

ADMET Property Prediction With Oloren ChemEngine

Andrew Li, David Huang

Oloren AI

Abstract

We present a set of model architectures built using the Oloren ChemEngine library for the Therapeutics Data Commons (TDC) benchmarks. We evaluate the models on ADMET benchmarks, representing various small molecule drug properties encompassing absorption, distribution, metabolism, excretion, and toxicity. The models predict values for these properties given the chemical structure of the molecule. The presented models draw from a wide set of state-of-art strategies, such as custom model gradient boosting, graph neural networks, molecular fingerprinting, and OlorenVec, a supervised learned molecular representation. We show that these models are superior to current leaderboard model approaches.

Introduction to OCE

Oloren ChemEngine¹ (OCE) is an open-source Python software package developed by Oloren AI for developing and using molecular property predictors. OCE implements current state-of-art methods such as message passing neural networks, graph neural networks, and molecular fingerprints and representations. It additionally allows for the creation of custom gradient boosting models consisting of user-defined submodels for each learner stage.

The TDC² ADMET benchmarks consist of 22 datasets composed of drug chemical structures and the property value for their respective ADMET property. Dataset sizes range from hundreds to ~15,000 molecules. We present the models which achieve a rank of first-place in 15 of these tasks, and top 3 in 4 tasks.

Model Architectures

We construct and test 9 unique models. This set of models consists of 1 random forest stacker, 2 random forests, and 6 custom gradient boosting models composed of various submodels.

The models are represented in the table below:

Model Name	Model Description	Number Parameters
BaseBoosting ADkCCrwJ	Random Forest (RDKit2DNormalized), Random Forest (Morgan Fingerprint), Random Forest (OlorenVec)	23
BaseBoosting OQRAYLPP	SPGNN, Random Forest (OlorenVec), BaseTorchGeometricModel (GINModel)	2,003,108
BaseBoosting xXOn1QFI	Random Forest (Morgan Fingerprint), Random Forest(Morgan Chiral Counts), Random Forest (Morgan Feature Counts), Random Forest (RDKit2DNormalized), Random Forest (OlorenVec)	27
BaseBoosting 1zpI0dIb	Random Forest (Morgan Fingerprint), Random Forest (Morgan Chiral Counts), Random Forest (Morgan Feature Counts), Random Forest (RDKit2DNormalized), Random Forest (OlorenVec)	27
BaseBoosting ktIIq91G	Random Forest (Morgan Fingerprint), Random Forest (Morgan Feature Counts), Random Forest (Morgan Chiral Counts), Random Forest (RDKit2DNormalized), Random Forest (Mordred Descriptor), Random Forest (OlorenVec)	40
BaseBoosting QjtOKx4i	Random Forest (RDKit2DNormalized), ChemProp	365713
RFStacker jTNhN7U7	Random Forest (RDKit2DNormalized) Random Forest (OlorenVec) Random Forest (Lipinski Descriptor) Random Forest (Morgan Feature Counts) SPGNN	1,858,225
RandomForestModel yjJ-ak-b	Random Forest (Morgan Feature Counts)	8

RandomForestModel gfFM673V	Random Forest (OlorenVec)	9
----------------------------	---------------------------	---

Model names are defined by OCE's uniquely generated name for a given model. BaseBoosting models are custom-made gradient boosting models. SPGNN represents the graph neural network implemented in Hu et al. (2020)³. Certain gradient boosting models have identical model and feature architectures but differ in terms of the parameterization of the constituent submodels, which can be observed in the source code. Random forest submodels are defined at a high-level by the molecular fingerprint or representation that they use. BaseTorchGeometricModel is an implementation of the graph neural network in Brossard et al. (2021)⁴. OlorenVec is a custom supervised molecular representation developed and trained by Oloren AI, available to use in OCE.

Results

Task	Model Name	Metric	Score	Stddev.	Leaderboard Rank
HIA	RFStacker jTNhN7U7	ROC-AUC	0.988	0.002	1
Pgp	BaseBoosting xXOn1QFI	ROC-AUC	0.946	0.001	1
BBB	BaseBoosting ADkCCrwJ	ROC-AUC	0.923	0.002	1
AMES	BaseBoosting 1zpI0dlb	ROC-AUC	0.865	0.002	1
DILI	BaseBoosting ADkCCrwJ	ROC-AUC	0.937	0.004	1
CYP3A4 Substrate	BaseBoosting ADkCCrwJ	ROC-AUC	0.679	0.033	1
Caco2	BaseBoosting QjtOKx4i	Mean Absolute Error	0.285	0.005	1
CYP2D6 Inhibition	BaseBoosting ktIIq91G	Average Precision	0.721	0.001	1
CYP3A4 Inhibition	BaseBoosting ktIIq91G	Average Precision	0.882	0.001	1
CYP2D6 Substrate	BaseBoosting OQRAYLPP	Average Precision	0.711	0.006	1
Clearance	RFStacker	Spearman	0.625	0.002	1

Microsome	jTNhN7U7				
Half Life	BaseBoosting ktIIq91G	Spearman	0.416	0.009	1
Clearance Hepatocyte	BaseBoosting xXOn1QFI	Spearman	0.491	0.006	1
hERG	RandomForest Model yjJ-ak-b	ROC-AUC	0.875	0.003	1
CYP2C9 Substrate	RandomForest Model gfFM673V	Average Precision	0.437	0.022	1
CYP2C9 Inhibition	BaseBoosting 1zpI0dlb	Average Precision	0.791	0.005	2
Bioavailability	BaseBoosting 1zpI0dlb	ROC-AUC	0.736	0.004	2
PPBR	BaseBoosting ktIIq91G	Mean Absolute Error	8.503	0.023	2
LD50	BaseBoosting ktIIq91G	Mean Absolute Error	0.613	0.001	3

Source Code

OCE is open-source and accessible at <https://github.com/Oloren-AI/olorenchemengine>. Code to define, train, and use these models is available publicly at <https://github.com/Oloren-AI/OCE-TDC>.

References

- [1] Oloren AI, Oloren ChemEngine, (2022), Github repository, <https://github.com/Oloren-AI/olorenchemengine>
- [2] Huang, Kexin, et al. Therapeutics Data Commons: Machine Learning Datasets and Tasks for Drug Discovery and Development. arXiv, 28 Aug. 2021. arXiv.org, <http://arxiv.org/abs/2102.09548>.
- [3] Hu, Weihua, et al. Strategies for Pre-Training Graph Neural Networks. arXiv, 18 Feb. 2020. arXiv.org, <https://doi.org/10.48550/arXiv.1905.12265>.

[4] Brossard, Rémy, et al. Graph Convolutions That Can Finally Model Local Structure. arXiv, 3 June 2021. arXiv.org, <https://doi.org/10.48550/arXiv.2011.15069>.