

Negative dataset (non-bioavailable)	Positive dataset (bioavailable)
<p>CASE 1.</p> <p>NIHMS negative, %F<50</p> <p>486 compounds</p> <p>Both datasets are downloadable from the article</p> <p>https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3955412/</p>	<p>NIHMS positive, %F>=50</p> <p>509 compounds</p>
<p>CASE 2.</p> <p>HSDB dataset, 3551 minus substances that are in approved drugs,</p> <p>3530 substances</p> <p><i>Data retrieval from pubchem API: pubchem_retrieval.ipynb</i></p> <p><i>Data FILE with data:hsdb-smiles-no-repeat.csv</i></p>	<p>ChEMBL approved and oral drugs, not withdrawn or discontinued + all drugs in trials</p> <p>3842 compounds</p> <p><i>Query FILE:</i></p> <p><i>chembl_ba_descriptors.sql</i></p>
<p>CASE 3.</p> <p>non-BA from curated data (DrugBank, CentralDrug combined),</p> <p>1451 compounds</p> <p><i>Query FILE: curated_non_ba_descriptors.sql</i></p>	<p>BA AND ORAL from curated data (DrugBank, CentralDrug combined),</p> <p>889 compounds</p> <p><i>Query FILE:</i></p> <p><i>curated_ba_descriptors.sql</i></p>
<p>CASE 4.</p> <p>Chembl compounds that are NOT in the clinical trials (and have descriptors)</p> <p>5590 compounds</p> <p><i>Query FILE: chembl_not_in_trials_descriptors.sql</i></p>	<p>ChEMBL approved and oral drugs, not withdrawn or discontinued + all drugs in trials</p> <p>3842 compounds</p> <p><i>Query FILE:</i></p> <p><i>chembl_ba_descriptors.sql</i></p>
<p>CASE 5.</p> <p>Chembl compounds with QED <0.3 and ro5_violations>=2</p> <p>1997 compounds</p> <p><i>Query FILE: chembl_qed_ro5.sql</i></p>	<p>ChEMBL approved and oral drugs, not withdrawn or discontinued + all drugs in trials</p> <p>3842 compounds</p> <p><i>Query FILE:</i></p> <p><i>chembl_ba_descriptors.sql</i></p>

Figure 1. Datasets used in training the model to predict on structure, descriptors, both

A dataset of around 13,000 compounds from curated data mart were selected to predict on (all_curated.sql query selects the data). The prediction was made on the Case 1 and Case 5 model by eliminating insignificant number of duplicates between the datasets to ensure that the prediction on the data of which of the model was not trained on. The predicted two probabilities (one probability per prediction) of the compound (to be bioavailable) are meant to be outputted to the user. This would happen in the application simultaneously with the bioavailability class and/or percent, if such existed in the data.

We also recorded the x and y coordinates of the Uniform Manifold Approximation & Projection (UMAP) from the data set on Case 1 model. The coordinates are meant to be used in the application to find the 10 closest compounds to the one that the user is interested in. Case1 dataset was obtained from: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3955412/> . Case 2 dataset was selected from the curated data mart.