Preprocess Data" request!

SAnDReS generated a new file named *scores4xtal.csv* with scaled data. From now on, the machine learning modeling will use the data in this file. The use the options *Choose PDBs for Training Set* and *Choose PDBs for Test Set* to open a file using the Fast Editor. We may use these options to enter the PDB access codes you chose. We do not use them in this tutorial.

Click on *Machine Learning Box->For Modeling->Automatic Generation of PDBs for Training and Test Sets*.

Click on the Yes option. After separating the PDB access codes in the files *pdb_codes_test_set.csv* and *pdb_codes_training_set.csv*, we have the following message:

SAnDReS finished the "Automatic Generation of PDBs for Training and Test Sets" request!

Click on *Machine Learning Box->For Modeling->Generate Training and Test Sets*.

Click on the Yes option. Now we split the dataset. SAnDReS creates two new files, named: *scores4xtal_test.csv* and *scores4xtal_training.csv*. We have the following message after the creation of both files:

SAnDReS finished the "Generate Training and Test Sets" request!

Click on *Machine Learning Box->For Modeling->Filter Data*.

Click on the Yes option three times for the following files: *scores4xtal.csv*, *scores4xtal_test.csv*, and *scores4xtal_training.csv*. In this part, SAnDReS eliminates any line for which we have features with *nan* (not a number).

After eliminating these lines, if necessary, SAnDReS shows the following message:

SAnDReS finished the "Filter Data" request!

Now we determine the correlation between potential features and $\log(\text{IC}_{50})$.

Click on *Machine Learning Box->For Modeling->Statistical Analysis (Features)*.

Click on the Yes option. After generating the *scores4xtal_training_stats4features.csv* file, we have the following message:

SAnDReS finished the "Statistical Analysis (Features)" request!

In this new file, we have the bivariate statistical analysis of potential features carried out using Scikit-Learn and SciPy. This file is sorted from the highest r^2 to the lowest. Figure 28 shows the *scores4xtal_training_stats4features.csv* file.

| Feature | r | p-value | r2 | rho | p-value1 | MSE | RMSE | RSS |
|-------------|-------------|-------------|-------------|------------|-------------|---------|---------|---------|
| C | -0.418426 | 0.000254526 | 0.17508 | -0.393667 | 0.000623522 | 45.0824 | 6.71434 | 3245.93 |
| Gauss 2 | -0.375868 | 0.00113905 | 0.141277 | -0.360388 | 0.00187285 | 45.1341 | 6.71819 | 3249.66 |
| Average Q | -0.373714 | 0.00122244 | 0.139662 | -0.482266 | 1.78911e-05 | 44.0001 | 6.63326 | 3168.01 |
| Q | -0.372286 | 0.00128075 | 0.138597 | -0.441963 | 0.000101619 | 43.9578 | 6.63006 | 3164.96 |
| Gauss 1 | -0.25894 | 0.0280684 | 0.0670498 | -0.246267 | 0.0370412 | 44.8517 | 6.69714 | 3229.32 |
| Torsional | -0.252748 | 0.0321922 | 0.0638817 | -0.181425 | 0.127213 | 44.1103 | 6.64156 | 3175.94 |
| N | -0.234687 | 0.0472174 | 0.0550782 | -0.197267 | 0.0967256 | 44.5865 | 6.67731 | 3210.23 |
| Hydrophobic | -0.224626 | 0.0578311 | 0.0504568 | -0.226039 | 0.0562333 | 44.11 | 6.64154 | 3175.92 |
| Torsions | -0.221948 | 0.0609616 | 0.0492609 | -0.181425 | 0.127213 | 44.0337 | 6.63579 | 3170.42 |
| S | 0.117337 | 0.326293 | 0.0137679 | 0.131486 | 0.270915 | 44.0421 | 6.63642 | 3171.03 |
| Hydrogen | -0.102113 | 0.393369 | 0.010427 | -0.0811614 | 0.497934 | 44.8853 | 6.69965 | 3231.74 |
| O | -0.0678274 | 0.571316 | 0.00460056 | -0.0602544 | 0.615116 | 43.7967 | 6.61791 | 3153.36 |
| F | -0.00800414 | 0.946797 | 6.40662e-05 | -0.0903209 | 0.450531 | 43.8646 | 6.62304 | 3158.25 |
| Repulsion | -0.00483638 | 0.967838 | 2.33906e-05 | -0.102447 | 0.391815 | 44.1789 | 6.64672 | 3180.88 |

Figure 28. View of the scores4xtal_training_stats4features.csv file.

We observe the highest r^2 for the following features: C, Gauss 2, Average Q, Q, Gauss 1, Torsional, N.

We have a high correlation between Q and average Q. So, it is better to choose one of them. We keep Average Q. So, the list of features is the following:

C, Gauss 2, Average Q, Gauss 1, Torsional, N

Click on *Machine Learning Box*->*Enter Machine Learning Parameters*.

We have the following view after inserting the selected features.

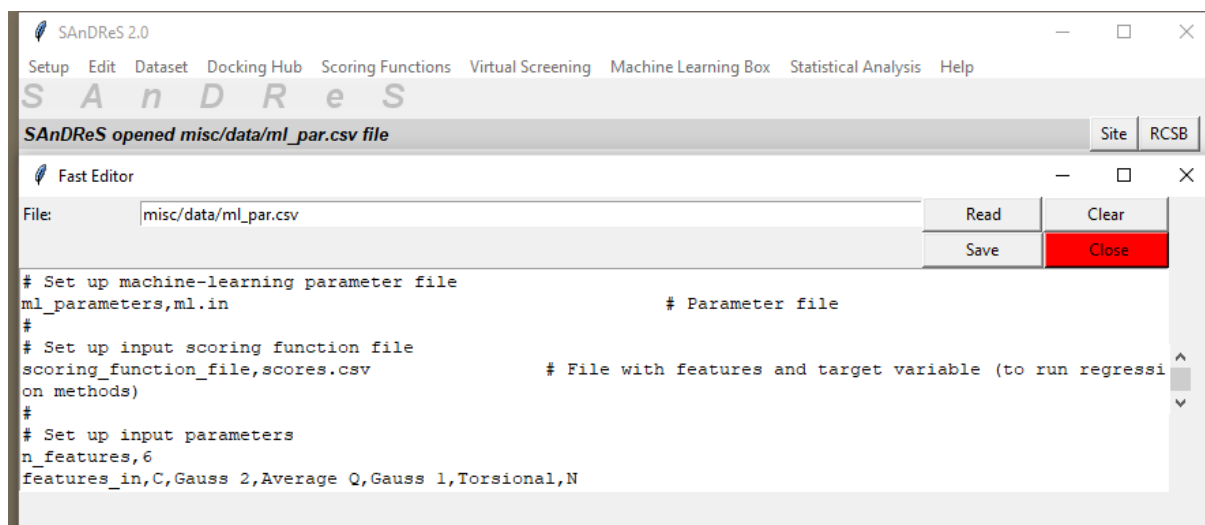


Figure 29. Fast Editor (Windows Version).

Click on the *Save* button and the *Close* button.

We should keep in mind that this selection of features does not mean that we will generate the best possible model. We may try different scenarios and choose the best-performing model. It is common to select the model that generalizes better (de Azevedo, 2021). It means that the model has the highest correlation between experimental and predicted affinity for the test set. For half of the regression methods available in SAnDReS, we have cross-validation (CV) approaches.

Click on *Machine Learning Box->For Modeling->Regression Methods*.

Click on the *Yes* option.

SAnDReS started our exploration of the SFS. It goes into a loop, generating machine learning models for all methods we kept in the previously edited *ml_par.csv* file. For this tutorial, we keep them all. For each regression method, SAnDReS generates a model stored in the *models* folder of the project directory. We have these models saved as *joblib* files. These files allow us to apply previously generated models (*.joblib* file) to any dataset presented as a CSV file. Once generated a machine learning model (*.joblib* file), we may copy it, keep it on a website, or send it by email. Then, we can use it to predict binding affinity using as input data any CSV file that has the same features used to generate the machine learning. In this tutorial, we used the following features: C, Gauss 2, Average Q, Gauss 1, Torsional, N.

After generating all regression models, we have the following message:

SAnDReS finished the "Regression Methods" request!

In Figure 30, we have part of the files found in the *models* folder.












| | |
|---|------------------|
|  features | 23/02/2022 12:27 |
|  model_AdaBoostRegressor.joblib | 23/02/2022 12:23 |
|  model_AdaBoostRegressorCV.joblib | 23/02/2022 12:23 |
|  model_ARDRRegression.joblib | 23/02/2022 12:23 |
|  model_ARDRRegression.json | 23/02/2022 12:23 |
|  model_ARDRRegression | 23/02/2022 12:23 |
|  model_ARDRRegressionCV.joblib | 23/02/2022 12:23 |
|  model_ARDRRegressionCV.json | 23/02/2022 12:23 |
|  model_ARDRRegressionCV | 23/02/2022 12:23 |
|  model_BaggingRegressor.joblib | 23/02/2022 12:25 |
|  model_BaggingRegressorCV.joblib | 23/02/2022 12:25 |

Figure 30. Part of the files found in the directory

C:\Users\Walter\sandres2_win\datasets\CDK2_IC50\models (Windows Version).

Next, we will use these regression models (.joblib files) and apply them to *scores4xtal_test.csv* and *scores4xtal_training.csv* files. SAnDReS will add the binding affinity values predicted using all 64 regression models as additional columns to *scores4xtal_test.csv* and *scores4xtal_training.csv* files. We name all added columns as regression methods, e.g., AdaBoostRegressor for a column with binding affinity determined using this method.

Click on *Machine Learning Box-> For Modeling->Apply Regression Model*.

We have the following pop-up window.

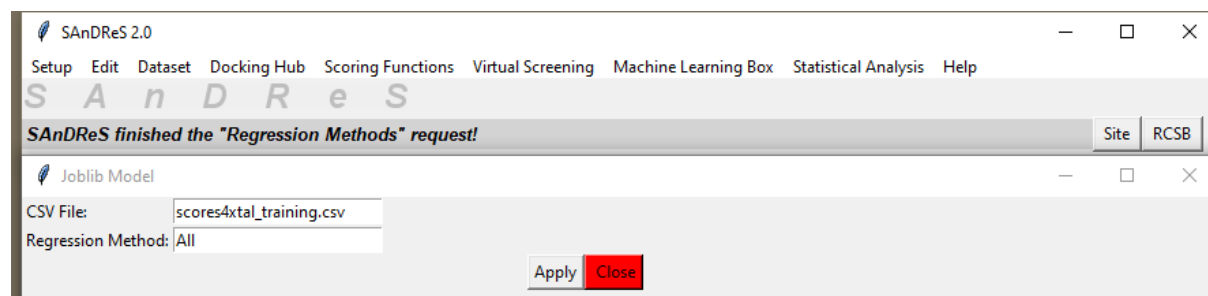


Figure 31. Joblib Model Menu (Windows Version).

Click on the *Apply* button. After applying all regression models to the *scores4xtal_training.csv* file, we have the following message:

SAnDReS finished the "Apply Regression Model" request!

Click on the *Close* button and repeat the same procedure for the *scores4xtal_test.csv* file. In the field *Regression Method* of the above pop-up window, we have the possibility of inserting a specific method, instead of asking to add for all methods. In this tutorial, we keep them all.

Click on *Machine Learning Box->For Modeling->Bivariate Analysis of Regression Models*.

Click on the *Yes* option.

After finishing the statistical analysis, we get the following message:

SAnDReS finished the "Bivariate Analysis of Regression Models" request!

SAnDReS performed the statistical analysis and generated the following files: *scores4xtal_test_stats_joblib_models.csv* and *scores4xtal_training_stats_joblib_models.csv*.

Since we seek models with generalization capacity, we check the *scores4xtal_test_stats_joblib_models.csv* file. Figure 32 shows part of this file.

| Method | r | p-value(r) | r2 ▼ | rho | p-value(rho) | MSE | RMSE |
|-----------------------------|-----------|------------|-----------|-------------|--------------|---------|---------|
| KemelRidge | 0.343069 | 0.0634535 | 0.117696 | 0.315087 | 0.0898899 | 41.4289 | 6.43653 |
| ExtraTreeRegressorCV | -0.249596 | 0.183463 | 0.0622979 | -0.195413 | 0.300735 | 3.90047 | 1.97496 |
| BaggingRegressor | 0.24368 | 0.194411 | 0.0593798 | 0.0745438 | 0.695435 | 1.73731 | 1.31807 |
| ElasticNet | 0.241245 | 0.199046 | 0.0581992 | 0.279745 | 0.134342 | 1.50202 | 1.22557 |
| OrthogonalMatchingPursuit | 0.241245 | 0.199046 | 0.0581992 | 0.279745 | 0.134342 | 1.45488 | 1.20618 |
| DecisionTreeRegressor | 0.21417 | 0.255763 | 0.0458689 | 0.203229 | 0.281424 | 2.91911 | 1.70854 |
| Affinity | 0.209736 | 0.265971 | 0.0439893 | 0.314197 | 0.0908504 | 2.60701 | 1.61462 |
| ElasticNetCV | 0.208295 | 0.269346 | 0.0433866 | 0.205054 | 0.277034 | 1.50054 | 1.22496 |
| TweedieRegressor | 0.206358 | 0.273924 | 0.0425836 | 0.239875 | 0.201686 | 1.48107 | 1.21699 |
| ExtraTreeRegressor | -0.204803 | 0.277636 | 0.0419441 | -0.247994 | 0.186383 | 4.03134 | 2.00782 |
| VotingRegressor | 0.182242 | 0.335112 | 0.033212 | 0.168669 | 0.372933 | 1.65289 | 1.28565 |
| ARDRegression | 0.17893 | 0.344117 | 0.0320159 | 0.252781 | 0.177749 | 1.58268 | 1.25805 |
| RandomForestRegressor | 0.174118 | 0.357459 | 0.0303172 | 0.093903 | 0.621613 | 1.86038 | 1.36396 |
| GradientBoostingRegressor | 0.165686 | 0.38157 | 0.0274517 | -0.00956831 | 0.959977 | 2.25801 | 1.50267 |
| GradientBoostingRegressorCV | 0.162782 | 0.390086 | 0.0264979 | 0.0658656 | 0.729488 | 2.10684 | 1.4515 |
| BayesianRidge | 0.161061 | 0.395185 | 0.0259405 | 0.232532 | 0.216255 | 1.55663 | 1.24765 |
| TweedieRegressorCV | 0.15926 | 0.40056 | 0.0253637 | 0.215621 | 0.252481 | 1.51642 | 1.23143 |
| AdaBoostRegressorCV | 0.154691 | 0.414382 | 0.0239294 | 0.139183 | 0.463239 | 1.98737 | 1.40974 |
| LassoLarsIC | 0.154333 | 0.415478 | 0.0238185 | 0.212728 | 0.259055 | 1.59118 | 1.26142 |
| LinearRegression | 0.151784 | 0.423315 | 0.0230383 | 0.227637 | 0.226355 | 1.63407 | 1.27831 |
| Lars | 0.151784 | 0.423315 | 0.0230383 | 0.227637 | 0.226355 | 1.63407 | 1.27831 |
| Ridge | 0.147166 | 0.437719 | 0.0216578 | 0.218514 | 0.246017 | 1.63138 | 1.27726 |
| RandomForestRegressorCV | 0.140335 | 0.459506 | 0.019694 | 0.0300401 | 0.874788 | 1.8961 | 1.37699 |

Figure 32. Partial view of the `scores4xtal_test_stats_joblib_models.csv` file.

Above, we sorted the data using r^2 as a sorting criterion. We see several models performing better than the AutoDock Vina scoring function (Affinity). We select the model created using the OrthogonalMatchingPursuit method since this model has good performance for RMSE (root-mean-square error).

Note. SAnDReS may show nan (not a number) for statistical analysis of some regression methods. It happens due to errors generated in the Scikit-Learn. If you have only a few methods for this problem, you may ignore them. Otherwise, you may have to change regression parameters. Please see [Appendix](#).

Now, we have our machine learning model (OrthogonalMatchingPursuit). We will test it against docking results.

5.7. Tutorial 1. Machine Learning Box (For Docking Results)

In this part, we will use the previously generated docking results, apply our machine learning models (KernelRidge) against them, and carry out the same bivariate statistical analysis performed for the models in the previous section. **Here, SAnDReS will not perform machine learning regression, only apply the previously determined models (.joblib files) to the docking results.**

Click on *Machine Learning Box->For Docking Results->Preprocess Data*.

Click on the Yes option. After finishing preprocessing the data, we get the following message:

SAnDReS finished the "Preprocess Data" request!

SAnDReS generated a new file named *scores4poses.csv* with scaled data.

Click on *Machine Learning Box->For Docking Results->Generate Training and Test Sets*.

Click on the Yes option. SAnDReS generated two new files, named: *scores4poses_test.csv* and *scores4poses_training.csv*. We have the following message after the creation of both files:

SAnDReS finished the "Generate Training and Test Sets" request!

Click on *Machine Learning Box->For Docking Results ->Filter Data*.

Click on the Yes option three times for the following files: *scores4poses.csv*, *scores4poses_test.csv*, and *scores4poses_training.csv*. Once finished, we have the following message:

SAnDReS finished the "Filter Data" request!

Click on *Machine Learning Box->For Docking Results->Apply Regression Model*.

As we saw in the previous section, we have a pop-up window. Initially, we apply the regression models to the *scores4poses.csv* file. We click on the *Apply* button and then on the *Close* button. Repeat the same procedure for the following files: *scores4poses_training.csv* and *scores4poses_test.csv*. Different from the previous section, we applied the models to three files. After the inclusion of all regression models, we have the following message:

SAnDReS finished the "Apply Regression Model Docking Results" request!

Click on *Machine Learning Box*->*For Docking Results*->*Bivariate Analysis of Regression Models*.

Click on the Yes option. After finishing the statistical analysis, we get the following message:

SANdReS finished the "Bivariate Analysis of Regression Models Applied Docking Results" request!

Figure 33 shows a partial view of the *scores4poses_test_stats_joblib_models.csv* file.

| Method | Mean RMSD(A) | Minimum RMSD(A) | Maximum RMSD(A) | DA1 | DA2 | r | p-value | r2 | rho | p-value1 | MSE | RMSE |
|-----------------------------|---------------|-----------------|-----------------|------------|-----------------|-----------------|-----------------|----------------|-----------------|------------------|----------------|----------------|
| OrthogonalMatchingPursuit | 4.3388 | 0.786 | 8.952 | 0.333333 | 0.35 | 0.241245 | 0.199046 | 0.0581992 | 0.279745 | 0.134342 | 1.45963 | 1.20815 |
| ElasticNet | 4.3388 | 0.786 | 8.952 | 0.333333 | 0.35 | 0.241245 | 0.199046 | 0.0581992 | 0.279745 | 0.134342 | 1.50209 | 1.2256 |
| LassoLarsICCV | 4.3388 | 0.786 | 8.952 | 0.333333 | 0.35 | 0.0488957 | 0.797498 | 0.00239079 | 0.0814419 | 0.668773 | 1.72722 | 1.31424 |
| ElasticNetCV | 4.3388 | 0.786 | 8.952 | 0.333333 | 0.35 | 0.209283 | 0.267028 | 0.0437995 | 0.205054 | 0.277034 | 1.50063 | 1.225 |
| OrthogonalMatchingPursuitCV | 4.3388 | 0.786 | 8.952 | 0.333333 | 0.35 | 0.0768178 | 0.686604 | 0.00590098 | 0.0407025 | 0.830896 | 1.62969 | 1.27659 |
| XGBRegressorCV | 4.3893 | 0.786 | 8.952 | 0.316667 | 0.333333 | -0.0755107 | 0.691675 | 0.00570187 | -0.0073918 | 0.969076 | 2.17656 | 1.47532 |
| DecisionTreeRegressorCV | 4.55907 | 0.786 | 7.409 | 0.266667 | 0.266667 | 0.388618 | 0.0338102 | 0.151024 | 0.36237 | 0.0490786 | 2.46213 | 1.56912 |
| AdaBoostRegressorCV | 5.28283 | 0.786 | 8.952 | 0.2 | 0.2 | 0.181878 | 0.336093 | 0.0330797 | 0.205656 | 0.275595 | 1.84508 | 1.35834 |
| Affinity | 4.9303 | 1.107 | 8.951 | 0.2 | 0.208333 | 0.231063 | 0.219254 | 0.05339 | 0.353995 | 0.0549602 | 7.55667 | 2.74876 |
| TwieedieRegressor | 5.25407 | 0.671 | 8.575 | 0.183333 | 0.191667 | 0.207206 | 0.271913 | 0.0429344 | 0.251224 | 0.180527 | 1.50589 | 1.22715 |
| DecisionTreeRegressor | 4.99233 | 0.786 | 7.469 | 0.183333 | 0.183333 | 0.343598 | 0.0630201 | 0.118059 | 0.291484 | 0.11809 | 2.53351 | 1.5917 |
| SVRCV | 5.39613 | 0.671 | 8.951 | 0.183333 | 0.191667 | 0.0152374 | 0.936303 | 0.000232178 | 0.103026 | 0.58798 | 2.02847 | 1.42424 |
| KernelRidge | 5.38403 | 0.671 | 8.951 | 0.183333 | 0.191667 | 0.34813 | 0.0594003 | 0.121194 | 0.308411 | 0.0972849 | 39.9541 | 6.32092 |
| XGBRegressor | 5.07797 | 0.786 | 8.952 | 0.183333 | 0.2 | -0.0885664 | 0.641644 | 0.007844 | 0.0138751 | 0.941988 | 2.19991 | 1.48321 |
| BayesianRidgeCV | 5.25993 | 0.671 | 8.575 | 0.183333 | 0.191667 | 0.0888081 | 0.640731 | 0.00788688 | 0.179795 | 0.34175 | 1.65975 | 1.28831 |
| ARDRegressionCV | 5.3873 | 0.671 | 8.575 | 0.15 | 0.158333 | 0.0414042 | 0.828024 | 0.00171431 | 0.0892301 | 0.639139 | 1.7567 | 1.32541 |
| LinearRegressionCV | 5.38593 | 0.671 | 8.575 | 0.15 | 0.158333 | 0.0207846 | 0.91319 | 0.000431999 | 0.104139 | 0.583934 | 1.83265 | 1.35376 |
| TwieedieRegressorCV | 5.3873 | 0.671 | 8.575 | 0.15 | 0.158333 | 0.156347 | 0.409342 | 0.0244444 | 0.198932 | 0.29194 | 1.54203 | 1.24179 |
| RidgeCV | 5.41377 | 0.671 | 8.575 | 0.15 | 0.158333 | 0.0271509 | 0.88675 | 0.000737169 | 0.117713 | 0.535586 | 1.81412 | 1.34689 |
| LarsCV | 5.35873 | 0.671 | 8.575 | 0.15 | 0.158333 | 0.0156111 | 0.934744 | 0.000243707 | 0.0772141 | 0.68507 | 1.83453 | 1.35445 |
| SGDRegressorCV | 5.36953 | 0.671 | 8.575 | 0.133333 | 0.141667 | 0.00559779 | 0.976579 | 3.13352e-05 | 0.105919 | 0.577485 | 1.84241 | 1.35735 |
| SGDRegressor | 5.631 | 0.786 | 8.951 | 0.133333 | 0.133333 | 0.159703 | 0.399234 | 0.0255051 | 0.228972 | 0.22357 | 1.62559 | 1.27499 |

Figure 33. Partial view of the *scores4poses_test_stats_joblib_models.csv* file.

As we can see, our machine learning model (OrthogonalMatchingPursuit) performs better than the Affinity function taking docking accuracy (DA1) (Xavier et al., 2016) as a sorting criterion. The performance is also better for r^2 and RMSE. Amongst all these metrics, we may say that the most demanding for a machine learning model trained against binding affinity data and crystal structures are the DAs, and our model has an improvement of almost 13.33 % for DA1. Additionally, we may take a machine learning model to predict binding affinity and a different model to sort docking poses. In summary, the overall performance of the OrthogonalMatchingPursuit model is better for sorting poses and predicting the binding affinity.

We finished our application of machine learning models to docking results.

5.8. Tutorial 1. Machine Learning Box (For Virtual Screening Results)

In this part, we will apply the `OrthogonalMatchingPursuit` model to the results of our VS simulation. We will employ one regression model only. Our goal is to use a machine learning model to sort VS poses, selecting the most promising ones (lowest score). Remember, we have only five ligands in the dataset used for VS. We intend only to show you how to apply our approach to this problem. For real VS projects, we need larger datasets, but the procedure is the same.

Click on *Machine Learning Box->For Virtual Screening Results->Preprocess Data*.

Click on the Yes option.

After finishing preprocessing the data, we get the following message:

SAnDReS finished the "Preprocess Data for Virtual Screening Results" request!

Click on *Machine Learning Box->For Virtual Screening Results->Apply Regression Model*.

In the new pop-up window, we change the regression method to *OrthogonalMatchingPursuit*. We have the screen shown below.

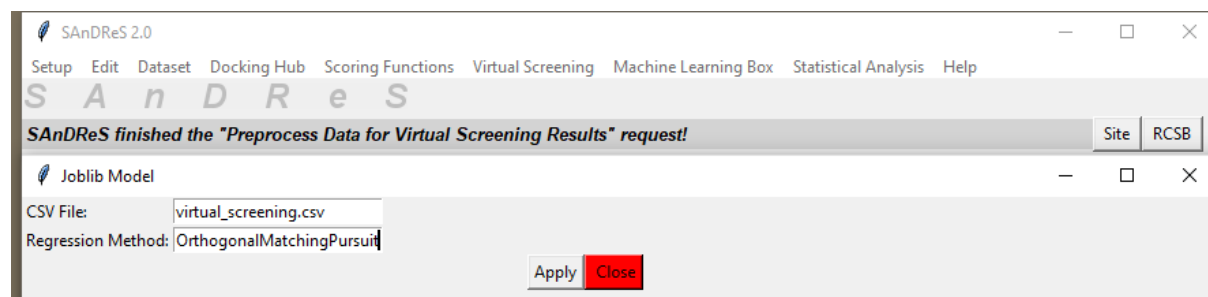


Figure 34. Joblib Model Menu (Windows Version).

Click on the *Apply* button and the Yes option. You get the following message:

SAnDReS finished the "Apply Regression Model Virtual Screening Results" request!

Click on the *Close* button.

Click on *Machine Learning Box->For Virtual Screening Results->Sort VS Results*.

Click on the Yes option then we get a new pop-up window. In the Sorting Criterion field, type `OrthogonalMatchingPursuit`

We have the following window.

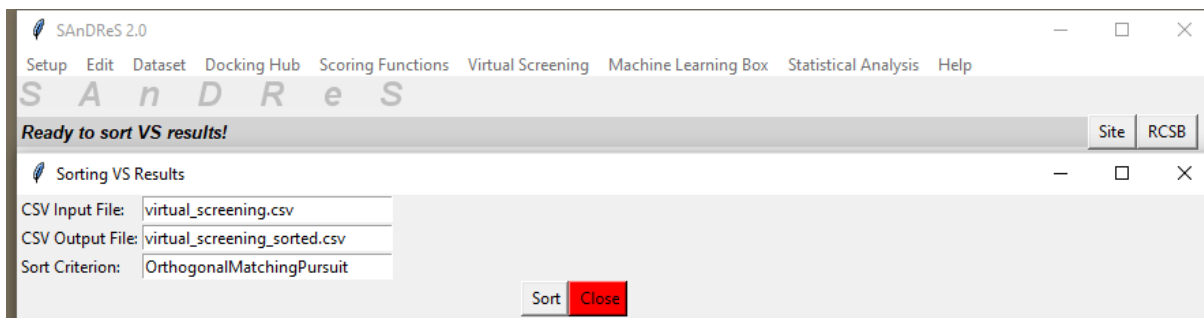


Figure 35. Sorting VS Results (Windows Version).

Click on the *Sort* button.

You get the following message:

Finished sorting VS results!

Click on the *Close* button. SAnDReS added the OrthogonalMatchingPursuit predicted values to the last column. Figure 36 shows part of the *virtual_screening_sorted_unified.csv* file.

| PDB | Ligand | Chain | Number | Affinity(kcal/mol) | Gauss1 | Gauss2 | Repulsion | Hydrophobic | Hydrogen | Torsional | OrthogonalMatchingPursuit |
|------|--------|--------------|--------|--------------------|------------|-----------|-----------|-------------|-----------|-----------|---------------------------|
| 2DS1 | lig_1 | ZINC08034121 | pose_1 | -2.166 | 1.00996 | -0.999006 | 0.22132 | -0.468982 | 0.464206 | 1.60357 | -6.81956 |
| 2DS1 | lig_3 | ZINC05224164 | pose_3 | -2.52 | 1.31823 | 0.294739 | -0.716334 | 1.95309 | -0.850125 | -1.06904 | -6.81956 |
| 2DS1 | lig_5 | ZINC05224188 | pose_8 | -2.014 | -0.0804938 | 0.692577 | -0.876834 | -0.468982 | -0.850125 | 0.267261 | -6.81956 |
| 2DS1 | lig_4 | ZINC08101126 | pose_9 | -2.146 | -0.190161 | 0.804652 | -0.585401 | -0.468982 | -0.557491 | -1.06904 | -5.74221 |
| 2DS1 | lig_2 | ZINC06827693 | pose_8 | -1.832 | -1.72682 | 1.30928 | -1.04156 | -0.468982 | -0.850125 | 0.267261 | -5.74221 |

Figure 36. Partial view of the *virtual_screening_sorted_unified.csv* file.

In the above figure, we see in the last column the predicted binding affinity generated with OrthogonalMatchingPursuit method.

We finished the application of the machine learning model to VS results.

5.9. Tutorial 1. Statistical Analysis

In the last part of tutorial 1, we will perform some statistical analysis and generate plots. SAnDReS keeps all graphic files in the *plots* folder of the project directory. To start, we will create a scatter plot.

Click on *Statistical Analysis->Scatter Plot->Set up Parameters*.

We have a new pop-up window. Add *OrthogonalMatchingPursuit* to the Y-axis Label field. We have the following window.

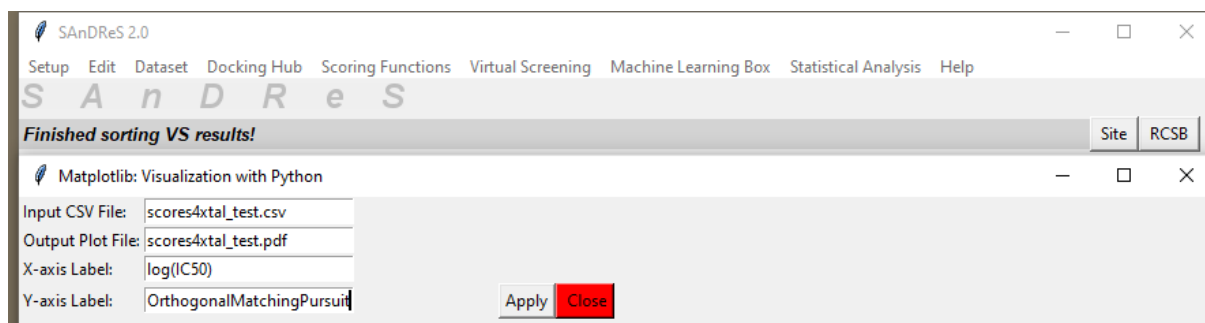


Figure 37. Matplotlib: Visualization with Python Menu (Windows Version).

Click on the *Apply* button. Click on the *Close* button.

To generate the plot, click on *Statistical Analysis->Scatter Plot->Generate*.

We have the following pop-up window.

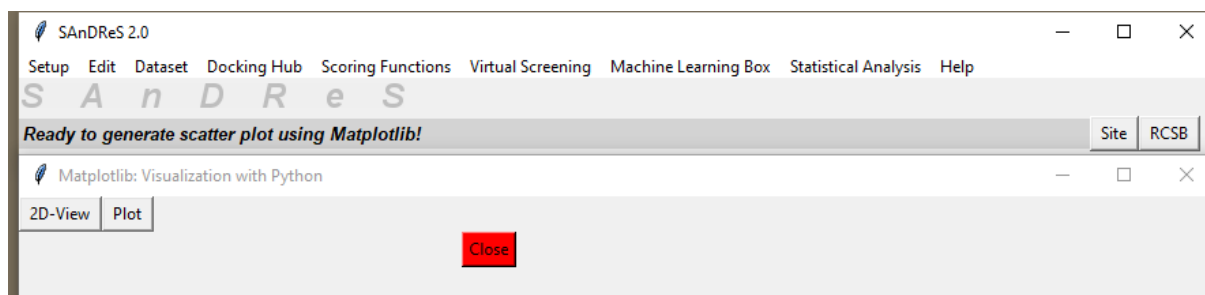


Figure 38. Matplotlib: Visualization with Python Menu (Windows Version).

Click on the *2D-View* button. Click on the *Close* button.

SAnDReS generated the following plot.

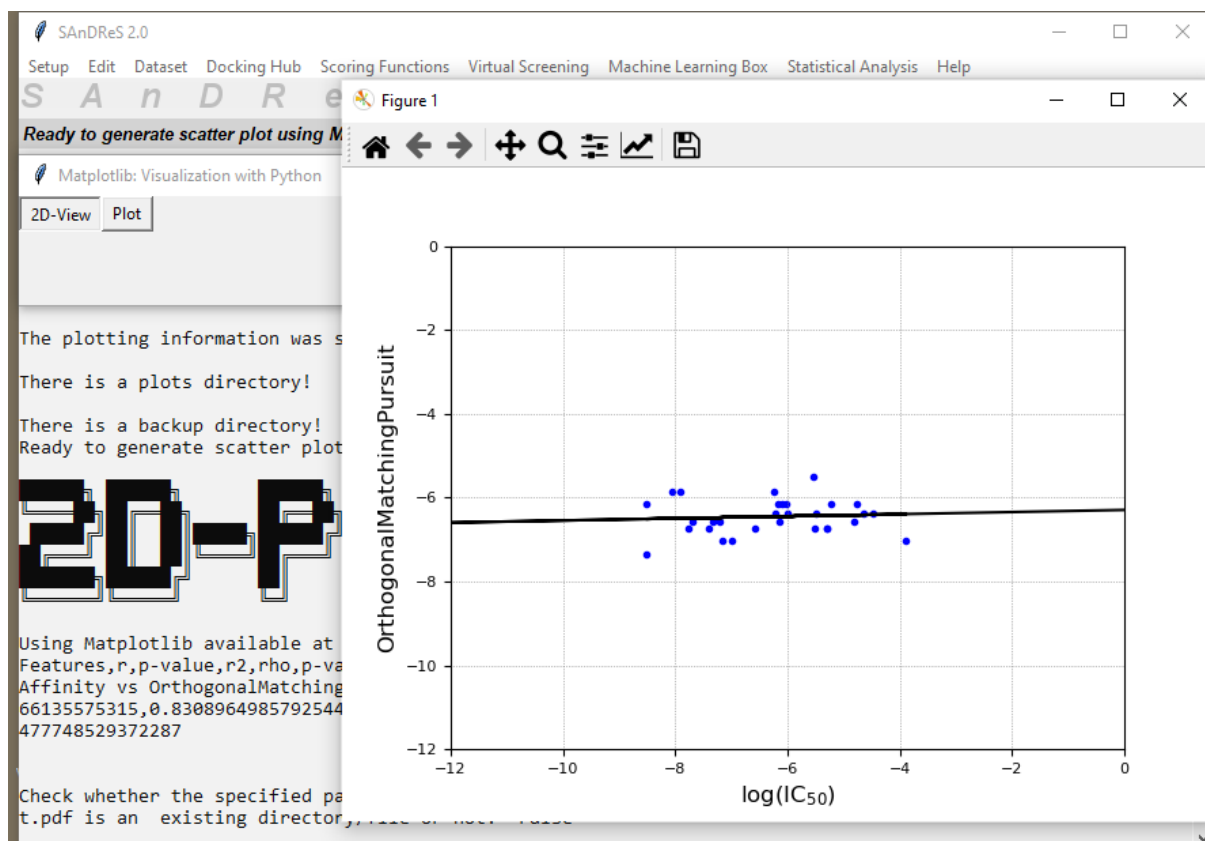


Figure 39. Scatter plot of OrthogonalMatchingPursuit vs $\log(\text{IC}_{50})$ for scores4xtal_test.csv file (Windows Version).

Click on the Close button. To assess the intermolecular contacts, click on *Statistical Analysis->2D Plot Intermolecular Interactions*.

We have the following pop-up window.

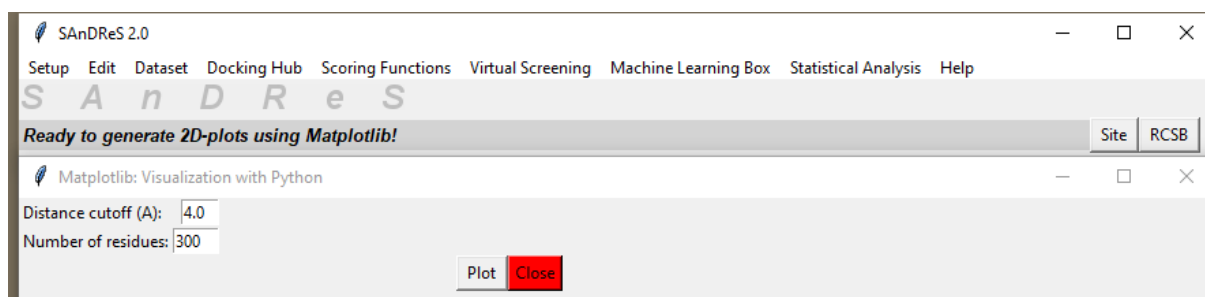


Figure 40. Matplotlib: Visualization with Python Menu (Windows Version).

Click on the Plot button. Click on the Yes button. It may take a few minutes to count all contacts. After finishing, we get the following message:

SAnDReS finished the "Add 2D Plot" request!

Click on the *Close* button. SAnDReS generated the following plots.

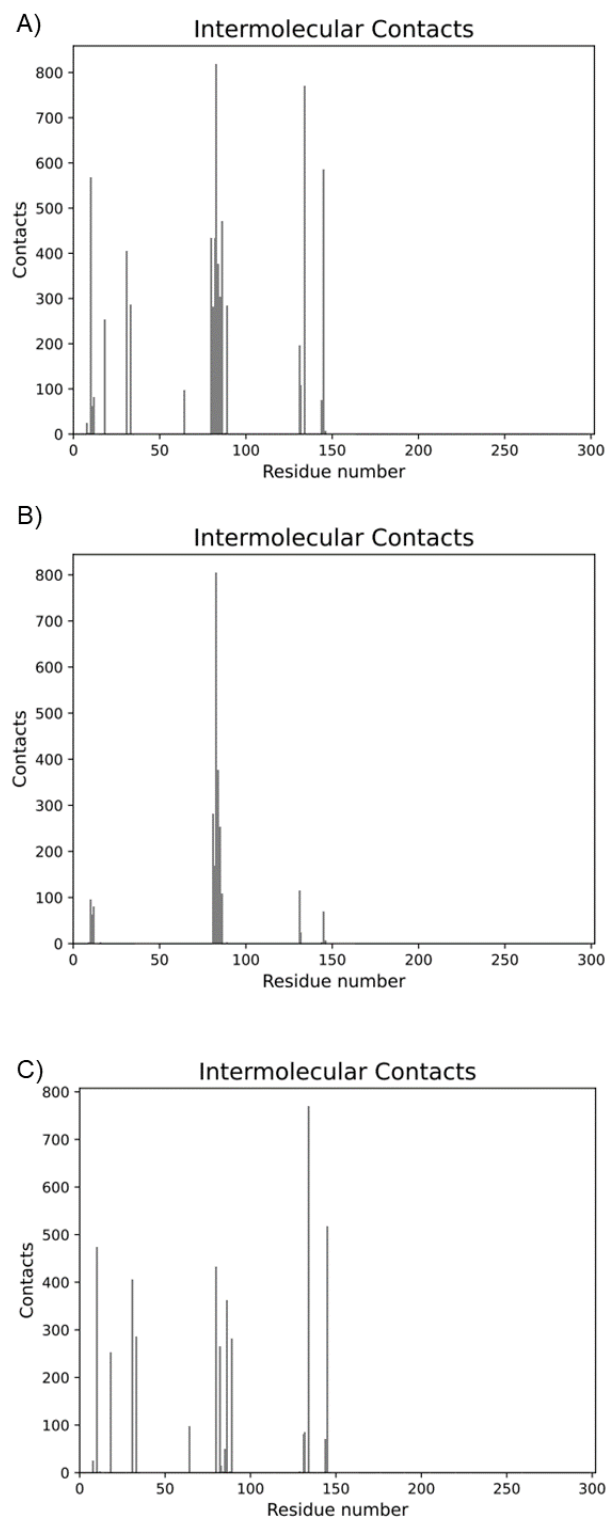


Figure 41. Intermolecular contacts taking all atoms in the protein (A), main-chain atoms (B), and side-chain atoms (C) (Windows Version).

To have a 3D-view of the intermolecular contacts, click on *Statistical Analysis->3D Plot Intermolecular Interactions*.

We have the following pop-up window.

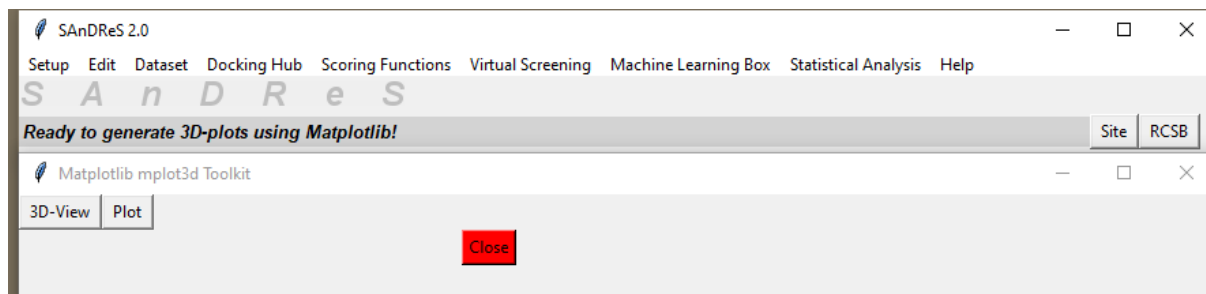


Figure 42. Matplotlib mplot3d Toolkit (Windows Version).

Click on the *3D-View* button. Click on the *Yes* button.

Now, we have an interactive 3D-view of the intermolecular contacts for all the structures in the dataset. You may left-click on the figure to rotate it.

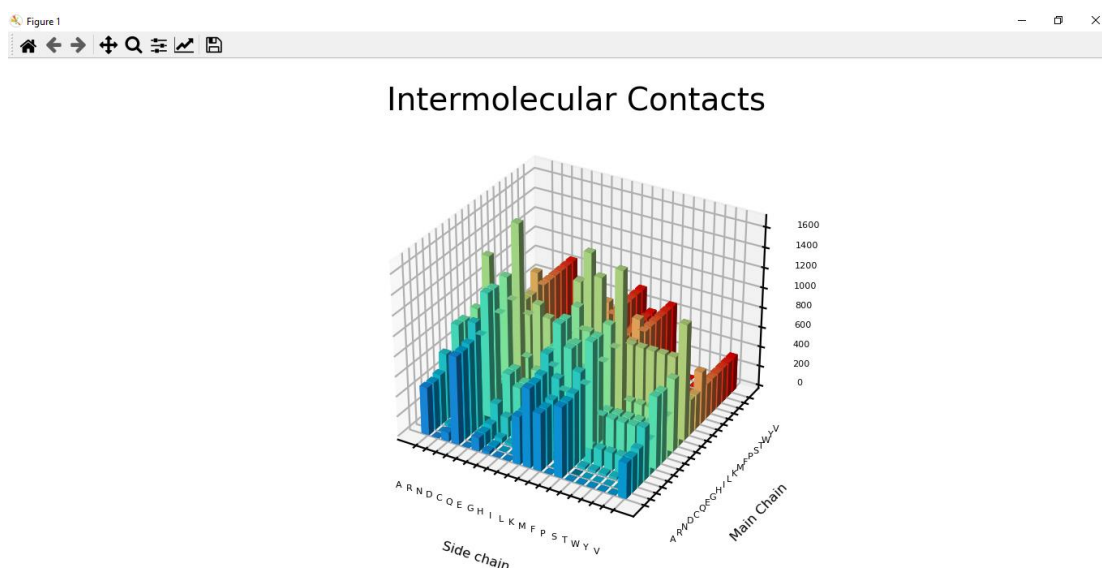


Figure 43. Interactive 3D-view of the intermolecular contacts (Windows Version).

To generate a plot, close the interactive window. Click on the *Plot* button and the *Yes* option. SAnDReS will generate the *bar_plot_3d.pdf* file in the *plots* directory. You get the following message:

SAnDReS finished the "Add 3D Plot" request!

Click on the *Close* button. We finished the statistical analysis. Now, let's save our results as a zipped folder. Click on *Setup->Project Directory->Backup Current Project*.

Click on the *Yes* option. It may take a few minutes. After backing up the current project directory, you get the following message:

*Successfully created a backup of the directory
C:/Users/Walter/sandres2_win/datasets/CDK2_IC50/*

You may delete or unzip this zipped folder using additional options in the Setup menu shown in Figure 44.

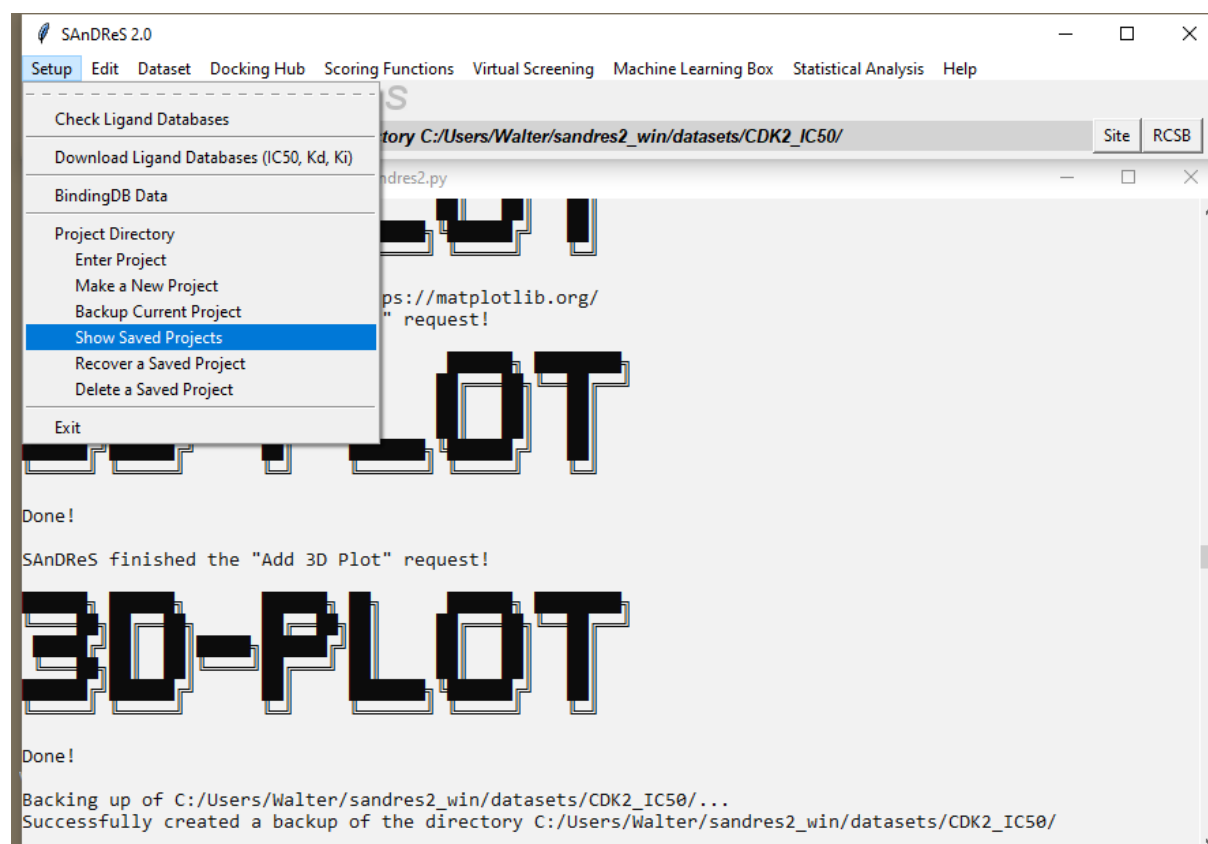


Figure 44. Setup menu (Windows Version).

To finish this session, click on *Setup->Exit*.

Click on the *Yes* option. We finished Tutorial 1.

6. Tutorial 2. Molegro Virtual Docker Data

Here, we will explore some new tools unseen in Tutorial 1. We will use virtual screening (VS) results generated with Molegro Virtual Docker (MVD) (Thomsen & Christensen, 2006) for a dataset with CDK2 inhibitors downloaded from the BindingDB (Gilson et al., 2016). We employed the PDB 2C5N without waters and ligands as our protein target in the VS simulation with the MVD. Our goal is to generate machine learning models with external VS results generated with MVD. We will not carry out docking simulations and calculations of using AutoDock Vina through the SAnDReS interface. We will convert the MVD data to SAnDReS format and use it to generate machine learning models. To run this tutorial, you will need the following files: *cdk2ki.sdf*, *cdk2ki.tsf*, and *vs2.csv*. We downloaded the first two files directly from the BindingDB on January 12, 2022. We carried out a search on the BindingDB for inhibitors of CDK2 with K_i data. We found a total of 131 ligands, but we have repeated entries. SAnDReS will filter this dataset to have unique ligands in our dataset. The *cdk2ki.tsf* file has the binding affinity data. SAnDReS uses this information to generate the *affinity_BindingDB_Ki.csv* file. This internal file furnishes binding data to generate the *bind_Ki.csv* file. We have the VS results generated with MVD for unique ligands in the *vs2.csv* file. You find all necessary files to run this tutorial in the *Ligands* folder of the *sandres2_win* directory (*sandres2* for Linux).

6.1. Tutorial 2. Setup

To learn how to start SAnDReS see [Tutorial 1](#). If you have already downloaded the ligand databases, you do not have to do that again. We must define a project directory different from the one used in Tutorial 1. Click on *Setup->Parameters->Enter Project*.

SAnDReS opens the *strdir.in* file. In the new pop-up window, add the following project directory: `C:/Users/Walter/sandres2_win/datasets/CDK2_MVD/`

We have the following view.

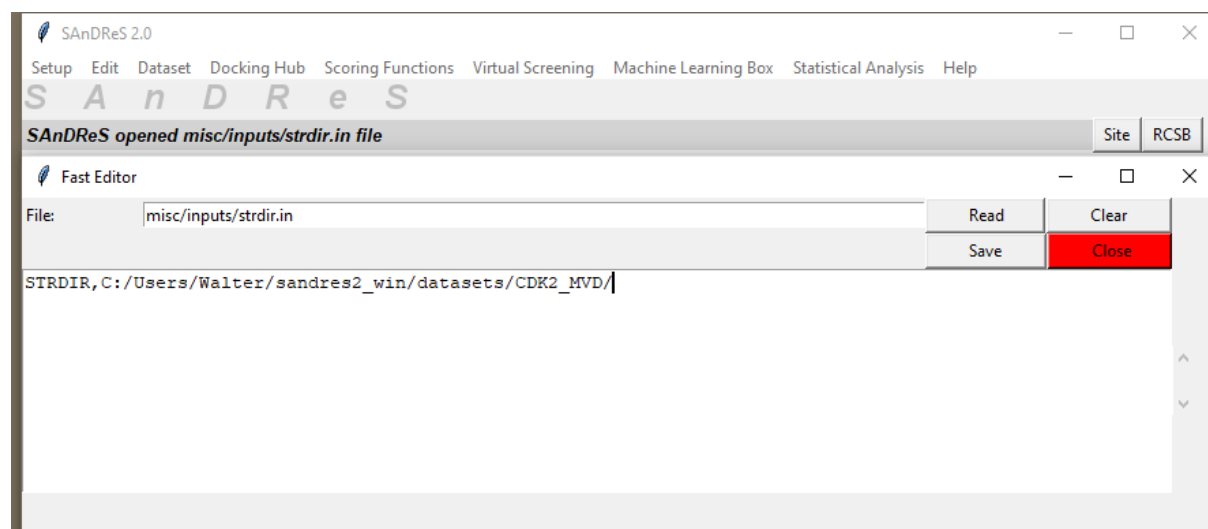


Figure 45. Fast Editor (Windows Version).

Click on the *Save* button and the *Close* button.

Click on *Setup->Project Directory->Make a New Project*.

Click on the *Yes* option.

SAnDReS creates a new directory and shows the following message:

Successfully created the directory
`C:/Users/Walter/sandres2_win/datasets/CDK2_MVD/`

It is necessary to copy to the project directory the following files: *cdk2ki.sdf*, *cdk2ki.tsf*, and *vs2.csv*. These files are in the *Ligands* folder. Now, we filter the ligands in the *cdk2ki.sdf* file and read binding data in the *cdk2ki.tsv* file.

Click on *Setup-BindingDB Data*.

Click on the *Yes* option.

SAnDReS shows the following pop-up window.

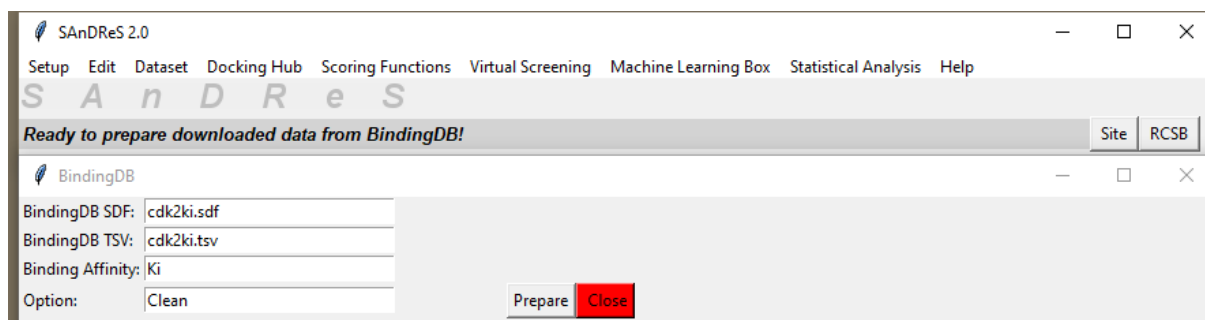


Figure 46. BindingDB Menu (Windows Version).

Both files (*cdk2ki.sdf* and *cdk2ki.tsv*) should be in the project directory. Click on the *Prepare* button.

SAnDReS will read binding affinity data and filter ligands then we get the following message:

Done! SAnDReS prepared BindingDB data for machine learning modeling!

Click on the *Close* button. SAnDReS generated two new files: *affinity_BindingDB_Ki.csv* and *cdk2ki_out.sdf*. The first file has the ligand IDs and binding affinity data; the last one has the atomic coordinates for all ligands after deleting the repeated ones.

We finished this part of Tutorial 2.

6.2. Tutorial 2. Machine Learning Box

We will use the VS results to generate a machine learning model. This information is in the *vs2.csv* file generated with MVD.

Click on *Machine Learning Box->For Modeling Using Molegro Results->Set up Molegro Parameters*.

SAnDReS opens the *molegro_par.csv* file. In this file, we have energy terms (*features_in*). Also, we have the scoring function used by MVD in the VS simulation. In this tutorial is the *PlantsScore*.

For this tutorial, leave the default values. Click on the *Close* button.

Click on *Machine Learning Box->For Modeling Using Molegro Results->Convert VS Data*.

Click on the *Yes* option.

SAnDReS shows the following pop-up window.

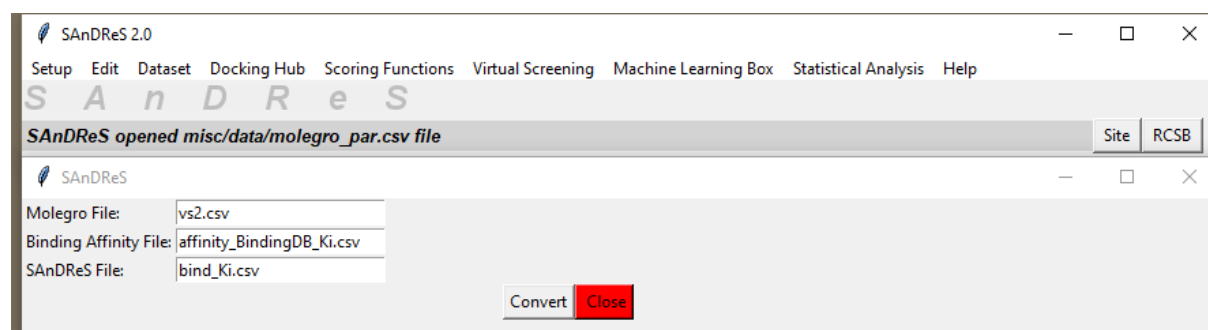


Figure 47. Menu to convert from MVD to SAnDReS format (Windows Version).

Click on the *Convert* button.

After converting, SAnDReS show the following message:

Done! Molegro data converted to SAnDReS format!

Click on the *Close* button.

SAnDReS generated *bind_Ki.csv* file.

Click on *Machine Learning Box->Enter Machine Learning Parameters*.

You should have MVD energy terms in the *features_in*, as shown in the following figure.

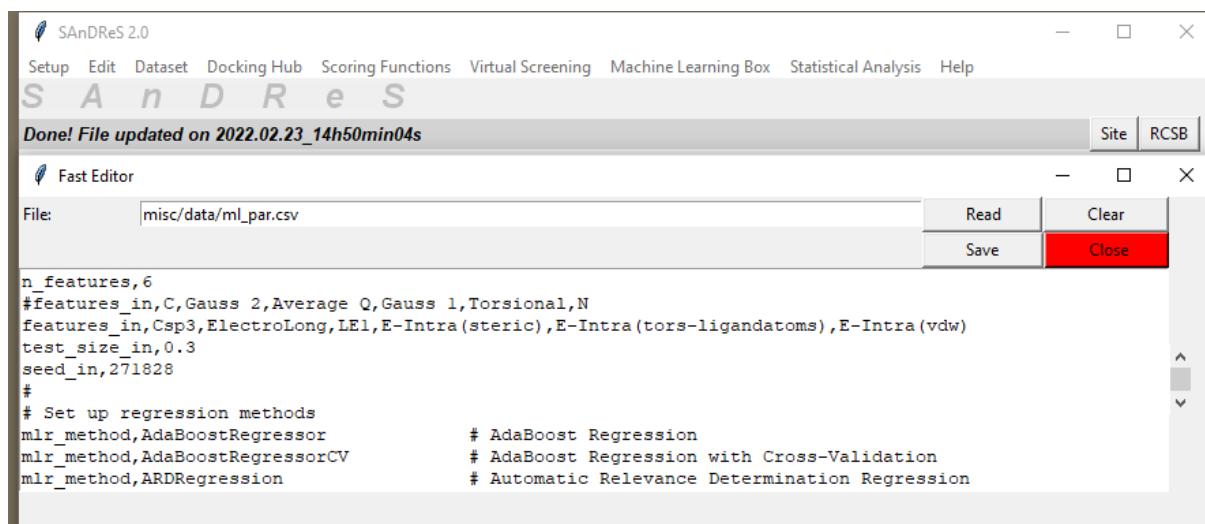


Figure 48. Fast Editor (Windows Version).

Click on the Save button and the Close button.

If you look at the Machine Learning Box, we will see that the sequence of commands is the same used in Tutorial 1 ([Section 5.6](#)). Please, follow the same steps to generate a machine learning model with MVD data.

Click on *Machine Learning Box->For Modeling Using Molegro Results->Preprocess*.

Click on *Machine Learning Box->For Modeling Using Molegro Results->Automatic Generation of PDBs for Training and Test Sets*.

Click on *Machine Learning Box->For Modeling Using Molegro Results->Generate Training and Test Sets*.

Click on *Machine Learning Box->For Modeling Using Molegro Results->Filter Data*.

Click on *Machine Learning Box->For Modeling Using Molegro Results->Statistical Analysis (Features)*.

Click on *Machine Learning Box->For Modeling Using Molegro Results->Regression Methods*.

Click on *Machine Learning Box->For Modeling Using Molegro Results->Apply Regression Model*.

Click on *Machine Learning Box->For Modeling Using Molegro Results->Bivariate Analysis of Regression Models*.

In the end, you will have new machine learning models based on the MVD energy terms. The following figure shows the `scores4xtal_test_stats_joblib_models.csv` file sorted by rho. The best model has rho = 0.6557 for ExtraTreesRegressorCV method. The correlation for the MVD scoring function (PlantsScore) is -0.297656. As we can see, our machine learning models present better overall predictive performance.

| Method | r | p-value(r) | r2 | rho | p-value(rho) | MSE | RMSE |
|-----------------------------|----------|------------|-----------|----------|--------------|---------|---------|
| ExtraTreesRegressorCV | 0.405254 | 0.061338 | 0.16423 | 0.655747 | 0.000921923 | 1.13176 | 1.06384 |
| GradientBoostingRegressorCV | 0.481403 | 0.0233062 | 0.231749 | 0.648404 | 0.00109928 | 1.05451 | 1.02689 |
| ExtraTreesRegressor | 0.346761 | 0.113872 | 0.120243 | 0.612256 | 0.00245638 | 1.2687 | 1.12637 |
| KNeighborsRegressorCV | 0.568208 | 0.0058 | 0.32286 | 0.594626 | 0.00351525 | 1.00939 | 1.00468 |
| KNeighborsRegressor | 0.359258 | 0.100577 | 0.129066 | 0.514689 | 0.0142481 | 1.17639 | 1.08462 |
| RandomForestRegressorCV | 0.278475 | 0.209507 | 0.0775481 | 0.506072 | 0.0162575 | 1.23831 | 1.11279 |
| XGBRegressor | 0.379586 | 0.0814408 | 0.144086 | 0.489794 | 0.0206774 | 1.20786 | 1.09903 |
| XGBRegressorCV | 0.315219 | 0.153022 | 0.0993631 | 0.429376 | 0.046129 | 1.2703 | 1.12708 |
| GradientBoostingRegressor | 0.271035 | 0.222443 | 0.0734598 | 0.408359 | 0.0591898 | 1.46492 | 1.21034 |
| AdaBoostRegressor | 0.449663 | 0.0357605 | 0.202197 | 0.40678 | 0.0602751 | 1.12965 | 1.06285 |
| BaggingRegressorCV | 0.284875 | 0.198789 | 0.0811537 | 0.379554 | 0.081469 | 1.21142 | 1.10065 |
| AdaBoostRegressorCV | 0.283421 | 0.20119 | 0.0803277 | 0.376378 | 0.0842643 | 1.22242 | 1.10563 |
| DecisionTreeRegressorCV | 0.12971 | 0.565076 | 0.0168246 | 0.375639 | 0.0849245 | 1.92835 | 1.38865 |
| RandomForestRegressor | 0.211695 | 0.344273 | 0.0448147 | 0.367693 | 0.0922725 | 1.34357 | 1.15913 |
| GaussianProcessRegressor | 0.195872 | 0.382336 | 0.0383658 | 0.366563 | 0.0933542 | 7.07942 | 2.66072 |
| TweedieRegressor | 0.296533 | 0.180229 | 0.0879319 | 0.334369 | 0.128273 | 1.20248 | 1.09658 |
| BayesianRidge | 0.290156 | 0.190228 | 0.0841905 | 0.331545 | 0.131729 | 1.20617 | 1.09826 |
| VotingRegressor | 0.239385 | 0.283283 | 0.0573053 | 0.324202 | 0.141027 | 1.31166 | 1.14528 |
| VotingRegressorCV | 0.247393 | 0.266996 | 0.0612032 | 0.321378 | 0.144723 | 1.28815 | 1.13497 |
| TheilSenRegressorCV | 0.240745 | 0.280474 | 0.0579583 | 0.298786 | 0.176785 | 1.27327 | 1.12839 |
| TweedieRegressorCV | 0.290362 | 0.189899 | 0.0843102 | 0.294267 | 0.18374 | 1.20341 | 1.097 |
| BayesianRidgeCV | 0.288693 | 0.192574 | 0.0833439 | 0.294267 | 0.18374 | 1.20553 | 1.09797 |

Figure 49. Partial view of the `scores4xtal_test_stats_joblib_models.csv` file.

Note. SAnDReS may show nan (not a number) for statistical analysis of some regression methods. It happens due to errors generated in the Scikit-Learn. If you have only a few methods for this problem, you may ignore them. Otherwise, you may have to change regression parameters. Please see [Appendix](#).

We finished the development of machine learning models of this tutorial.

6.3. Tutorial 2. Statistical Analysis

We may now generate a scatter plot for our best machine learning model.

Click on *Statistical Analysis->Scatter Plot->Set up Parameters*.

In the new pop-up window, insert in the field *Y-axis Label* the best model:
ExtraTreesRegressorCV

We have the following view.

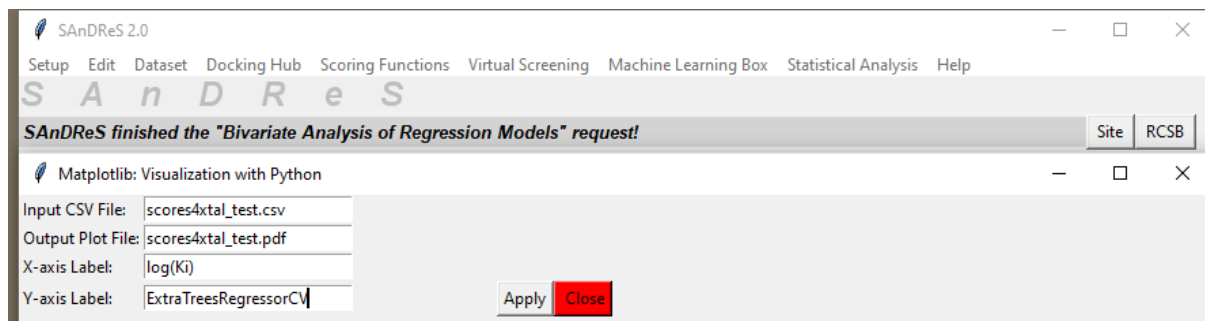


Figure 50. Matplotlib: Visualization with Python Menu (Windows Version).

Click on the *Apply* button.

SAnDReS generate the following message:

The plotting information was successfully updated!

Click on the *Close* button.

Click on *Statistical Analysis->Scatter Plot->Generate*.

We have the following pop-up window.

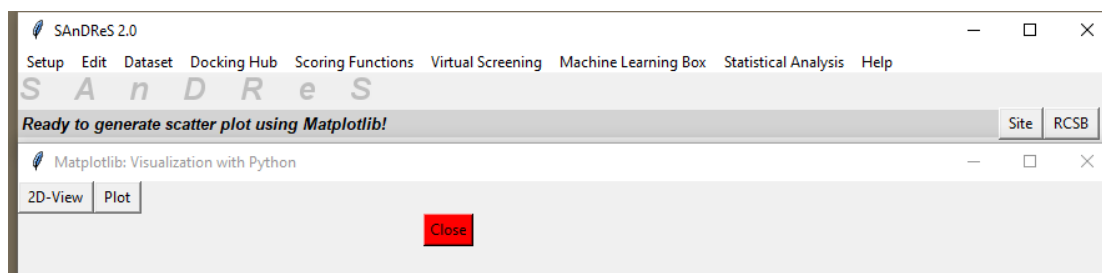


Figure 51. Matplotlib: Visualization with Python Menu (Windows Version).

Click on the *2D-View* button. SAnDReS generates the following plot.

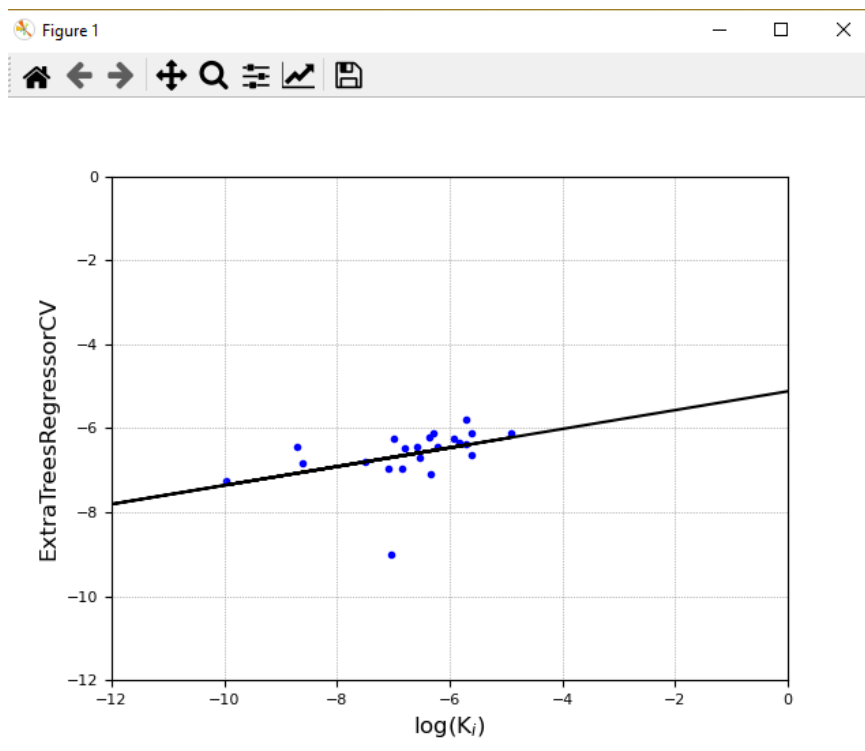


Figure 52. Scatter plot of *ExtraTreesRegressorCV* vs $\log(K_i)$ for *scores4xtal_test.csv* file (Windows Version).

To finish Tutorial 2, we save our results as a zipped folder. Click on *Setup->Project Directory->Backup Current Project*.

Click on the Yes option. After finishing the backup, SAnDReS shows the following message:

*Successfully created a backup of the directory
C:/Users/Walter/sandres2_win/datasets/CDK2_MVD/*

Click on *Setup->Exit*.

Choose the Yes option.

We finished Tutorial 2.

7. Tutorial 3. Docking of One Structure

Here, we will show how to use SAnDReS to run docking simulations for one crystallographic structure (redocking). We developed SAnDReS to explore the Scoring Function Space, which demands as many structures as possible. For instance, in [Tutorial 1](#), we worked with 102 crystal structures. Nevertheless, we can enjoy the SAnDReS integration and run docking simulations for one PDB.

We will perform docking simulations for the human purine nucleoside phosphorylase (PNP) complexed with immucillin-H (PDB access code: 1PF7) (de Azevedo et al., 2003). It is possible to highlight one additional aspect of SAnDReS here. For PNP, the asymmetric unit is a monomer, and the biological unit is a trimer with the inhibitor (immucillin-H) bound at the interface between monomers (de Azevedo et al., 2003). So if you download the PDB without investigating these differences, you might end up carrying out docking simulations for a structure not adequate for analysis of intermolecular interactions (the monomer structure). SAnDReS automatically carries out this analysis and generates the biological unit if necessary. Docking simulations focused on the PNP monomer would miss important intermolecular contacts found in the trimeric structure.

7.1. Tutorial 3. Setup

We consider that you have already installed SAnDReS and downloaded the ligand databases. We show the details about starting SAnDReS and downloading ligand databases in [Tutorial 1](#).

Now, we must define a new project directory.

Click on *Setup->Project Directory->Enter Project*.

SAnDReS opened the *strdir.in* file. Add the following line to this file.

```
STRDIR,C:/Users/Walter/sandres2_win/datasets/HsPNP/
```

The following figure shows the editing.

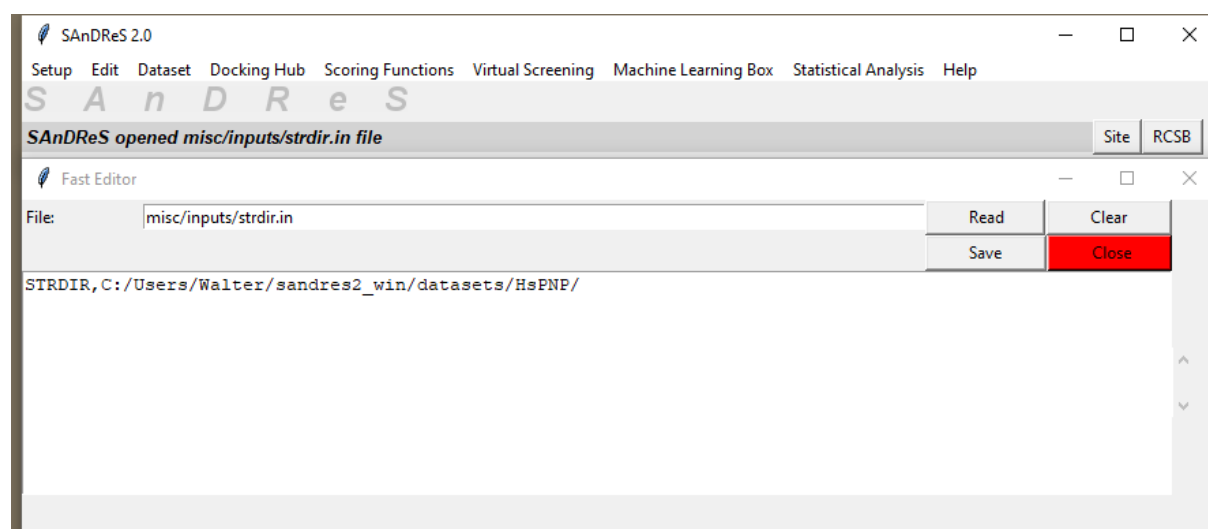


Figure 53. Fast Editor (Windows Version).

Click on the Save button and the Close button.

Click on *Setup->Project Directory->Make a New Project*.

Click on the Yes option. SAnDReS generates the new folder and shows the following message:

```
Successfully created the directory C:/Users/Walter/sandres2_win/datasets/HsPNP/
```

We finished the Setup part of this tutorial.

7.2. Tutorial 3. Dataset

Now, we will add the PDB access code.

Click on *Dataset->Enter PDB Access Codes*.

In the new pop-up window, we add the PDB access code: *1PF7*

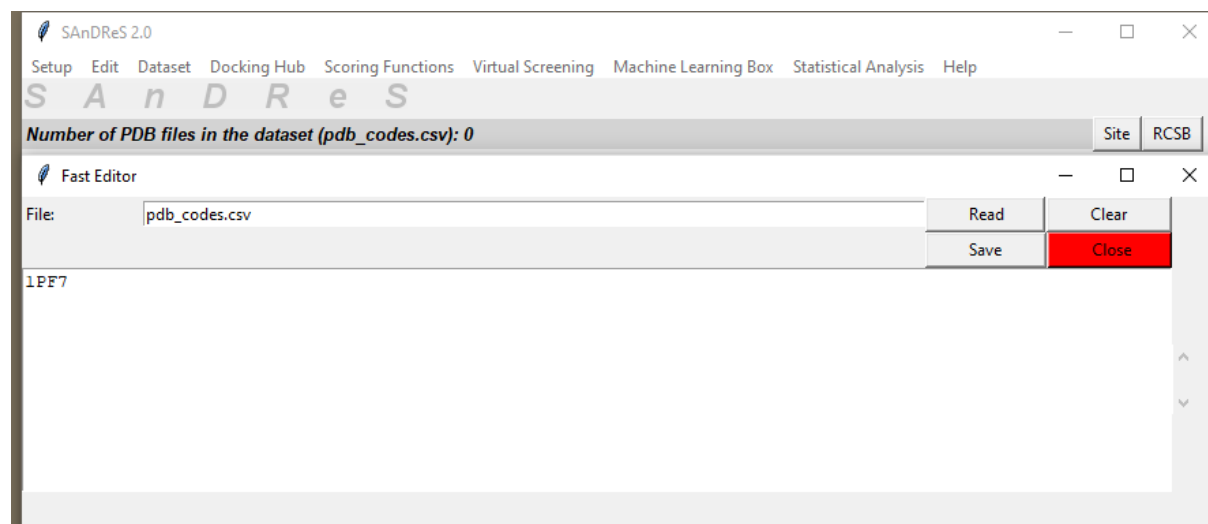


Figure 54. Fast Editor (Windows Version).

Click on the *Save* button and the *Close* button.

Click on *Dataset->Add->Binding Affinity Data*.

Click on the *Ki* button and the *Start* button. After finishing, we have the following message:

SAnDReS finished the "Add Binding Affinity Data" request!

Click on the *Done* button and the *Close* button.

Click on *Dataset->Add-Structures (PDB)*.

Click on the *with Ki data* button and the *Start* button.

SAnDReS will download the PDB file. After downloading it, we have the following message:

SAnDReS finished the "Add Structures (PDB)" request!

Click on the *Done* button and the *Close* button.

Click on *Dataset->Add-Structures (PDBQT)*.

Click on the *Yes* option.

We have the following pop-up window.

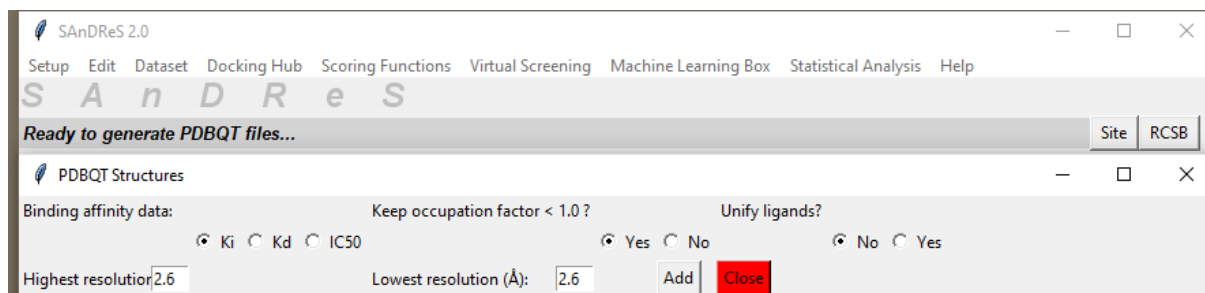


Figure 55. Menu to generate PDBQT structures (Windows Version).

Click on the *Add* button and the *Start* button.

SAnDReS generates a PDBQT file and shows the following message:

SAnDReS finished the "Add Structures (PDBQT)" request!

Click on the *Done* button and the *Close* button.

Click on *Dataset->Check Directories for Current Dataset*.

We get the following message:

Done! There are no missing PDBQT directories!

Click on *Dataset->Check Missing PDBQT files*.

We have the following message:

Done! No missing PDBQT files in the dataset!

It is not necessary to update the dataset. We finished the dataset part of this tutorial.

7.3. Tutorial 3. Docking Hub

In this part, we perform docking simulations.

Click on *Docking Hub*->*Enter Vina Parameters*.

We may edit the *vina_par.csv*. For this tutorial, leave it as it is. Click on the *Close* button.

Click on *Docking Hub*->*One-Click Docking Simulation with Vina (Centering: CM, EC, GC)*->*Run Simulation*.

Click on the *Yes* option and the *Run* button. Click on the *Start* button.

After finishing docking simulations, we have the following message:

SAnDReS finished the "One Click Docking with Vina" request using AutoDock Vina!

Click on the *Done* button and the *Close* button.

Note. We should not perform a statistical analysis of the docking results since we have only one structure. This statistical analysis determines docking accuracy for an ensemble of structures.

Your results are on *C:\Users\Walter\sandres2_win\datasets\HsPNP\pdbqt\1PF7* folder as shown in Figure 56.



















| | |
|---|------------------|
|  1PF7 | 23/02/2022 15:49 |
|  biomatrix | 23/02/2022 15:49 |
|  config_electric | 23/02/2022 15:52 |
|  config_geometric | 23/02/2022 15:52 |
|  config_mass | 23/02/2022 15:52 |
|  lig | 23/02/2022 15:49 |
|  lig.pdbqt | 23/02/2022 15:49 |
|  lig_out_electric.pdbqt | 23/02/2022 15:53 |
|  lig_out_geometric.pdbqt | 23/02/2022 15:52 |
|  lig_out_mass.pdbqt | 23/02/2022 15:52 |
|  receptor | 23/02/2022 15:49 |
|  receptor.pdbqt | 23/02/2022 15:49 |
|  vina_results_electric | 23/02/2022 15:53 |
|  vina_results_geometric | 23/02/2022 15:52 |
|  vina_results_mass | 23/02/2022 15:52 |
|  vina_simulation_electric | 23/02/2022 15:53 |
|  vina_simulation_geometric | 23/02/2022 15:52 |
|  vina_simulation_mass | 23/02/2022 15:52 |

Figure 56. Files found in the directory *C:\Users\Walter\sandres2_win\datasets\HsPNP\pdbqt\1PF7* (Windows Version).

SAnDReS wrote the RMSD values in the following files: *vina_results_electric.csv*, *vina_results_geometric.csv*, and *vina_results_electric.csv*.

The geometric centering generated the lowest RMSD (0.8129 Å) against 0.879 Å (electric) and 3.096 Å (mass).

Now, we save our results as a zipped folder. Click on *Setup->Project Directory->Backup Current Project*.

Click on the Yes option. After finishing the backup of the current project, SAnDReS shows the following message:

*Successfully created a backup of the directory
C:/Users/Walter/sandres2_win/datasets/HsPNP/*

Click on *Setup->Exit*.

Click on the Yes option.

We finished this tutorial.

8. Tutorial 4. Cyclin-Dependent Kinase 2 with K_i Data

In [Tutorial 1](#), we generated a machine learning model to predict $\log(IC_{50})$. Here, we focus on CDK2 with inhibition constant (K_i) data. Our goal is to develop a machine learning model to predict $\log(K_i)$.

8.1. Tutorial 4. Setup

We begin this part considering that you started SAnDReS and downloaded the ligand databases. Any doubts related to these tasks, please see [Tutorial 1](#).

Let's define the project directory. Click on *Setup->Project Directory->Enter Project*.

In the editor, enter the following string:

```
C:/Users/Walter/sandres2_win/datasets/CDK2_Ki/
```

Click on the *Save* button and the *Close* button.

Click on *Setup->Project Directory->Make a New Project*.

Click on the *Yes* option. SAnDReS created a new project directory. We have the following message:

Successfully created the directory C:/Users/Walter/sandres2_win/datasets/CDK2_Ki/

Now, on the main menu, click on the *RCSB* button. SAnDReS shows the following pop-up window.

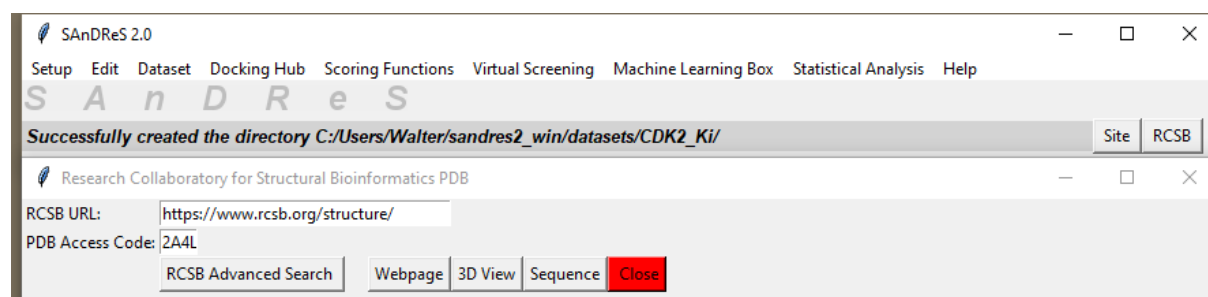


Figure 57. RCSB Menu (Windows Version).

In this menu, we insert *2A4L* (CDK2 in complex with roscovitine) (de Azevedo et al., 1997). Click on the *Sequence* button to start Google Chrome at the PDB site with the sequence for the structure 2A4L. If you have an error message here, please check the installation of Google Chrome on your computer. If everything is fine, you can Ctrl C the amino acid sequence of the CDK2 structure. Back to the RCSB menu, click on the *RCSB Advanced Search* button. SAnDReS will open the PDB advanced search tool with Google Chrome. Set up the parameters as shown in the following figure. Do not forget to Ctrl C the CDK2 sequence to the field sequence.

Advanced Search Query Builder

Full Text

Structure Attribute

Binding Affinity Value - Ki x exists nM + NOT Count

Add Attribute Add Subquery Remove Subquery

Add Subquery

Chemical Attribute

Sequence

AND KLADFGGLARAFGVPRITYTHEVVTWYRAPEILLGCKYYSTAVDIWSLGCIFAEMVTRRALFPGDSEIDQLFRIFRTLGTPEVVWPGVTSMPDYKPSFKWARQDFSKVVPPLDEEDGRSLLSQMLHYDPNKRISA
KAALAHFFQDVTKPVPHRL

PDB ID 1MBN Target Protein E-Value Cutoff 0.001 Identity Cutoff 98 % (Integer only) Count Clear

Sequence Motif

Structure Similarity

Structure Motif

Chemical

Return Structures grouped by No Grouping Count Clear

Figure 58. Advanced Search Menu (Windows Version).

Click on the magnifying glass symbol on the right and click on the *OK* option. Click on the download icon on the bottom right and Ctrl C the PDB access codes. You may download the *pdb_codes.csv* file by clicking [here](#). You have to save this file in the project directory. We have the following PDB access codes for this tutorial:

1JSV, 1H1S, 4ACM, 1PF8, 1H1P, 1PYE, 2FVD, 2EXM, 2CLX, 2XMY, 2C6O, 1E1V, 4BCM, 4BCK, 4BCP, 1E1X, 1PXM, 1PXN, 1PXO, 4BCO, 1PXP, 4BCN, 4EOS, 4EOR, 4EOP, 4EOI, 4EON, 4EOL, 4EOK

Close Google Chrome. Click on the *Close* button of the RCSB menu.

We finished this part of Tutorial 4.

8.2. Tutorial 4. Dataset

Click on *Edit->Parameters->Project Summary*.

In the editor, add the following sentence:

CDK2 with Ki data

Click on the *Save* button and the *Close* button.

Click on *Dataset->Enter PDB Access Codes*.

Ctrl C the PDB access codes. Click on the *Save* button and the *Close* button.

Click on *Dataset->Unify PDB Access Codes*.

SAnDReS will show the following message:

Done! Number of repeated PDB files: 0.

Click on *Dataset->Add->Binding Affinity Data*.

On the new pop-up window, click on the *Ki* button and the *Start* button.

After finishing adding K_i data, SAnDReS shows the following message:

SAnDReS finished the "Add Binding Affinity Data" request!

Click on the *Done* button and the *Close* button.

Click on *Dataset->Add->Structures (PDB)*.

On the new pop-up window, click on the *Ki* button and the *Start* button.

After completing the downloading, we get the following message:

SAnDReS finished the "Add Structures (PDB)" request!

Click on the *Done* button and the *Close* button.

Click on *Dataset->Add->Structures (PDBQT)*.

Click on the *Yes* option. Click on the *Add* button and the *Start* button.

Once finished, we have the following message:

SAnDReS finished the "Add Structures (PDBQT)" request!

Click on the *Done* button and the *Close* button.

Click on *Dataset->Check Directories for Current Dataset*.

SAnDReS shows the following message:

Done! There are no missing PDBQT directories!

Click on *Dataset->Check Missing PDBQT files*.

We get the following message:

Done! No missing PDBQT files in the dataset!

We do not need to update the dataset. We finished the dataset part of this tutorial.

8.3. Tutorial 4. Docking Hub

Now, we perform docking simulations.

Click on *Docking Hub*->*One-Click Docking with Vina (Centering: CM, EC, GC)*->*Run Simulation*.

Follow the same steps shown in Tutorial 1 ([Section 5.3](#)). Once finished, we have the best RMSD for each structure in the *vina_rmsd_results.csv* file. SAnDReS deleted one PDB from the dataset due to an unsuccessful docking simulation (PDB access code: 4BCP). Now, our dataset has 28 structures. Below, we have part of the summary of the docking results (*docking_summary.txt*).

Docking accuracies

DA1 = 0.286

DA2 = 0.295

Docking RMSD

Mean RMSD (A) = 4.666

Maximum RMSD (A) = 8.366

Minimum RMSD (A) = 1.034

No docking results for the following PDB(s): ['4BCP']

We finished this part of Tutorial 4.

8.4. Tutorial 4. Scoring Functions

Now we calculate the energy terms for the crystal structures and docked poses. We follow the same procedure described in Tutorial 1 ([Section 5.4](#)), only keeping in mind that the binding affinity for this tutorial is K_i . Here, we generate the *scores4poses.csv* file with the energy terms for all poses. We have the energy terms added to the *bind_Ki.csv* file for the crystal structures.

We finished the calculation of the energy terms. We will use this data to generate machine learning models.

8.5. Tutorial 4. Machine Learning Box

Here, we generate machine learning models based on the energy terms calculated for the crystal structures. Since we have only 28 PDBs in the dataset, we need extra care in the machine learning modeling. Some authors define five as the minimum number of observations for each feature in the modeling (Gramatica, 2013). Considering that we take 70 % of the dataset as the training set, we have 20 structures (observations). But for cross-validation, we take 80 % of the training set. We ended up with 16 entries. So, we can barely build a model with three features.

We repeat the steps described in Tutorial 1 for the machine learning part ([Section 5.6](#)). We take the following features: O, Torsional, Gauss 1, Gauss 2.

After finishing the generation of the machine learning models, we may check the predictive performance by analyzing the *scores4xtal_test_stats_joblib_models.csv* file. Figure 59 shows part of the results.

| Method | r | p-value(r) | r2 | rho | p-value(rho) | MSE | RMSE |
|------------------------------|----------|------------|-----------|----------|--------------|---------|---------|
| GaussianProcessRegressorCV | 0.419643 | 0.300663 | 0.1761 | 0.512348 | 0.194223 | 2.19599 | 1.48189 |
| GaussianProcessRegressor | 0.332622 | 0.420807 | 0.110638 | 0.439155 | 0.276327 | 3.52633 | 1.87785 |
| DecisionTreeRegressorCV | 0.225384 | 0.591498 | 0.0507979 | 0.273297 | 0.512512 | 2.48651 | 1.57687 |
| ExtraTreesRegressorCV | 0.224615 | 0.592798 | 0.0504519 | 0.414758 | 0.306912 | 2.44803 | 1.56462 |
| MLPRegressorCV | 0.210381 | 0.61702 | 0.0442603 | 0.317168 | 0.443989 | 3.48123 | 1.8658 |
| RANSACRegressor | 0.196385 | 0.641136 | 0.038567 | 0.365963 | 0.372625 | 3.06544 | 1.75084 |
| PassiveAggressiveRegressor | 0.194978 | 0.643576 | 0.0380165 | 0.365963 | 0.372625 | 1.65285 | 1.28563 |
| SVRCV | 0.192942 | 0.647111 | 0.0372267 | 0.365963 | 0.372625 | 2.06293 | 1.43629 |
| MLPRegressor | 0.1926 | 0.647706 | 0.0370949 | 0.121988 | 0.773532 | 2.11687 | 1.45495 |
| SVR | 0.192141 | 0.648505 | 0.036918 | 0.365963 | 0.372625 | 1.80235 | 1.34251 |
| LinearSVRCV | 0.191566 | 0.649504 | 0.0366977 | 0.463553 | 0.247323 | 1.68465 | 1.29794 |
| NuSVRCV | 0.191518 | 0.649588 | 0.0366792 | 0.463553 | 0.247323 | 1.94074 | 1.3931 |
| RANSACRegressorCV | 0.191485 | 0.649645 | 0.0366665 | 0.317168 | 0.443989 | 2.70793 | 1.64558 |
| TheilSenRegressor | 0.190028 | 0.652182 | 0.0361107 | 0.463553 | 0.247323 | 1.54587 | 1.24333 |
| TweedieRegressorCV | 0.189321 | 0.653414 | 0.0358425 | 0.463553 | 0.247323 | 1.4201 | 1.19168 |
| PassiveAggressiveRegressorCV | 0.189212 | 0.653604 | 0.0358012 | 0.317168 | 0.443989 | 2.26334 | 1.50444 |
| BayesianRidgeCV | 0.18763 | 0.656364 | 0.035205 | 0.463553 | 0.247323 | 1.44061 | 1.20025 |
| LarsCV | 0.186977 | 0.657504 | 0.0349603 | 0.463553 | 0.247323 | 1.5377 | 1.24004 |
| TweedieRegressor | 0.186895 | 0.657647 | 0.0349297 | 0.365963 | 0.372625 | 1.30757 | 1.14349 |
| ExtraTreeRegressor | 0.186558 | 0.658235 | 0.0348039 | 0.253185 | 0.545176 | 2.99684 | 1.73114 |
| LinearSVR | 0.186361 | 0.65858 | 0.0347304 | 0.463553 | 0.247323 | 1.68322 | 1.29739 |
| TheilSenRegressorCV | 0.184789 | 0.661328 | 0.0341469 | 0.365963 | 0.372625 | 2.50186 | 1.58173 |
| BayesianRidge | 0.184484 | 0.66186 | 0.0340344 | 0.365963 | 0.372625 | 1.31805 | 1.14806 |

Figure 59. Partial view of the *scores4xtal_test_stats_joblib_models.csv* file.

Here, our machine learning models generalize better than the AutoDock Vina scoring function (Affinity), which has an $r^2 = 0.00909986$. Figure 60 shows the scattering plot of the GaussianProcessRegressorCV model against the experimental $\log(K_i)$ for the test set ($r^2 = 0.1761$). To check how to create the scatter plots, see Tutorial 1 ([Section 5.9](#)).

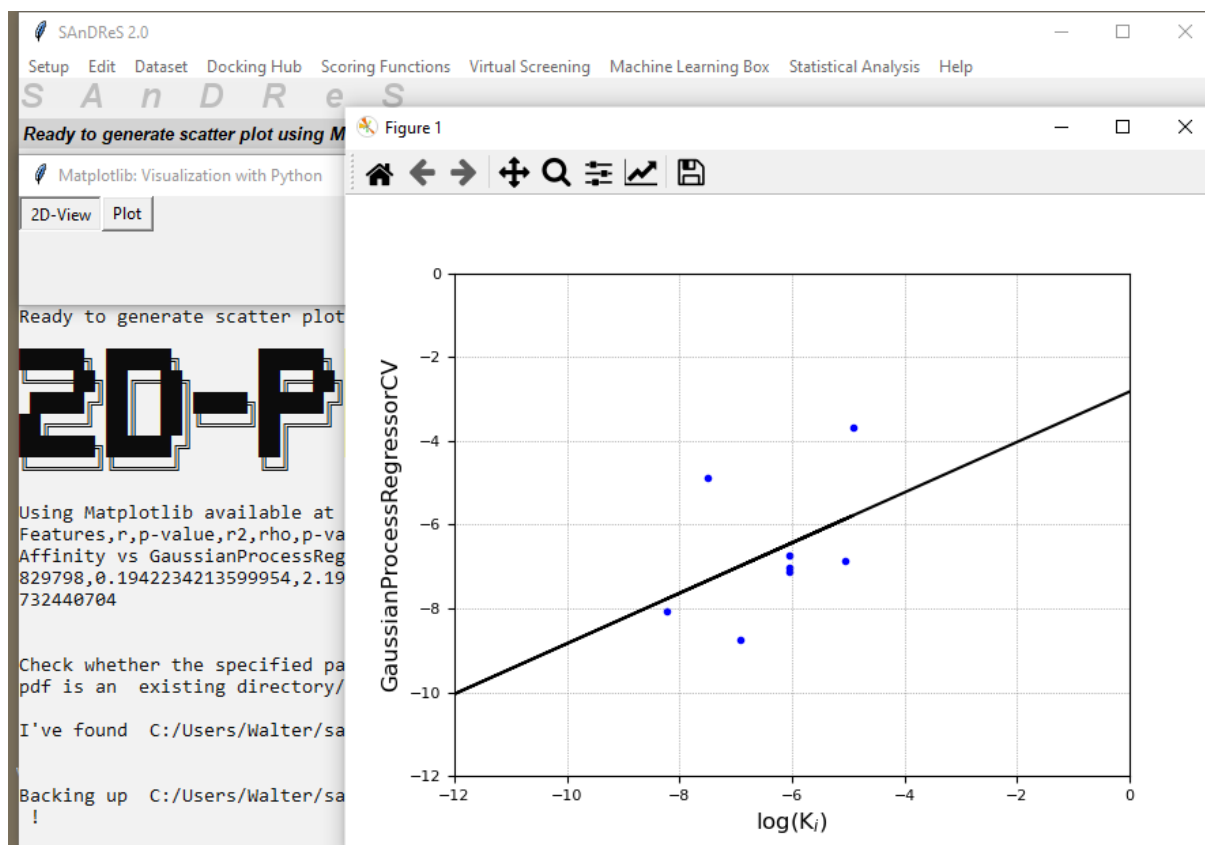


Figure 60. Scatter plot of GaussianProcessRegressorCV vs $\log(K_i)$ for scores4xtal_test.csv file (Windows Version).

We may apply our regression model to sort docking results, as shown in Tutorial 1 (Section 5.7). We will not carry out this task here.

Finally, we will save these results as a zipped folder. Click on *Setup->Project Directory->Backup Current Project*.

Click on the Yes option. After finishing the backup of the current project, SAnDReS shows the following message:

*Successfully created a backup of the directory
C:/Users/Walter/sandres2_win/datasets/CDK2_Ki/*

Click on *Setup->Exit*.

Click on the Yes option.

We finished Tutorial 4. Any question related to this tutorial, or any other material related to SAnDReS feel free to email us <mailto:sandres@azevedolab.net>.

9. Scoring Function Space

We envisage protein-ligand interaction as a result of the relation between the protein space (Smith, 1970) and the chemical space (Bohacek et al., 1996), and we propose to approach these sets as a unique complex system, where the application of computational methodologies could contribute to understand the structural basis for the specificity of ligands for proteins. Such approaches have the potential to create novel scoring functions to predict binding affinity with superior predictive power when compared with standard methodologies. We propose to use the abstraction of a mathematical space composed of infinite computational models to predict ligand-binding affinity, named here as scoring function space (Heck et al., 2017; Bitencourt-Ferreira & de Azevedo Jr., 2019).

Using supervised machine learning techniques, we can explore this scoring function space to build a computational model targeted to a specific biological system. SAnDReS explores the SFS using 64 machine learning methods, as highlighted in Tutorials [1](#), [2](#), and [4](#).

10. FAQ. Frequently Asked Questions

1) Why does SAnDReS only use crystallographic structures in the ligand databases?

Answer. A recent review (Veit-Acosta & de Azevedo, 2021) showed that the X-ray diffraction crystallography contributed to over 99 % of the protein structures for which K_i , K_d , and IC_{50} data is available. That is the main reason SAnDReS focused exclusively on crystallographic structures.

2) Can I use SAnDReS to analyze RNA/DNA datasets?

Answer. No, the current version of SAnDReS does not handle nucleic acids.

3) Can I include peptides (or polysaccharides) in the ligand database?

Answer. SAnDReS does not handle peptides (or polysaccharides) as ligands. The major problem is how SAnDReS identifies active ligands (with binding data), with one set of three characters (e.g., RRC) and only one number for each ligand. If you have a peptide (or polysaccharide) satisfying these limitations, it should work. Otherwise, the answer is no.

4) Can I use AutoDock4 or Vinardo scoring functions in SAnDReS?

Answer. No. AutoDock Vina 1.2.3 handles these scoring functions, but the current version of SAnDReS considers only the AutoDock Vina energy terms (Gauss 1, Gauss 2, Repulsion, Hydrophobic, Hydrogen, and Torsional). Considering the built-in descriptors and AutoDock Vina energy terms combined with 64 regression methods, we have a wide range of possibilities of generating scoring functions with SAnDReS. Nevertheless, it is always possible to add external energy terms editing the scoring function files.

5) I generated PDBQT files for my dataset using *Dataset->Add->Structures (PDBQT) using AutoDockTools4*. I got several messages reporting a not successful generation of PDBQT files. Is that a bug?

Answer. No. SAnDReS integrates several packages such as MGLTools, AutoDock Vina, Scikit-Learn, and with such huge integration, we lost some specificity of each of the packages that SAnDReS handles. Most notably for MGLTools 1.5.7 (AutoDockTools). It is always possible to generate the PDBQT files outside SAnDReS. If you keep the same file names and directories, SAnDReS will recognize them all.

6) Structures are missing in the ligand databases. Can I add them?

Answer: Sure. We intend to update the ligand databases every year. And I will be glad to receive an email with any addition you have (<mailto:sandres@azevedolab.net>). You may edit the CSV files (*bind_IC50.csv*, *bind_Kd.csv*, and *bind_Ki.csv*) on the *sandres2_win\misc\data* directory. It is also necessary to add ligand PDBQT files in the *sandres2_win\misc\data\pdbqt* directory.

Alternatively, you may use a built-in function of SAnDReS to add ligands to the databases. It is necessary to have a CSV file with the data you want to add. You may find a template (*add_ligand_data.csv*) in the *Ligands* folder. You also need the structures in the *pdbqt* folder. For each structure in the CSV file, you need a *lig.pdbqt* file in the *pdbqt* folder. To add the new ligand data, do as follows:

a) Click on *Edit->Parameters->Ligands to Add*.

b) In the new pop-up window, enter the type of binding affinity to add (*bind_in*), the ligand directory (*ligand_dir*), and the CSV file (*ligand_data*). Click on the *Save* button and *Close* button.

c) Click on *Edit->Manually Add Data to Ligand Databases*. Click on the *Yes* option. If everything goes fine, you will get the following message: *Done! Successfully updated Ligand Database!*

For the template (*add_ligand_data.csv*), we added data for structure 4I3K. You may check the *bind_Ki.csv* file of the Ligand Database. Click on *Edit->Binding Affinity Ligand Data (Ki)*. You will see in the first line after the header: *4I3K, 1BX, A, 505, 3.31, 1.0, 2.8e-07, -6.552841968657781, 6.552841968657781*

7) My machine learning model does not generalize well. How can I improve it?

Answer. There is no easy answer to this question. The main strength of SAnDReS is the freedom to generate a machine learning model adequate to your protein system. If the model is not what you expected, try to play with the features. Try alternative sets of features that might generate a model that generalizes better. Keep trying. With SAnDReS, you have the liberty to explore different regions of the SFS.

8) Can I use SAnDReS to generate machine learning models using energy terms or descriptors from other sources?

Answer. Certainly, yes. Please see [Tutorial 2](#). If you have your data as a CSV file, you can easily deceive SAnDReS, and it will take your dataset to generate machine learning models.

9) How do I cite SAnDReS?

Answer. Thank you very much for your question. SAnDReS 2.0.0 became operational on January 12, 2022. And we are about to submit it (February 25, 2022).

If you are in a hurry, please cite the paper describing SAnDReS 1.0. It is the following:

Xavier MM, Heck GS, Avila MB, Levin NMB, Pintro VO, Carvalho NL, Azevedo WF. SAnDReS a Computational Tool for Statistical Analysis of Docking Results and Development of Scoring Functions. Comb Chem High Throughput Screen. 2016;19(10):801-812. doi: 10.2174/1386207319666160927111347. PMID: 27686428.

10) Can I use part of the SAnDReS code in my program?

Answer. Yes. But, you need to give credit to those who developed SAnDReS. See question 9. SAnDReS is free software (GNU General Public License v3.0).

11) Can I run SAnDReS on a Mac OS?

Answer. There are some issues related to the Mac OS. The most critical one is related to MGLTools. MGLTools is NOT working under the Catalina OS. Please see the following note: <https://ccsb.scripps.edu/mgltools/>.

If you substitute the executables in *misc\linux_third_party_software\vina* folder for Mac ones, and follow the same steps described for the Linux installation, I expect that SAnDReS would run on some of the Mac OS. But there is a problem with Catalina OS.

12) Can I use SAnDReS to run docking simulation for one structure?

Answer. Yes. See [Tutorial 3](#). But we did not develop SAnDReS for this. Our main goal is to provide a free software to explore the Scoring Function Space. In doing so, SAnDReS generates an adequate scoring function to predict binding and ranking poses. There is one limitation to running SAnDReS for one structure. This structure needs to have experimental binding affinity in the ligand databases. You may edit the ligand database yourself and add any missing ligand data.

13) Can I export a machine learning model to another project directory?

Answer. Yes. It is necessary to copy the joblib file to the *models* folder in your project directory. For instance, in [Tutorial 1](#), we selected the OrthogonalMatchingPursuit model. We can find it in the *C:\Users\Walter\sandres2_win\datasets\CDK2_IC50\models* folder. Just copy the *model_OrthogonalMatchingPursuit.joblib* file to the new project directory (in the *models* folder). To use the model, follow Tutorial 1 ([Section 5.8](#)).

14) How can I make available a machine learning model generated with SAnDReS?

Answer. One way to do that is to make the joblib file available on GitHub. If you want to cite it in a paper, include the link to the model in your manuscript. If you are not familiar with GitHub or did not want to keep a link for it, we would be glad to make your model available in the SAnDReS GitHub

(<https://github.com/azevedolab/sandres>). Just send me an email <mailto:sandres@azevedolab.net> that I upload your model.

We are working on the development of a database of SAnDReS models. Once we have this database, users may submit their models. We will include in future versions of SAnDReS a tool to download previously generated machine learning models.

15) Can I use another version of Scikit-Learn?

Answer. No. You should stick with Scikit-Learn version 1.0.2. SAnDReS might work with a newer version of Scikit-Learn, but we know it works with version 1.0.2.

16) Are there any executables for SAnDReS?

Answer. No. During the development of this version, we generated working executables for Linux and Windows using auto-py-to-exe. But as soon as we moved to Scikit-Learn version 1.0.2, the executables we generated with this new version did not work anymore. We apologize, but you should follow the installation procedure described here ([Section 3](#) and [Section 4](#)). We used auto-py-to-exe to generate our executables, but it did not work for Scikit-Learn 1.0.2.

17) Can I apply SAnDReS to create a general scoring function?

Answer. Yes. You may use as your dataset thousands of structures for which inhibition constant data is available, for instance. You may see [Tutorial 4](#). Instead of limiting the search to structures to a specific protein, you may search the PDB for those with K_i data and insert the PDB access codes in SAnDReS and proceed as shown to generate machine learning models.

18) Can I contribute to the development of SAnDReS?

Answer. Yes. We welcome anyone interested in participating in this project. Please send an email <mailto:sandres@azevedolab.net>.

19) Are there people using SAnDReS?

Answer. Yes. SAnDReS 1.0 had over 90 citations (<https://scholar.google.com.br/scholar?oi=bibs&hl=pt-BR&cites=12867895474248966104>).

20) What is the 3D plot shown in Tutorial 1 ([Section 5.9](#))?

Answer. This plot shows intermolecular contacts involving ligands and protein in all structures in the dataset. We split into main-chain and side-chain contacts. Our goal is to identify possible patterns of ligand binding by analyzing these types of interactions. For instance, looking at side-chain atoms, we see Asp as the top residue, but with substantial contacts involving: Ala, Ile, Leu, Lys, Phe, and Val. We see a concentration of hydrophobic residues for contacts with side-chain atoms.

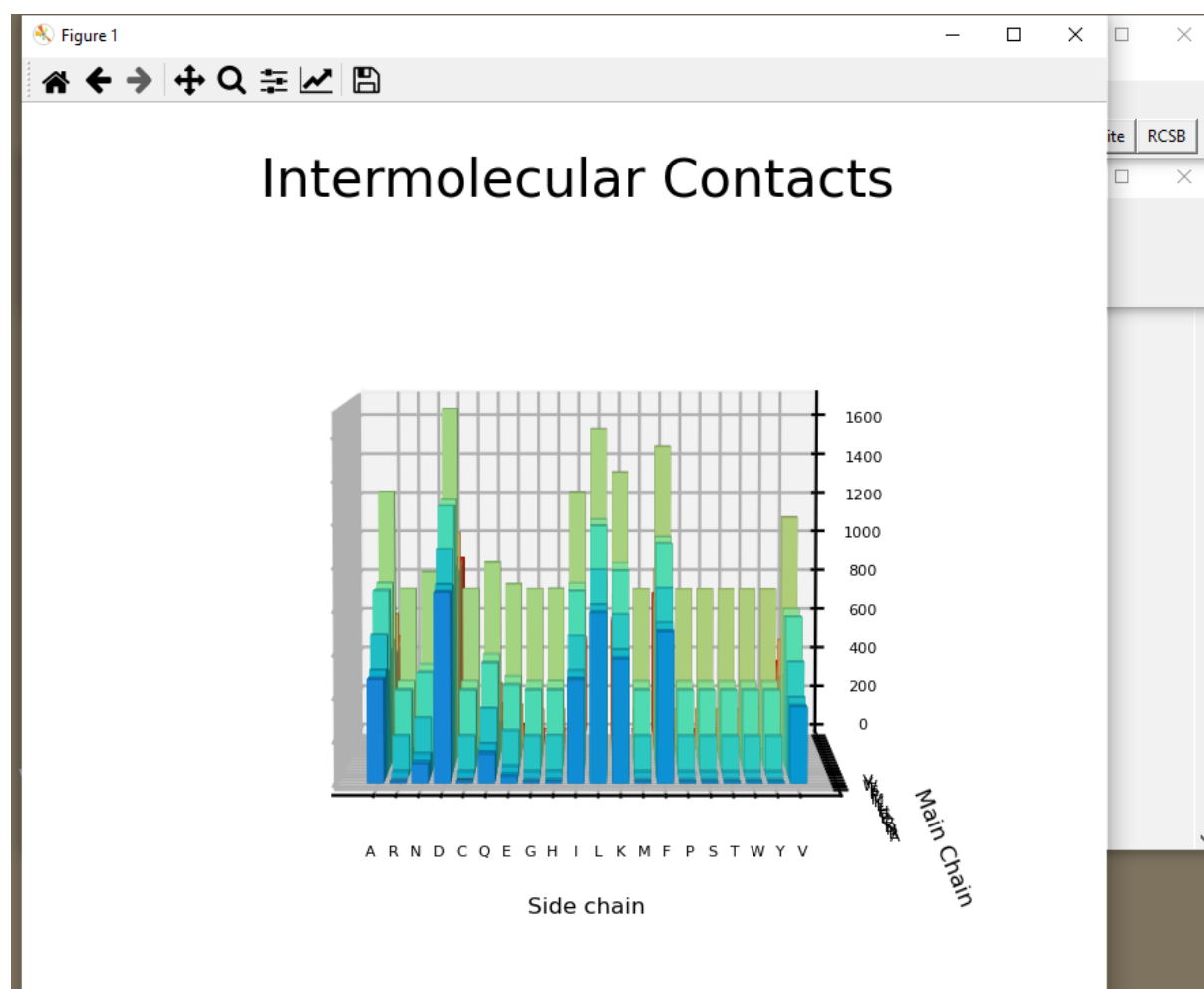


Figure 61. Intermolecular contacts (Windows Version).

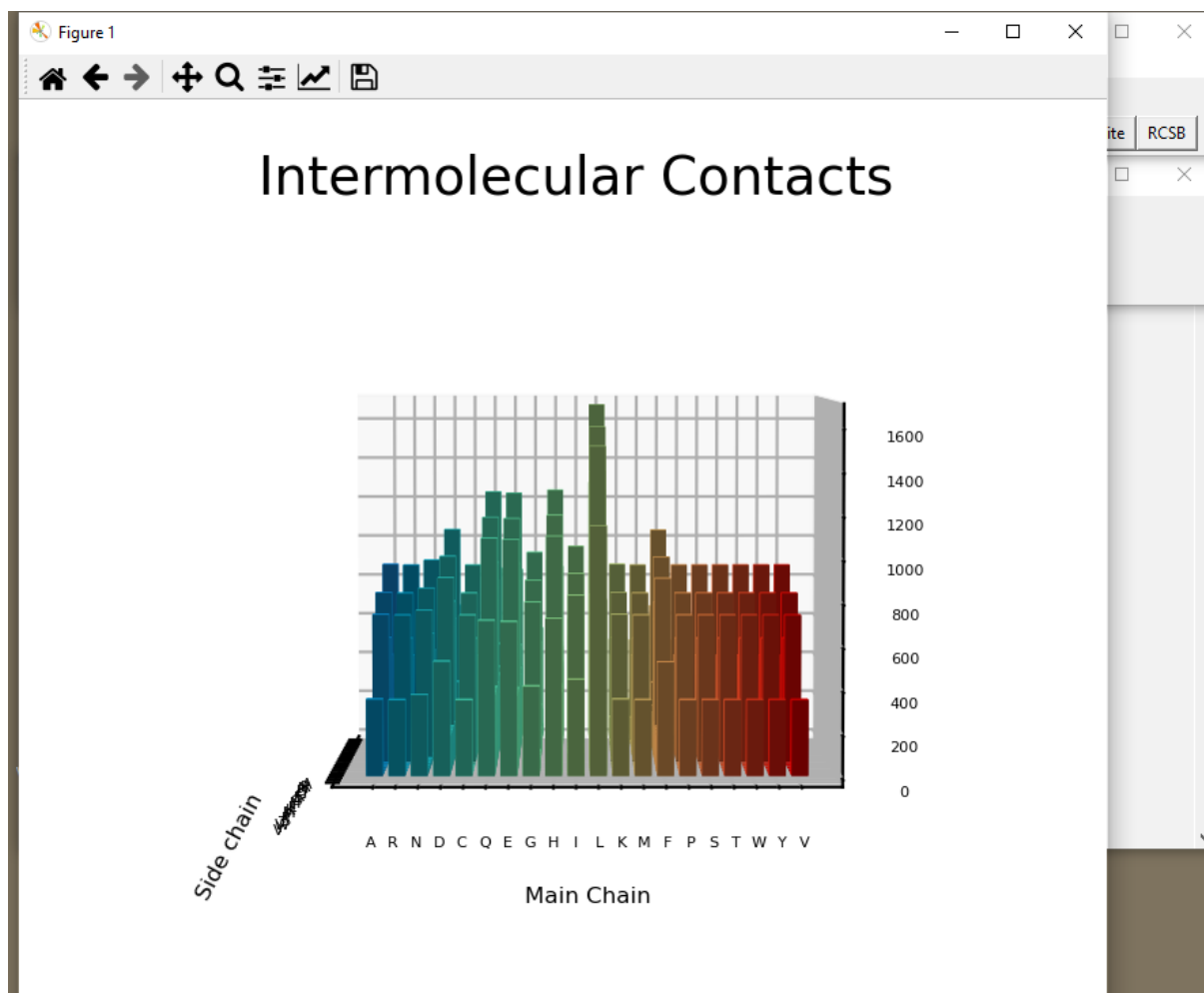


Figure 62. Intermolecular contacts (Windows Version).

Looking at main-chain atoms, we see a peak for Leu. Such analysis helps to capture overall aspects of the intermolecular interactions.

21) What is new in the SAnDReS 2.0?


Answer. It is a new software with the same goal, to explore the Scoring Function Space using machine learning methods. SAnDReS 1.0 generated machine learning models based on nine regression models. Now we have 64 methods. SAnDReS 1.0 used AutoDock Vina 1.1.2. Now we have the most recent version, AutoDock Vina 1.2.3. There are a lot of minor modifications, but these previously highlighted are the most important ones.

11. GitHub

SAnDReS source code is available to download on GitHub:

<https://github.com/azevedolab/sandres>

12. Authors



Gabriela Bitencourt-Ferreira

Doctorate Student - Pontifícia Universidade Católica do Rio Grande do Sul

Web of Science ResearcherID[®]
O-3423-2018

| PUBLICATIONS | TOTAL TIMES CITED | H-INDEX |
|--------------|-------------------|-----------------|
| 23 | 304 | 10 [®] |



Walter F de Azevedo, Jr.

"Walter Azevedo Jr."
[Show more](#)

Pontifícia Universidade Católica do Rio Grande do Sul

Web of Science ResearcherID[®]
N-5585-2018

| PUBLICATIONS | TOTAL TIMES CITED | H-INDEX | VERIFIED REVIEWS | VERIFIED EDITOR RECORDS |
|--------------|-------------------|-----------------|------------------|-------------------------|
| 203 | 6.541 | 44 [®] | 227 | 29 |

SAnDReS 2.0 was developed by Ms. [Gabriela Bitencourt-Ferreira](#) and [Dr. Walter F. de Azevedo, Jr.](#) This project is part of the doctoral thesis of Ms. Bitencourt-Ferreira.

13. Acknowledgements

The authors express their gratitude to those who participated in the development and testing of SAnDReS 1.0: Mariana Morrone Xavier, Gabriela Sehnem Heck, Mauricio Boff de Avila, Nayara Maria Bernhardt Levin, Val Oliveira Pinto, Nathalia Lemes Carvalho, Ariadne de Castro Silvério, Bruna Boldrini de Mattos, and Amauri Duarte da Silva. We also would like to thank Prof. Nelson José Freitas da Silveira, Prof. Alexandra Martins Dos Santos Soares, and Wallyson André Dos Santos Bezerra for their valuable suggestions. GB-F acknowledges support from Capes fellowship. WFA is a researcher for CNPq (Brazil) (Process Number: 309029/2018-0). This study was financed in part by the Coordenação de Aperfeiçoamento de Pessoal de Nível Superior - Brasil (CAPES) – Finance Code 001.

14. Colophon

We performed protein-ligand docking simulations and machine learning modeling reported on this User Guide using a Notebook with 4GB of memory, a 1 TB hard disk, and an Intel® Core® i3-7020U @ 2.30 GHz processor running Windows 10.

15. References

- Ballester PJ. Selecting machine-learning scoring functions for structure-based virtual screening. *Drug Discov Today Technol.* 2019; 32-33: 81–87.
- Bitencourt-Ferreira G, de Azevedo Jr. WF. Development of a machine-learning model to predict Gibbs free energy of binding for protein-ligand complexes. *Biophys Chem.* 2018; 240: 63–69.
- Bitencourt-Ferreira G, de Azevedo WF Jr. SAnDReS: A Computational Tool for Docking. *Methods Mol Biol.* 2019; 2053: 51–65.
- Bitencourt-Ferreira G, de Azevedo WF Jr. Machine Learning to Predict Binding Affinity. *Methods Mol Biol.* 2019; 2053: 251–273.
- Bitencourt-Ferreira G, de Azevedo WF Jr. Exploring the Scoring Function Space. *Methods Mol Biol.* 2019; 2053: 275–281.
- Bitencourt-Ferreira G, Duarte da Silva A, Filgueira de Azevedo W Jr. Application of Machine Learning Techniques to Predict Binding Affinity for Drug Targets: A Study of Cyclin-Dependent Kinase 2. *Curr Med Chem.* 2021; 28(2): 253–265.
- Bitencourt-Ferreira G, Rizzotto C, de Azevedo Junior WF. Machine Learning-Based Scoring Functions, Development and Applications with SAnDReS. *Curr Med Chem.* 2021; 28(9):1746–1756.
- Bohacek RS, McMartin C, Guida WC. The art and practice of structure-based drug design: a molecular modeling perspective. *Med Res Rev.* 1996;16(1): 3-50.
- da Silva AD, Bitencourt-Ferreira G, de Azevedo WF Jr. Taba: A Tool to Analyze the Binding Affinity. *J Comput Chem.* 2020; 41(1): 69–73.
- de Ávila MB, Xavier MM, Pinto VO, de Azevedo WF. Supervised machine learning techniques to predict binding affinity. A study for cyclin-dependent kinase 2. *Biochem Biophys Res Commun.* 2017; 494: 305–310.
- de Azevedo WF Jr, Mueller-Dieckmann HJ, Schulze-Gahmen U, Worland PJ, Sausville E, Kim SH. Structural basis for specificity and potency of a flavonoid inhibitor of human CDK2, a cell cycle kinase. *Proc Natl Acad Sci U S A.* 1996; 93(7): 2735–2740.
- de Azevedo WF, Leclerc S, Meijer L, Havlicek L, Strnad M, Kim SH. Inhibition of cyclin-dependent kinases by purine analogues: crystal structure of human cdk2 complexed with roscovitine. *Eur J Biochem.* 1997; 243(1-2): 518–526.
- de Azevedo W Jr, Canduri F, Marangoni dos Santos D, Pereira JH, Dias MV, Silva RG, Mendes MA, Basso LA, Palma MS, Santos DS. Structural basis for inhibition of human PNP by immucillin-H. *Biochem Biophys Res Commun.* 2003; 309(4): 917–922.
- de Azevedo WF Jr, dos Santos GC, dos Santos DM, Olivieri JR, Canduri F, Silva RG, Basso LA, Renard G, da Fonseca IO, Mendes MA, Palma MS, Santos DS.

- Docking and small angle X-ray scattering studies of purine nucleoside phosphorylase. *Biochem Biophys Res Commun.* 2003; 309(4): 923–928.
- de Azevedo Junior WF, Bitencourt-Ferreira G, Godoy JR, Adriano HMA, Dos Santos Bezerra WA, Dos Santos Soares AM. Protein-ligand Docking Simulations with AutoDock4 Focused on the Main Protease of SARS-CoV-2. *Curr Med Chem.* 2021; 28(37): 7614–7633.
- de Azevedo Junior WF. Application of Machine Learning Techniques for Drug Discovery. *Curr Med Chem.* 2021; 28(38): 7805–7807.
- Eberhardt J, Santos-Martins D, Tillack AF, Forli S. AutoDock Vina 1.2.0: New Docking Methods, Expanded Force Field, and Python Bindings. *J Chem Inf Model.* 2021; 61(8):3891–3898.
- Gilson MK, Liu T, Baitaluk M, Nicola G, Hwang L, Chong J. BindingDB in 2015: A public database for medicinal chemistry, computational chemistry and systems pharmacology. *Nucleic Acids Res.* 2016; 44(D1): D1045–1053.
- Gramatica P. On the development and validation of QSAR models. *Methods Mol Biol.* 2013; 930: 499–526.
- Halevy A, Norvig P, Pereira F. The Unreasonable Effectiveness of Data. *IEEE Intell. Syst.* 2009; 24(2): 8–12.
- Heck GS, Pinto VO, Pereira RR, de Ávila MB, Levin NMB, de Azevedo WF. Supervised Machine Learning Methods Applied to Predict Ligand- Binding Affinity. *Curr Med Chem.* 2017; 24(23): 2459–2470.
- Levin NMB, Pinto VO, Bitencourt-Ferreira G, Mattos BB, Silvério AC, de Azevedo Jr. WF. Development of CDK-targeted scoring functions for prediction of binding affinity. *Biophys Chem.* 2018; 235: 1–8.
- Morris GM, Huey R, Lindstrom W, Sanner MF, Belew RK, Goodsell DS, Olson AJ. AutoDock4 and AutoDockTools4: Automated docking with selective receptor flexibility. *J Comput Chem.* 2009; 30(16): 2785–2791.
- Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, Blondel M, Prettenhofer P, Weiss R, Dubourg V, Verplas J, Passos A, Cournapeau D, Brucher M, Perrot M, Duchesnay E. Scikit-learn: Machine Learning in Python. *J Mach Learn Res.* 2011; 12: 2825–2830.
- Pinto VO, Azevedo WF. Optimized Virtual Screening Workflow. Towards Target-Based Polynomial Scoring Functions for HIV-1 Protease. *Comb Chem High Throughput Screen.* 2017; 20(9): 820–827.
- Smith JM. Natural selection and the concept of a protein space. *Nature.* 1970; 225(5232): 563–564.
- Thomsen R, Christensen MH. MolDock: a new technique for high-accuracy molecular docking. *J Med Chem.* 2006; 49(11): 3315–3321.

Trott O, Olson AJ. AutoDock Vina: improving the speed and accuracy of docking with a new scoring function, efficient optimization, and multithreading. *J Comput Chem*. 2010; 31(2):455–461.

Veit-Acosta M, de Azevedo Junior WF. The Impact of Crystallographic Data for the Development of Machine Learning Models to Predict Protein-Ligand Binding Affinity. *Curr Med Chem*. 2021; 28(34): 7006–7022.

Wang R, Fang X, Lu Y, Wang S. The PDBbind database: collection of binding affinities for protein-ligand complexes with known three-dimensional structures. *J Med Chem*. 2004; 47(12): 2977–2980.

Wójcikowski M, Siedlecki P, Ballester PJ. Building Machine-Learning Scoring Functions for Structure-Based Prediction of Intermolecular Binding Affinity. *Methods Mol Biol*. 2019;2053:1–12.

Xavier MM, Heck GS, Avila MB, Levin NMB, Pintro VO, Carvalho NL, Azevedo WF. SAnDReS a Computational Tool for Statistical Analysis of Docking Results and Development of Scoring Functions. *Comb Chem High Throughput Screen*. 2016; 19(10): 801–812.

Appendix. Regression Methods

On the following page, we have a description of all parameters used for each regression method available in SAnDReS. The values for these parameters are in the file *misc/data/ml.in*. Mostly, we keep the same names used in the Scikit-Learn manual. The exceptions are the parameters: *rand_in+* and *cv_in**. In the following table, we provide links for complete descriptions of the parameters for each regression method.

Half of the regression methods available SAnDReS use cross-validation (CV). We implemented the *Kfold* class from Scikit-Learn to perform cross-validation. *Kfold* class builds an *n*-fold cross-validation loop and tests the generalization ability of regression. Cross-validation is a better estimate of how well we could generalize to predict unseen data.

Scikit-Learn provides some regression classes with built-in cross-validation implementation, e.g., *ElasticNetCV*. But this inclusion of built-in CV is not available for all regression methods (e.g., *AdaBoostRegressor*). Therefore, we adopted the same CV approach (Coelho & Richert, 2015) for the regression methods in SAnDReS 2.0. The *MLRegMPy* package has a class (*ValidationLoop*) that carries out CV for the 32 CV methods.

Reference: Coelho LP, Richert W. (2015) Building Machine Learning Systems with Python. 2nd ed. Packt Publishing Ltd. Birmingham UK. 301 pp. See page 162 (Cross-validation for regression).

Before applying the regression methods, SAnDReS scales the data using the following methods available in Scikit-Learn: [StandardScaler](#), [MaxAbsScaler](#), [MinMaxScaler](#), and [RobustScaler](#).

Regression methods available in SAnDReS.

| Method | Cross Validation | Parameters in <i>ml.in</i> file | Default values for parameters in <i>ml.in</i> file | Link |
|---------------------|------------------|--|--|---|
| AdaBoostRegressor | None | <i>Regression method, base_estimator, n_estimators, learning_rate, loss, random_state</i> | <i>AdaBoostRegressor, None, 50, 1.0, linear, None</i> | https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.AdaBoostRegressor.html |
| AdaBoostRegressorCV | Kfold class | <i>Regression method, base_estimator, n_estimators, learning_rate, loss, random_state, cv_in*</i> | <i>AdaBoostRegressorCV, None, 50, 1.0, linear, None, 5</i> | https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.AdaBoostRegressor.html |
| ARDRegression | None | <i>Regression method, n_iter, tol, alpha_1, alpha_2, lambda_1, lambda_2, compute_score, threshold_lambda, fit_intercept, copy_X, verbose, rand_in*</i> | <i>ARDRegression, 1000, 1e-3, 1e-6, 1e-6, 1e-6, 1e-6, False, 10000, True, True, False, 1123581321</i> | https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.ARDRegression.html |
| ARDRegressionCV | Kfold class | <i>Regression method, n_iter, tol, alpha_1, alpha_2, lambda_1, lambda_2, compute_score, threshold_lambda, fit_intercept, copy_X, verbose, rand_in*, cv_in*</i> | <i>ARDRegressionCV, 1000, 1e-3, 1e-6, 1e-6, 1e-6, 1e-6, False, 10000, True, True, False, 1123581321, 5</i> | https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.ARDRegression.html |
| BaggingRegressor | None | <i>Regression method, base_estimator, n_estimators, max_samples, max_features, bootstrap, bootstrap_features, oob_score, warm_start,</i> | <i>BaggingRegressor, None, 10, 1.0, 1.0, True, False, False, False, -1, None, 0</i> | https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.BaggingRegressor.html |

| | | | | |
|-----------------------|-------------|--|---|---|
| | | <i>n_jobs, random_state, verbose</i> | | |
| BaggingRegressorCV | Kfold class | <i>Regression method, base_estimator, n_estimators, max_samples, max_features, bootstrap, bootstrap_features, oob_score, warm_start, n_jobs, random_state, verbose, cv_in*</i> | <i>BaggingRegressorCV, None, 10, 1.0, 1.0, True, False, False, False, -1, None, 0, 5</i> | https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.BaggingRegressor.html |
| BayesianRidge | None | <i>Regression method, n_iter, tol, alpha_1, alpha_2, lambda_1, lambda_2, alpha_init, lambda_init, compute_score, fit_intercept, copy_X, verbose, rand_in*</i> | <i>BayesianRidge, 1000, 1e-3, 1e-6, 1e-6, 1e-6, 1e-6, None, None, False, True, True, False, 1123581321</i> | https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.BayesianRidge.html |
| BayesianRidgeCV | Kfold class | <i>Regression method, n_iter, tol, alpha_1, alpha_2, lambda_1, lambda_2, alpha_init, lambda_init, compute_score, fit_intercept, copy_X, verbose, rand_in*, cv_in*</i> | <i>BayesianRidgeCV, 1000, 1e-3, 1e-6, 1e-6, 1e-6, 1e-6, None, None, False, True, True, False, 1123581321, 5</i> | https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.BayesianRidge.html |
| DecisionTreeRegressor | None | <i>Regression method, criterion, splitter, max_depth, min_samples_split, min_samples_leaf, min_weight_fraction_leaf, max_features, random_state, max_leaf_nodes,</i> | <i>DecisionTreeRegressor, squared_error, best, None, 2, 1, 0.0, None, None, None, 0.0, 0.0</i> | https://scikit-learn.org/stable/modules/generated/sklearn.tree.DecisionTreeRegressor.html#sklearn.tree.DecisionTreeRegressor |

| | | | | |
|-------------------------|-------------|---|--|---|
| | | <i>min_impurity_decrease, ccp_alpha</i> | | |
| DecisionTreeRegressorCV | Kfold class | <i>Regression method, criterion, splitter, max_depth, min_samples_split, min_samples_leaf, min_weight_fraction_leaf, max_features, random_state, max_leaf_nodes, min_impurity_decrease, ccp_alpha, cv in*</i> | <i>DecisionTreeRegressorCV, squared_error, best, None, 2, 1, 0.0, None, None, None, 0.0, 0.0, 5</i> | https://scikit-learn.org/stable/modules/generated/sklearn.tree.DecisionTreeRegressor.html#sklearn.tree.DecisionTreeRegressor |
| ElasticNet | None | <i>Regression method, alpha, l1_ratio, fit_intercept, precompute, max_iter, copy_X, tol, warm_start, positive, random_state, selection, rand_in*</i> | <i>ElasticNet, 1.0, 0.5, True, False, 1000, True, 1e-4, False, False, None, cyclic, 123581321</i> | https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.ElasticNet.html |
| ElasticNetCV | Kfold class | <i>Regression method, alpha, l1_ratio, fit_intercept, precompute, max_iter, copy_X, tol, warm_start, positive, random_state, selection, rand_in*, cv in*</i> | <i>ElasticNetCV, 1.0, 0.5, True, False, 1000, True, 1e-4, False, False, None, cyclic, 123581321, 5</i> | https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.ElasticNet.html |
| ExtraTreeRegressor | None | <i>Regression method, criterion, splitter, max_depth, min_samples_split, min_samples_leaf, min_weight_fraction_leaf, max_features, random_state, min_impurity_decrease, max_leaf_nodes, ccp_alpha</i> | <i>ExtraTreeRegressor, squared_error, random, None, 2, 1, 0.0, auto, None, 0.0, None, 0.0</i> | https://scikit-learn.org/stable/modules/generated/sklearn.tree.ExtraTreeRegressor.html |

| | | | | |
|-----------------------|-------------|---|--|---|
| ExtraTreeRegressorCV | Kfold class | <i>Regression method, criterion, splitter, max_depth, min_samples_split, min_samples_leaf, min_weight_fraction_leaf, max_features, random_state, min_impurity_decrease, max_leaf_nodes, ccp_alpha, cv in*</i> | <i>ExtraTreeRegressorCV, squared_error, random, None, 2, 1, 0.0, auto, None, 0.0, None, 0.0, 5</i> | https://scikit-learn.org/stable/modules/generated/sklearn.tree.ExtraTreeRegressor.html |
| ExtraTreesRegressor | None | <i>Regression method, n_estimators, criterion, max_depth, min_samples_split, min_samples_leaf, min_weight_fraction_leaf, max_features, max_leaf_nodes, min_impurity_decrease, bootstrap, oob_score, n_jobs, random_state, verbose, warm_start, ccp_alpha, max_samples</i> | <i>ExtraTreesRegressor, 142, squared_error, None, 2, 1, 0.0, auto, None, 0.0, False, False, -1, 1123581321, 0, False, 0.0, None</i> | https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.ExtraTreesRegressor.html |
| ExtraTreesRegressorCV | Kfold class | <i>Regression method, n_estimators, criterion, max_depth, min_samples_split, min_samples_leaf, min_weight_fraction_leaf, max_features, max_leaf_nodes, min_impurity_decrease, bootstrap, oob_score, n_jobs, random_state, verbose, warm_start,</i> | <i>ExtraTreesRegressorCV, 142, squared_error, None, 2, 1, 0.0, auto, None, 0.0, False, False, -1, 1123581321, 0, False, 0.0, None, 5</i> | https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.ExtraTreesRegressor.html |

| | | | | |
|-----------------------------|-------------|---|---|---|
| | | <code>ccp_alpha, max_samples, cv_in*</code> | | |
| GaussianProcessRegressor | None | <code>Regression method, kernel, alpha, optimizer, n_restarts_optimizer, normalize_y, copy_X_train, random_state</code> | <code>GaussianProcessRegressor, None, 1e-10, fmin_l_bfgs_b, 0, False, True, None</code> | https://scikit-learn.org/stable/modules/generated/sklearn.gaussian_process.GaussianProcessRegressor.html |
| GaussianProcessRegressorCV | Kfold class | <code>Regression method, kernel, alpha, optimizer, n_restarts_optimizer, normalize_y, copy_X_train, random_state, cv_in*</code> | <code>GaussianProcessRegressorCV, None, 1e-10, fmin_l_bfgs_b, 0, False, True, None, 5</code> | https://scikit-learn.org/stable/modules/generated/sklearn.gaussian_process.GaussianProcessRegressor.html |
| GradientBoostingRegressor | None | <code>Regression method, loss, learning_rate, n_estimators, subsample, criterion, min_samples_split, min_samples_leaf, min_weight_fraction_leaf, max_depth, min_impurity_decrease, init, random_state, max_features, alpha, verbose, max_leaf_nodes, warm_start, validation_fraction, n_iter_no_change, tol, ccp_alpha</code> | <code>GradientBoostingRegressor, squared_error, 0.1, 100, 1.0, friedman_mse, 2, 1, 0.0, 3, 0.0, None, None, None, 0.9, 0, None, False, 0.1, None, 1e-4, 0.0</code> | https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.GradientBoostingRegressor.html |
| GradientBoostingRegressorCV | Kfold class | <code>Regression method, loss, learning_rate, n_estimators, subsample, criterion, min_samples_split, min_samples_leaf,</code> | <code>GradientBoostingRegressorCV, squared_error, 0.1, 100, 1.0, friedman_mse, 2, 1, 0.0, 3, 0.0, None, None, None, 0.9, 0, None, False, 0.1, None, 1e-4, 0.0, 5</code> | https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.GradientBoostingRegressor.html |

| | | | | |
|---------------------|-------------|--|---|---|
| | | <i>min_weight_fraction_leaf, max_depth, min_impurity_decrease, init, random_state, max_features, alpha, verbose, max_leaf_nodes, warm_start, validation_fraction, n_iter_no_change, tol, ccp_alpha, cv</i> in* | | |
| HuberRegressor | None | <i>Regression method, epsilon, max_iter, alpha, warm_start, fit_intercept, tol, rand_in*</i> | <i>HuberRegressor, 1.35, 1000, 0.0001, False, True, 1e-5, 1123581321</i> | https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.HuberRegressor.html |
| HuberRegressorCV | Kfold class | <i>Regression method, epsilon, max_iter, alpha, warm_start, fit_intercept, tol, rand_in*, cv_in*</i> | <i>HuberRegressorCV, 1.35, 1000, 0.0001, False, True, 1e-5, 1123581321, 5</i> | https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.HuberRegressor.html |
| KernelRidge | None | <i>Regression method, alpha, kernel, gamma, degree, coef0, kernel_params</i> | <i>KernelRidge, 1.0, linear, None, 3.0, 1.0, None</i> | https://scikit-learn.org/stable/modules/generated/sklearn.kernel_ridge.KernelRidge.html |
| KernelRidgeCV | Kfold class | <i>Regression method, alpha, kernel, gamma, degree, coef0, kernel_params, cv_in*</i> | <i>KernelRidgeCV, 1.0, linear, None, 3.0, 1.0, None, 5</i> | https://scikit-learn.org/stable/modules/generated/sklearn.kernel_ridge.KernelRidge.html |
| KNeighborsRegressor | None | <i>Regression method, n_neighbors, weights, algorithm, leaf_size, p, metric, metric_params, n_jobs</i> | <i>KNeighborsRegressor, 5, uniform, auto, 30, 2, minkowski, None, -1</i> | https://scikit-learn.org/stable/modules/generated/sklearn.neighbors.KNeighborsRegressor.html |

| | | | | |
|-----------------------|-------------|---|--|---|
| KneighborsRegressorCV | Kfold class | <i>Regression method, n_neighbors, weights, algorithm, leaf_size, p, metric, metric_params, n_jobs, cv_in*</i> | <i>KNeighborsRegressorCV, 5, uniform, auto, 30, 2, minkowski, None, -1, 5</i> | https://scikit-learn.org/stable/modules/generated/sklearn.neighbors.KNeighborsRegressor.html |
| Lars | None | <i>Regression method, fit_intercept, verbose, precompute, n_nonzero_coefs, eps_0#, eps_1#, n_samples#, copy_X, fit_path, jitter, random_state</i> | <i>Lars, True, False, auto, 500, 0.25, 1.00, 50, True, True, None, 1123581321</i> | https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.Lars.html |
| LarsCV | Kfold class | <i>Regression method, fit_intercept, verbose, precompute, n_nonzero_coefs, eps_0#, eps_1#, n_samples#, copy_X, fit_path, jitter, random_state, cv_in*</i> | <i>LarsCV, True, False, auto, 500, 0.25, 1.00, 50, True, True, None, 1123581321, 5</i> | https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.Lars.html |
| Lasso | None | <i>Regression method, alpha, fit_intercept, precompute, copy_X, max_iter, tol, warm_start, positive, random_state, selection</i> | <i>Lasso, 0.65, True, False, True, 1000, 1e-4, False, False, 1123581321, cyclic</i> | https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.Lasso.html |
| LassoCV | Kfold class | <i>Regression method, alpha, fit_intercept, precompute, copy_X, max_iter, tol, warm_start, positive, random_state, selection, cv_in*</i> | <i>LassoCV, 0.65, True, False, True, 1000, 1e-4, False, False, 1123581321, cyclic, 5</i> | https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.Lasso.html |

| | | | | |
|--------------------|-------------|--|---|---|
| LassoLars | None | <i>Regression method, alpha, fit_intercept, verbose, precompute, max_iter, eps, copy_X, fit_path, positive, jitter, random_state</i> | <i>LassoLars, 0.1, True, False, auto, 500, 0.25, True, True, False, None, 1123581321</i> | https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.LassoLars.html |
| LassoLarsCV | Kfold class | <i>Regression method, alpha, fit_intercept, verbose, precompute, max_iter, eps, copy_X, fit_path, positive, jitter, random_state, cv in*</i> | <i>LassoLarsCV, 0.1, True, False, auto, 500, 0.25, True, True, False, None, 1123581321, 5</i> | https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.LassoLars.html |
| LassoLarsIC | None | <i>Regression method, criterion, fit_intercept, verbose, precompute, max_iter, eps, copy_X, positive, rand_in*</i> | <i>LassoLarsIC, aic, True, False, auto, 500, 0.25, True, False, 1123581321</i> | https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.LassoLarsIC.html |
| LassoLarsICCV | Kfold class | <i>Regression method, criterion, fit_intercept, verbose, precompute, max_iter, eps, copy_X, positive, rand_in*, cv_in*</i> | <i>LassoLarsICCV, aic, True, False, auto, 500, 0.25, True, False, 1123581321, 5</i> | https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.LassoLarsIC.html |
| LinearRegression | None | <i>Regression method, fit_intercept, copy_X, n_jobs, positive, rand_in*</i> | <i>LinearRegression, True, True, -1, False, 1123581321</i> | https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.LinearRegression.html |
| LinearRegressionCV | Kfold class | <i>Regression method, fit_intercept, copy_X, n_jobs, positive, rand_in*, cv_in*</i> | <i>LinearRegressionCV, True, True, -1, False, 1123581321, 5</i> | https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.LinearRegression.html |
| LinearSVR | None | <i>Regression, epsilon, tol, C, loss, fit_intercept,</i> | <i>LinearSVR, 1e-2, 1e-8, 1.0, epsilon_insensitive, T</i> | https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.LinearSVR.html |

| | | | | |
|----------------|-------------|--|--|---|
| | | <code>intercept_scaling, dual, verbose, random_state, max_iter</code> | <code>rue,1.0,True,0,1123581321,1000</code> | s/generated/sklearn.svm.LinearSVR.html |
| LinearSVRCV | Kfold class | <code>Regression, epsilon, tol, C, loss, fit_intercept, intercept_scaling, dual, verbose, random_state, max_iter, cv_in"</code> | <code>LinearSVRCV,1e-2,1e-8,1.0,epsilon_insensitive,True,1.0,True,0,1123581321,1000,5</code> | https://scikit-learn.org/stable/modules/generated/sklearn.svm.LinearSVR.html |
| MLPRegressor | None | <code>Regression, hidden_layer_sizes, activation, solver, alpha, batch_size, learning_rate, learning_rate_init, power_t, max_iter, shuffle, random_state, tol, verbose, warm_start, momentum, nesterovs_momentum, early_stopping, validation_fraction, beta_1, beta_2, epsilon, n_iter_no_change, max_fun</code> | <code>MLPRegressor,75,logistic,lbfgs,0.01,auto,adaptive,0.001,0.5,200,True,1123581321,5e-3,False,False,0.9,True,False,0.1,0.9,0.999,1e-8,10,15000</code> | https://scikit-learn.org/stable/modules/generated/sklearn.neural_network.MLPRegressor.html |
| MLPRegressorCV | Kfold class | <code>Regression, hidden_layer_sizes, activation, solver, alpha, batch_size, learning_rate, learning_rate_init, power_t, max_iter, shuffle, random_state, tol, verbose, warm_start, momentum, nesterovs_momentum, early_stopping, validation_fraction,</code> | <code>MLPRegressorCV,75,logistic,lbfgs,0.01,auto,adaptive,0.001,0.5,200,True,1123581321,5e-3,False,False,0.9,True,False,0.1,0.9,0.999,1e-8,10,15000,5</code> | https://scikit-learn.org/stable/modules/generated/sklearn.neural_network.MLPRegressor.html |

| | | | | |
|------------------------------|-------------|--|---|---|
| | | <i>beta_1, beta_2, epsilon, n_iter_no_change, max_fun, cv in*</i> | | |
| NuSVR | None | <i>Regression method, nu, C, kernel, degree, gamma, coef0, shrinking, tol, cache_size, verbose, max_iter, rand_in*</i> | <i>NuSVR, 0.5, 1.0, linear, 1, auto, 0.0, True, 0.001, 200.0, False, -1, 1123581321</i> | https://scikit-learn.org/stable/modules/generated/sklearn.svm.NuSVR.html |
| NuSVRCV | Kfold class | <i>Regression method, nu, C, kernel, degree, gamma, coef0, shrinking, tol, cache_size, verbose, max_iter, rand_in*, cv in*</i> | <i>NuSVRCV, 0.5, 1.0, linear, 1, auto, 0.0, True, 0.001, 200.0, False, -1, 1123581321, 5</i> | https://scikit-learn.org/stable/modules/generated/sklearn.svm.NuSVR.html |
| OrthogonalMatchingPursuit | None | <i>Regression method, n_nonzero_coefs, tol, fit_intercept, precompute, rand_in*</i> | <i>OrthogonalMatchingPursuit, None, None, True, auto, 1123581321</i> | https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.OrthogonalMatchingPursuit.html |
| OrthogonalMatchingPursuitCV | Kfold class | <i>Regression method, n_nonzero_coefs, tol, fit_intercept, precompute, rand_in*, cv_in*</i> | <i>OrthogonalMatchingPursuitCV, None, None, True, auto, 1123581321, 5</i> | https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.OrthogonalMatchingPursuit.html |
| PassiveAggressiveRegressor | None | <i>Regression method, C, fit_intercept, max_iter, tol, early_stopping, validation_fraction, n_iter_no_change, shuffle, verbose, loss, epsilon, random_state, warm_start, average</i> | <i>PassiveAggressiveRegressor, 1.0, True, 1000, 1e-3, False, 0.1, 5, True, 0, epsilon, 1e-4, 1123581321, True, True</i> | https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.PassiveAggressiveRegressor.html |
| PassiveAggressiveRegressorCV | Kfold class | <i>Regression method, C, fit_intercept, max_iter, tol, early_stopping,</i> | <i>PassiveAggressiveRegressorCV, 1.0, True, 1000, 1e-3, False, 0.1, 5, True, 0, epsilon,</i> | https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.PassiveAggressiveRegressor.html |

| | | | | |
|-------------------------|-------------|---|---|---|
| | | <i>validation_fraction,</i> <i>n_iter_no_change,</i> <i>shuffle, verbose, loss,</i> <i>epsilon, random_state,</i> <i>warm_start, average,</i> <i>cv_in*</i> | <i>n_insensitive,1e-</i> <i>4,1123581321,True,True,5</i> | ear_model.PassiveAggressiveRegressor.html |
| RandomForestRegressor | None | <i>Regression method,</i> <i>n_estimators, criterion,</i> <i>max_depth,</i> <i>min_samples_split,</i> <i>min_samples_leaf,</i> <i>min_weight_fraction_leaf,</i> <i>max_features,</i> <i>max_leaf_nodes,</i> <i>min_impurity_decrease,</i> <i>bootstrap, oob_score,</i> <i>n_jobs, random_state,</i> <i>verbose, warm_start,</i> <i>ccp_alpha, max_samples</i> | <i>RandomForestRegressor,142,s</i> <i>quared_error,None,2,1,0.0,a</i> <i>uto,None,0.0,True,False,-</i> <i>1,1123581321,0,False,0.0,No</i> <i>ne</i> | https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestRegressor.html |
| RandomForestRegressorCV | Kfold class | <i>Regression method,</i> <i>n_estimators, criterion,</i> <i>max_depth,</i> <i>min_samples_split,</i> <i>min_samples_leaf,</i> <i>min_weight_fraction_leaf,</i> <i>max_features,</i> <i>max_leaf_nodes,</i> <i>min_impurity_decrease,</i> <i>bootstrap, oob_score,</i> <i>n_jobs, random_state,</i> <i>verbose, warm_start,</i> <i>ccp_alpha, max_samples,</i> <i>cv_in*</i> | <i>RandomForestRegressorCV,142</i> <i>,squared_error,None,2,1,0.0</i> <i>,auto,None,0.0,True,False,-</i> <i>1,1123581321,0,False,0.0,No</i> <i>ne,5</i> | https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestRegressor.html |
| RANSACRegressor | None | <i>Regression method,</i> <i>base_estimator,</i> <i>min_samples,</i> | <i>RANSACRegressor,None,None,N</i> <i>one,None,None,100,np.inf,np</i> | https://scikit-learn.org/stable/modules/generated/sklearn.lin |

| | | | | |
|-------------------|-------------|---|--|---|
| | | <i>residual_threshold,</i> <i>is_data_valid,</i> <i>is_model_valid,</i> <i>max_trials, max_skips,</i> <i>stop_n_inliers,</i> <i>stop_score,</i> <i>stop_probability, loss,</i> <i>random_state</i> | <i>.inf,np.inf,0.99,absolute_e</i> <i>rror,1123581321</i> | ear_model.RANSACRe gressor.html |
| RANSACRegressorCV | Kfold class | <i>Regression method,</i> <i>base_estimator,</i> <i>min_samples,</i> <i>residual_threshold,</i> <i>is_data_valid,</i> <i>is_model_valid,</i> <i>max_trials, max_skips,</i> <i>stop_n_inliers,</i> <i>stop_score,</i> <i>stop_probability, loss,</i> <i>random_state, cv_in*</i> | <i>RANSACRegressorCV,None,None</i> <i>,None,None,None,100,np.inf,</i> <i>np.inf,np.inf,0.99,absolute</i> <i>_error,1123581321,5</i> | https://scikit- learn.org/stable/module s/generated/sklearn.lin ear_model.RANSACRe gressor.html |
| Ridge | None | <i>Regression method, alpha,</i> <i>fit_intercept, copy_X,</i> <i>max_iter, tol, solver,</i> <i>positive, random_state</i> | <i>Ridge,1.0,True,True,None,1e</i> <i>-3,auto,False,None</i> | https://scikit- learn.org/stable/module s/generated/sklearn.lin ear_model.Ridge.html |
| RidgeCV | Kfold class | <i>Regression method, alpha,</i> <i>fit_intercept, copy_X,</i> <i>max_iter, tol, solver,</i> <i>positive, random_state,</i> <i>cv_in*</i> | <i>RidgeCV,1.0,True,True,None,</i> <i>1e-3,auto,False,None,5</i> | https://scikit- learn.org/stable/module s/generated/sklearn.lin ear_model.Ridge.html |
| SGDRegressor | None | <i>Regression method, loss,</i> <i>penalty, alpha, l1_ratio,</i> <i>fit_intercept, max_iter,</i> <i>tol, shuffle, verbose,</i> <i>epsilon, random_state,</i> <i>learning_rate, eta0,</i> <i>power_t, early stopping,</i> <i>validation_fraction,</i> | <i>SGDRegressor,squared_error,</i> <i>12,0.001,0.15,True,20000000</i> <i>00,1e-</i> <i>3,True,0,0.1,1123581321,inv</i> <i>scaling,0.01,0.25,False,0.1</i> <i>,5,False,False</i> | https://scikit- learn.org/stable/module s/generated/sklearn.lin ear_model.SGDRegres sor.html |

| | | | | |
|---------------------|-------------|---|--|---|
| | | <i>n_iter_no_change</i> , <i>warm_start</i> , <i>average</i> | | |
| SGDRegressorCV | Kfold class | <i>Regression method</i> , <i>loss</i> , <i>penalty</i> , <i>alpha</i> , <i>l1_ratio</i> , <i>fit_intercept</i> , <i>max_iter</i> , <i>tol</i> , <i>shuffle</i> , <i>verbose</i> , <i>epsilon</i> , <i>random_state</i> , <i>learning_rate</i> , <i>eta0</i> , <i>power_t</i> , <i>early_stopping</i> , <i>validation_fraction</i> , <i>n_iter_no_change</i> , <i>warm_start</i> , <i>average</i> , <i>cv</i> in* | <i>SGDRegressorCV</i> , <i>squared_error</i> , <i>12</i> , <i>0.001</i> , <i>0.15</i> , <i>True</i> , <i>200000</i> <i>0000</i> , <i>1e-</i> <i>3</i> , <i>True</i> , <i>0</i> , <i>0.1</i> , <i>1123581321</i> , <i>inv</i> <i>scaling</i> , <i>0.01</i> , <i>0.25</i> , <i>False</i> , <i>0.1</i> <i>5</i> , <i>False</i> , <i>False</i> , <i>5</i> | https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.SGDRegressor.html |
| SVR | None | <i>Regression method</i> , <i>kernel</i> , <i>degree</i> , <i>gamma</i> , <i>coef0</i> , <i>tol</i> , <i>C</i> , <i>epsilon</i> , <i>shrinking</i> , <i>cache_size</i> , <i>verbose</i> , <i>max_iter</i> , <i>rand_in</i> † | <i>SVR</i> , <i>linear</i> , <i>1</i> , <i>scale</i> , <i>0.0</i> , <i>1e-</i> <i>3</i> , <i>1.0</i> , <i>0.1</i> , <i>True</i> , <i>200.0</i> , <i>False</i> , <i>-1</i> , <i>1123581321</i> | https://scikit-learn.org/stable/modules/generated/sklearn.svm.SVR.html |
| SVRCV | Kfold class | <i>Regression method</i> , <i>kernel</i> , <i>degree</i> , <i>gamma</i> , <i>coef0</i> , <i>tol</i> , <i>C</i> , <i>epsilon</i> , <i>shrinking</i> , <i>cache_size</i> , <i>verbose</i> , <i>max_iter</i> , <i>rand_in</i> †, <i>cv</i> in* | <i>SVRCV</i> , <i>linear</i> , <i>1</i> , <i>scale</i> , <i>0.0</i> , <i>1e-</i> <i>3</i> , <i>1.0</i> , <i>0.1</i> , <i>True</i> , <i>200.0</i> , <i>False</i> , <i>-1</i> , <i>1123581321</i> , <i>5</i> | https://scikit-learn.org/stable/modules/generated/sklearn.svm.SVR.html |
| TheilSenRegressor | None | <i>Regression method</i> , <i>Regression method</i> , <i>fit_intercept</i> , <i>copy_X</i> , <i>max_subpopulation</i> , <i>n_subsamples</i> , <i>max_iter</i> , <i>tol</i> , <i>random_state</i> , <i>n_jobs</i> , <i>verbose</i> | <i>TheilSenRegressor</i> , <i>True</i> , <i>True</i> <i>10000</i> , <i>None</i> , <i>300</i> , <i>1e-</i> <i>3</i> , <i>1123581321</i> , <i>-1</i> , <i>False</i> | https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.TheilSenRegressor.html |
| TheilSenRegressorCV | Kfold class | <i>Regression method</i> , <i>fit_intercept</i> , <i>copy_X</i> , <i>max_subpopulation</i> , <i>n_subsamples</i> , <i>max_iter</i> , | <i>TheilSenRegressorCV</i> , <i>True</i> , <i>Tr</i> <i>ue</i> , <i>10000</i> , <i>None</i> , <i>300</i> , <i>1e-</i> <i>3</i> , <i>1123581321</i> , <i>-1</i> , <i>False</i> , <i>5</i> | https://scikit-learn.org/stable/modules/generated/sklearn.lin |

| | | | | |
|--------------------|-------------|--|--|---|
| | | <i>tol, random_state, n_jobs, verbose, cv_in*</i> | | ear_model.TheilSenRegressor.html |
| TweedieRegressor | None | <i>Regression method, power, alpha, fit_intercept, link, max_iter, tol, warm_start, verbose, rand_in*</i> | <i>TweedieRegressor,0.0,1.0,True,auto,100,1e-4,False,0,1123581321</i> | https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.TweedieRegressor.html |
| TweedieRegressorCV | Kfold class | <i>Regression method, power, alpha, fit_intercept, link, max_iter, tol, warm_start, verbose, rand_in*, cv_in*</i> | <i>TweedieRegressorCV,0.0,1.0,True,auto,100,1e-4,False,0,1123581321,5</i> | https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.TweedieRegressor.html |
| VotingRegressor | None | <i>Regression method, fit_intercept, copy_X, n_jobs, n_estimators, criterion, max_depth, min_samples_split, min_samples_leaf, min_weight_fraction_leaf, max_features, max_leaf_nodes, min_impurity_decrease, bootstrap, oob_score, n_jobs, random_state, verbose, warm_start, ccp_alpha, max_samples, weights, n_jobs, verbose</i> | <i>VotingRegressor,True,True,-1,142,squared_error,None,2,1,0.0,auto,None,0.0,True,False,None,1123581321,0,False,0.0,None,None,None,False</i> | https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.VotingRegressor.html |
| VotingRegressorCV | Kfold class | <i>Regression method, fit_intercept, copy_X, n_jobs,</i> | <i>VotingRegressorCV,True,True,-1,142,squared_error,None,2,1,0.0,auto,None,0.0,True,fa</i> | https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.VotingRegressor.html |

| | | | | |
|----------------|-------------|--|--|---|
| | | <code>n_estimators, criterion, max_depth, min_samples_split, min_samples_leaf, min_weight_fraction_leaf, max_features, max_leaf_nodes, min_impurity_decrease, bootstrap, oob_score, n_jobs, random_state, verbose, warm_start, ccp_alpha, max_samples, weights, n_jobs, verbose, cv_in*</code> | <code>lse, None, 1123581321, 0, False, 0.0, None, None, None, False, 5</code> | semble.VotingRegressor.html |
| XGBRegressor | None | <code>Regression method, booster, verbosity, validate_parameters, disable_default_eval_metric, eta, gamma, max_depth, min_child_weight, subsample, sampling_method, colsample_bytree, colsample_bylevel, colsample_bynode, lambda, alpha, tree_method, scale_pos_weight, refresh_leaf, process_type, predictor, objective, n_estimators</code> | <code>XGBRegressor, gbtrees, 1, True, False, 0.3, 0.0, 6, 1.0, 1.0, uniform, 1.0, 1.0, 1.0, 1.0, 10.0, auto, 1.0, 1, default, auto, reg:squarederror, 100</code> | https://xgboost.readthedocs.io/en/latest/parameter.html#general-parameters and https://xgboost.readthedocs.io/en/latest/python/python_api.html |
| XGBRegressorCV | Kfold class | <code>Regression method, booster, verbosity, validate_parameters,</code> | <code>XGBRegressorCV, gbtrees, 1, True, False, 0.3, 0.0, 6, 1.0, 1.0, uniform, 1.0, 1.0, 1.0, 1.0, 10.0</code> | https://xgboost.readthedocs.io/en/latest/parameter.html |

| | | | | |
|--|--|--|--|---|
| | | <code>disable_default_eval_metric, eta, gamma, max_depth, min_child_weight, subsample, sampling_method, colsample_bytree, colsample_bylevel, colsample_bynode, lambda, alpha, tree_method, scale_pos_weight, refresh_leaf, process_type, predictor, objective, n_estimators, cv_in*</code> | <code>,auto,1.0,1,default,auto,reg:squarederror,100,5</code> | eter.html#general-parameters and https://xgboost.readthedocs.io/en/latest/python/python_api.html |
|--|--|--|--|---|

**cv_in* variable holds an integer for the number of subsets used in the cross-validation process.

+*rand_in* holds a dummy integer seed to be used to generate a Molegro Data Modeller (MDM) format file.

#*eps_0*, *eps_1*, and *n_samples* define an array (eps) as follows:

```
eps_array = np.linspace(eps_0, eps_1, num=n_samples_in)
```

This is applied for machine-precision regularization in the computation of the Cholesky diagonal factors.