

# **Qptuna: easy, automatic QSAR modelling**

Lewis Mervin Molecular Al, Discovery Sciences, AstraZeneca R&D, Cambridge, UK **Company Restricted** 



## **Overview - Present**

- What is Qptuna?
- Current algo's/descriptors/features
- GUI



## **Overview - Future**

- Automated model building
- Publication/public release



# What is **Qptuna**?

- QPTUNA: QSAR using Optimization for Hyperparameter Tuning
- Build predictive models for CompChem with hyperparameters optimized by Optuna
- Qptuna library searches for the best ML algorithm and molecular descriptor for the given data
- Automatically apply best practices
- Docs soon to be available



Why?



Thomas M, Boardman A, Garcia-Ortegon M, Yang H, de Graaf C, Bender A. Applications of Artificial Intelligence in Drug Design: Opportunities and Challenges. *Methods Mol Biol*. 2022



#### **Qptuna**

#### • Why:

- We have assay data, e.g. 1k points.
- We want a QSAR model to predict assay response for new molecules.
- Which ML algorithm shall we use? Which molecular representation/fingerprint shall we use?

#### • What – looking for the best algorithms and fingerprints:

- "Traditional" methods: Grid search, Random search
- This work: Sequential model-based optimization (SMBO): sequentially construct models to approximate the objective based on historical measurements.

#### • How:

 Use Optuna library for hyperparameter optimization combined with packages such as Rdkit and e.g. ChemProp



optuna.readthedocs.io/



scikit-learn.org/stable/



chemprop.readthedocs.io/





# **Available algorithms and descriptors**

#### • ML Algorithms:

- Logistic Regression (aka logit, MaxEnt) classifier
- SVR (Epsilon-Support Vector Regression)
- Random Forest
- Ridge (Linear least squares with L2 regularization)
- Lasso (Linear model trained with L1 prior as regularizer)
- Cross decomposition using partial least squares (PLS)
- XGBregressor (XGBoost: gradient boosting trees algorithm)
- AdaBoost Classifier SVC (C-Support Vector Classification)
- SVC (C-Support Vector Classification)
- ChemProp regression/classification (D-MPNN/FFN)
- Probabilistic Random Forest

#### Molecular Descriptors:

- Avalon (Fingerprint from Avalon Cheminformatics Toolkit)
- ECFP (Extended Connectivity Fingerprint MorganFingerprint from RDKit)
- ECFP\_counts (ECFP with counts)
- MACCS\_keys (RDKit SMARTS-based implementation of the 166 public MACCS keys)
- RDKit Physchem descriptors
- Special:

7

- Scaled Descriptor (takes another descriptor and scales it on training data)
- Composite Descriptor (concatenates multiple descriptors into one)
- Precomputed Descriptor from file (reads values from CSV file)
- SMILES (also side information related tasks)

Jianping Huang and Xiaohui Fan*		> Brief Bioinform. 2021 Jul 20;22(4):bbaa321. doi: 10.1093/bib/bbaa321.					
View Author Information $^{\sim}$		Do we need different machine learning algorithms					
		16 machine learning algorithms on 14 QSAR data sets					
Research article   Open Access   Published: 14 Augu	st 2017	ing Wu <sup>1</sup> , Minfeng Zhu <sup>2</sup> , Yu Kang <sup>1</sup> , Elaine Lai-Han Leung <sup>3</sup> , Tailong Lei <sup>1</sup> , Chao Shen <sup>1</sup> ,					
Beyond the hype: deep neur established methods using a benchmark set	al networks outperfor a ChEMBL bioactivity	Imag <sup>1</sup> , 2he Wang <sup>1</sup> , Dongsheng Cao <sup>4</sup> , Tingjun Hou <sup>6</sup> Solinformatics. 2018 Jan 1;34(1):72-79. doi: 10.1093/bioinformatics/btx525.     Orthologue chemical space and its influence on     target prediction					
Eelke B. Lenselink, Niels ten Dijke, Brandon Bongers, Wojtek Kowalczyk, Adriaan P. IJzerman & Gerard J. P.	George Papadatos, Herman W. T. van Vlij van Westen ⊠						
Journal of Cheminformatics 9, Article number: 45 (2017)   Cite this article 12k Accesses   162 Citations   33 Altmetric   <u>Metrics</u>		Lewis H Mervin <sup>1</sup> , Krishna C Bulusu <sup>1,2</sup> , Leen Kalash <sup>1</sup> , Avid M Afzal <sup>1</sup> , Fredrik Svensson <sup>1</sup> , Mike A Firth <sup>3</sup> , lan Barrett <sup>3</sup> , Ola Engkvist <sup>4</sup> , Andreas Bender <sup>1</sup>					
> Int J Mol Sci. 2020 Oct 22;21(21):7828. c	loi: 10.3390/ijms21217828.	Affiliations + expand PMID: 28961699 PMCID: PMC5870859 DOI: 10.1093/bioinformatics/btx525 Free PMC article					
Free DMO anticle	DOI: 10.3390/ijms21217828						
Free PMC article	201: 10.3390/ijms21217828						
Free PMC article	201: 10.3390/ijms21217828						
Free PMC article	Research article   Open Acc An automated	ess   <u>Published: 16. January 2018</u> framework for QSAR model building					
Free PMC article RDKit) Jblic MACCS keys)	Research article   Open Acc An automated   Samina Kausar & Andre O. F.	ess   Published: 16 January 2018 framework for QSAR model building alcao ⊡					
Free PMC article RDKit) ublic MACCS keys)	Research article   Open Acc An automated Samina Kausar & Andre O. F. Journal of Cheminformatics 14k Accesses   37 Citation	ess   Published: 16 January 2018 framework for QSAR model building alcao 10, Article number: 1 (2018) Is   22 Altmetric   <u>Metrics</u> Published online 2021 May 26. doi: <u>10.1186/s13321-021-0051</u>					
Free PMC article RDKit) Jublic MACCS keys)	Research article   Open Acc An automated   Samina Kausar & Andre O. Fl Journal of Cheminformatics 14k Accesses   37 Citation	ess   <u>Published: 16 January 2018</u> framework for QSAR model building alcao ⊠ 10, Article number: 1 (2018) Is   22 Altmetric   <u>Metrics</u> Benchmarks for interpretation of QSAR models					
Free PMC article RDKit) ublic MACCS keys)	Research article   Open Acc An automated : Samina Kausar & Andre O. F. Journal of Cheminformatics 14k Accesses   37 Citation	ess   <u>Published: 16 January 2018</u> framework for QSAR model building alcao 10. Article number: 1 (2018) s   22 Altmetric   <u>Metrics</u> Benchmarks for interpretation of QSAR models <u>Maria Matveleva and Pavel Polishchuk<sup>a</sup></u>					
Free PMC article RDKit) ublic MACCS keys) to n training data) to into one)	Research article   Open Acc An automated Samina Kausar & Andre O. F Journal of Cheminformatics 14k Accesses   37 Citation Better Informatics Descriptors	ess Published: 16 January 2018 framework for QSAR model building alcao 10, Article number: 1 (2018) 12 Altmetric   Metrics 10 Deminform, 2021; 13: 41. Published online 2021 May 26. doi: 10.1186/s13321-021-0051! Benchmarks for interpretation of QSAR models Maria Matveieva and Pavel Polishchuk <sup>20</sup>					

### **Probabilistic Threshold Representation (PTR) learning**

- · Model the experimental uncertainty in the input data
- Employ best practice
- Somewhere between regression/classification
- Easy to perform. Data transformed and only threshold and Stdev needed





Mervin LH, Trapotsi MA, Afzal AM, Barrett IP, Bender A, Engkvist O. Probabilistic Random Forest improves bioactivity predictions close to the classification threshold by taking into account experimental uncertainty. *J Cheminform*. 2021;13(1):62.

# The **Qptuna 3-step process**

QPTUNA is structured around three steps:

- **Hyperparameter Optimization**: Train many models with different parameters using Optuna. Only the training dataset is used here. Training is done with cross-validation.
- **Build (Training):** Pick the best model from Optimization. Optionally evaluate its performance on the test dataset.
- "**Prod-build:**" Re-train the best-performing model on the merged training and test datasets. Drawback: no data left to evaluate the resulting model, but large benefit that the final model is trained on all available data







## GUI

- GUI for public release available soon
- Improvements to GUI coming soon
- Automatically suggest configurations as good starting point
- Apply sensible defaults to obtain the best model in reasonable timeframe

Name			
Nume			
Description			
1 Dataset			
🛯 🚯 Main dataset file			
D360 Choose file	No file choser	1	
Input column*			
Response column*			
Deduplication str	ategy		
Keep all			
Data splitting stra	ategy		
No splitting			
Perform probabilistic	representation tr	ansform	
Algorithms			
Descriptors			

# CLI

- Greater flexibility/customizability
- · Less user friendly
- SCP singularity images available
- Local installation possible from bitbucket
- Possible to use run a local CLI/GUI server with MLFlow

```
[4]: # Prepare hyperparameter optimization configuration.
     config = OptimizationConfig(
        data=Dataset(
             input column="canonical", # Typical names are "SMILES" and "smiles".
             response column="molwt". # Often a specific name (like here), or just "activity".
             training_dataset_file="../tests/data/DRD2/subset-50/train.csv",
             test_dataset_file="../tests/data/DRD2/subset-50/test.csv" # Hidden during optimize
         ),
         descriptors=[
             ECFP.new().
             ECFP counts.new(),
             MACCS kevs.new().
         ],
         algorithms=[
             SVR.new(),
             RandomForestRegressor.new(n estimators={"low": 5, "high": 10}),
             Ridge.new(),
             Lasso.new(),
             PLSRegression.new().
             XGBRegressor.new(),
         1.
         settings=OptimizationConfig.Settings(
             mode=ModelMode.REGRESSION,
             cross_validation=3,
             n trials=100. # Total number of trials.
             n_startup_trials=50, # Number of startup ("random") trials.
             direction=OptimizationDirection.MAXIMIZATION,
        ),
```

#### **Run optimization**

```
[5]: # Setup basic logging.
import logging
from importlib import reload
reload(logging)
logging.basicConfig(level=logging.INF0)
```

#### [6]: # Run Optuna Study.

- study = optimize(config, study\_name="my\_study")
- # Optuna will log it's progress to sys.stderr
- # (usually rendered in red in Jupyter Notebooks).

[I 2023-01-06 11:00:56,695] A new study created in memory with name: my\_study



### **MLFIo**

Flow		$\leftarrow \  \   \rightarrow \  \   G$	() localho	ost:5000/#/exp	eriments/4						Q 🕁	U		0 *	+ 8
		mlflo	OW Expe	riments Mode	ls								GitH	lub D	Docs
		Experime	ents +	<b>≺</b> /h	ome/kfwm779/	code/remote	-pycharm/optu	una_az/confiç	gs/examples/o	ptimizat	tion/regre	ssion_	_drd2_5	0.json	1
		Search E	xperiments	Ex	periment ID: 4		Artifact L	ocation: file:///ho	me/kfwm779/code/	remote-pyc	charm/optuna	a_az/mlru	ins/4		
$\leftrightarrow$ $\rightarrow$ C (i) local	lhost:5000/#/compare	e-runs?runs=["1d9	9e2a61e70f4b2f	f8d8d873fd943	:466","570a9a30l	b2a54bc7a728a	ab0 Q 🏠	•	• * •	:					
test_neg_mean_squa	red_log_error -0.005	-0.034	-0.028	-0.01	-0.003	-0.003	-0.003	-0.003	-0.C	•					
test_neg_median_abs	solute_error -22.65	-54.48	-55.87	-37.19	-17.41	-16.65	-17.57	-16.65	-16.	в:	Active -	s	learch	Clea	ear
test_r2 🗠	0.88	0.189	0.33	0.781	0.922	0.921	0.922	0.921	0.92			≡ 8	<b>a</b>	Column	ns
4									۱.			Tags >			
Scatter Plot	Contour Plot Pa	rallel Coordinates Plo	t							test	_exp test_ma	o number	r datetim	datetim	ne
X-axis:										0.88	2 -55	99	202	202	<b>^</b>
trial_number	~									0.27	5 -17	98	202	202	
V avic:		1								0.34	3 -11	97	202	202	
ontimization object	ctive cymean r2	2		000,000 0000						0.78	i2 -72	96	202	202	
opumization_obje		ean	••							0.92	.4 -39 3 -42.6	95	202	202	
		Č 0.5								0.92	4 -39	93	202	202	
		Ke	•• •	•		•	•			0.92	3 -42	92	202	202	
		ject								0.92	4 -41	91	202	202	
			•		•	•	•			1.7					
		atio			-										
		imiz				•		•							
		td -0.5													
					•										
		0	20	)	40	60	80	100							
					trial number	r									
										-					

:

# **Optuna settings**

- Num trials
- Num of random trials

#### Pruning now implemented

Trials



### **Real CLS example**



# **Build (Training): Split for evaluation**





#### **PTR example**





# **Additional features: Uncertainty (Bias)**

#### VennABERS (Classification)

Large-scale probabilistic prediction with and without validity guarantees

Vladimir Vovk, Ivan Petei, and Valentina Fedorova

#### **Venn-ABERS Predictor**

(Preliminary documentation)

The VennABERS.pv file is a pure Python implementation of the fast Venn-ABERS Predictor described in Vovk2019

A Venn-ABERS predictor outputs two probability predictions for every test object. In particular, the Venn-ABERS predictor implemented here is the inductive form of probability predictor, which relies on a calibration set. In a nutshell, the Venn-ABERS predictor can be viewed as a distribution-free calibration function that maps scores output by a scoring classifier to well-calibrated probabilities. A gentle introduction can be found in this tutorial.

The function that implements the Venn-ABERS Predictor is ScoresToMultiProbs()

#### p0,p1 = ScoresToMultiProbs(calibrPts,testScores)

> J Chem Inf Model. 2020 Oct 26;60(10):4546-4559. doi: 10.1021/acs.jcim.0c00476. Epub 2020 Sep 21.

Comparison of Scaling Methods to Obtain Calibrated Probabilities of Activity for Protein-Ligand Predictions

ROYAL

HOLLOWAY

Lewis H Mervin<sup>1</sup>, Avid M Afzal<sup>2</sup>, Ola Engkvist<sup>3</sup>, Andreas Bender<sup>4</sup>

PMID: 32865408 DOI: 10.1021/acs.icim.0c00476

#### **MAPIE** (Regression)



#### **MAPIE - Model Agnostic Prediction Interval Estimator**

codecov 100%

C Unit tests

MAPIE allows you to easily estimate prediction intervals (or prediction sets) using your favourite scikit-learncompatible model for single-output regression or multi-class classification settings

Prediction intervals output by MAPIE encompass both aleatoric and epistemic uncertainties and are backed by strong theoretical guarantees thanks to conformal prediction methods [1-7].







## **Additional features: Uncertainty (variance)**

When the data, representation, architecture, etc. is kept constant.

Only variance is the initialisation

#### ChemProp (task type agnostic)

#### Uncertainty Estimation

The uncertainty of predictions made in Chemprop can be estimated by several different methods. Uncertainty estimation is carried out alongside model value prediction and reported in the predictions csv file when the argument --uncertainty\_method <method> is provided. If no uncertainty method is provided, then only the model value predictions will be carried out. The available methods are:

ensemble For a prediction using an ensemble of models. Returns the variance of predictions made by each
of the ensemble submodels. Ensemble variance can be used with any dataset type, but the results are only
usable for calibration or evaluation with regression datasets.



## **Additional features: Explainability**



game theory and their related extensions (see papers for details and citations)

- Auto select appropriate explainer
   (Tree/Linear/Permutation) for compatible algorithms
- Model agnostic example with KernelExplainer to explain unsupported algorithms
- Output: importance for features in a chemical series/ input
- 20 ECFP's explained for bits turned on

#### ChemProp

#### Interpreting

It is often helpful to provide explanation of model prediction (i.e., this molecule is toxic because of this substructure). Given a trained model, you can interpret the model prediction using the following command:

chemprop\_interpret --data\_path data/tox21.csv --checkpoint\_dir tox21\_checkpoints/fold\_0/ --property

If installed from source, chemprop\_interpret can be replaced with python interpret.py .

The output will be like the following:

- The first column is a molecule and second column is its predicted property (in this case NR-AR toxicity).
- The third column is the smallest substructure that made this molecule classified as toxic (which we call rationale).
- The fourth column is the predicted toxicity of that substructure.

As shown in the first row, when a molecule is predicted to be non-toxic, we will not provide any rationale for its prediction.

smiles	NR- AR	rationale	rationale_score
[C(F)(F)F)cc([N+](=O)[O-])c1Cl	0.014		
]2C[C@H]3[C@@H]4C[C@H] =C[C@]5(C)[C@H]4[C@@H] }(C)[C@]2(C(=O)CO)O1	0.896	C[C@]12C=CC(=0)C=C1[CH2:1]C[CH2:1] [CH2:1]2	0.769
2CC[C@H]3[C@@H] }@H]4C(O)=C(C#N)C[C@]35C) }H]1CC[C@@H]2O	0.941	C[C@]12C[CH:1]=[CH:1] [C@H]3O[C@]31CC[C@@H]1[C@@H]2CC[C:1] [CH2:1]1	0.808
≿@H](O)[C@H]3[C@@H] ≿C(=O)CC[C@@]43C)  2(O)C(=O)COP(=O)[(O-])[O-]	0.957	C1C[CH2:1][C:1] [C@@H]2[C@@H]1[C@@H]1CC[C:1] [C:1]1C[CH2:1]2	0.532

## **Qptuna to enable AL integration within the DMTA cycle**



Janet JP, Mervin L, Engkvist O "Artificial Intelligence in Molecular de novo Design - Integration with Experiment" Current Opinion in Structural Biology (accepted) Fill up plates with compounds from PLS that are predicted to improve the model (uncertainty/performance)

> Suggest compounds using AL



#### Live demo



### **Future Plans**

- GUI improvements
- Documentation improvements
- Active learning



#### **Confidentiality Notice**

This file is private and may contain confidential and proprietary information. If you have received this file in error, please notify us and remove it from your system and note that you must not copy, distribute or take any action in reliance on it. Any unauthorized use or disclosure of the contents of this file is not permitted and may be unlawful. AstraZeneca PLC, 1 Francis Crick Avenue, Cambridge Biomedical Campus, Cambridge, CB2 0AA, UK, T: +44(0)203 749 5000, www.astrazeneca.com

