

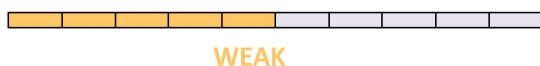


ROBERT v 1.0.4 2023/10/01 14:29:42

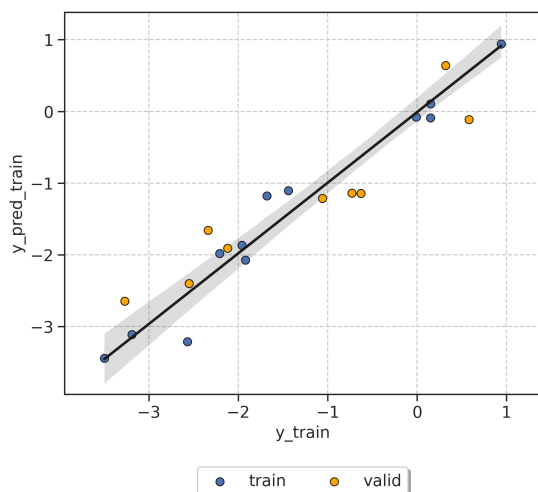
**How to cite:** ROBERT v 1.0.4, Dalmau, D.; Alegre-Requena, J. V., 2023. <https://github.com/jvalegre/robert>**ROBERT SCORE***This score is designed to analyze the predictive ability of the models using different metrics.***No PFI (all descriptors):**

ML model: NN

Proportion Train:Validation = 57:43

**The model has a score of 5/10**

- The valid. set shows an  $R^2$  of 0.89
- The valid. set has 33.3% of outliers
- Using 21:91 points(train+valid.):descriptors
- The valid. set passes 3 VERIFY tests



Train :  $R^2 = 0.96$ , MAE = 0.2, RMSE = 0.28  
 Valid. :  $R^2 = 0.89$ , MAE = 0.42, RMSE = 0.47

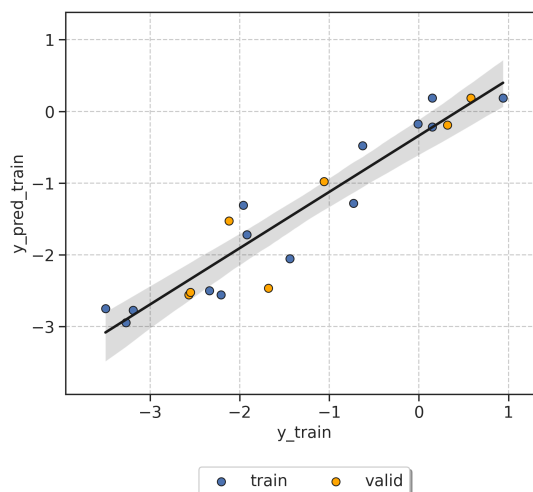
**PFI (only important descriptors):**

ML model: RF

Proportion Train:Validation = 67:33

**The model has a score of 8/10**

- The valid. set shows an  $R^2$  of 0.88
- The valid. set has 0.0% of outliers
- Using 21:4 points(train+valid.):descriptors
- The valid. set passes 3 VERIFY tests



Train :  $R^2 = 0.91$ , MAE = 0.39, RMSE = 0.45  
 Valid. :  $R^2 = 0.88$ , MAE = 0.34, RMSE = 0.45

**Score thresholds** (detailed in <https://robert.readthedocs.io/en/latest/Score/score.html>) **$R^2$**  \_\_\_\_\_

- $R^2 > 0.85$
- $0.85 > R^2 > 0.70$
- $R^2 < 0.70$

**Outliers** \_\_\_\_\_

- < 7.5% of outliers
- 7.5% < outliers < 15%
- > 15% of outliers

**Points:descriptors** \_\_\_\_\_

- > 10:1 p:d ratio
- 10:1 > p:d ratio > 3:1
- p:d ratio < 3:1

**VERIFY tests** \_\_\_\_\_

- Up to ●●●● (tests pass)
- (all tests failed)

### Some tips to improve the score

- ⚠ The model uses only 21 datapoints, adding meaningful datapoints might help to improve the model.
- ⚠ One of your models have more than 7.5% of outliers (5% is expected for a normal distribution with the t-value of 2 that ROBERT uses), using a more homogeneous distribution of results might help. For example, avoid using many points with similar y values and only a few points with distant y values.
- ⚠ Adding meaningful descriptors or replacing/deleting the least useful descriptors used might help. Feature importances are gathered in the SHAP and PFI sections of the /PREDICT/PREDICT\_data.dat file.

### How to predict new values with these models?

1. Create a CSV database with the new points, including the necessary descriptors.
2. Place the CSV file in the parent folder (i.e., where the module folders were created)
3. Run the PREDICT module as 'python -m robert --predict --csv\_test FILENAME.csv'.
4. The predictions will be stored in the last column of two CSV files called MODEL\_SIZE\_test(\_No)\_PFI.csv, which are stored in the PREDICT folder.



## REPRODUCIBILITY

*This section provides all the instructions to reproduce the results presented.*

### **1. Download these files (*the authors should have uploaded the files as supporting information!*):**

- Report with results (ROBERT\_report.pdf)
- CSV database (solubility\_short.csv)

### **2. Install and adjust the versions of the following Python modules:**

- Install ROBERT and its dependencies: `conda install -c conda-forge robert`
- Adjust ROBERT version: `pip install robert==1.0.4`
- Install scikit-learn-intelex: `pip install scikit-learn-intelex==2023.2.1`

*(if scikit-learn-intelex is not installed, slightly different results might be obtained)*

- Install AQME and its dependencies: `conda install -c conda-forge aqme`
- Adjust AQME version: `pip install aqme==1.5.1`
- Install xTB: `conda install -c conda-forge xtb`
- Adjust xTB version (if possible): `conda install -c conda-forge xtb=6.6.1`

### **3. Run ROBERT using this command line in the folder with the CSV database:**

```
python -m robert --aqme --y "solubility" --csv_name "solubility_short.csv"
```

### **4. Execution time, Python version and OS:**

Originally run in Python 3.10.12 using Linux #1 SMP Tue Dec 21 19:02:23 UTC 2021

Total execution time: 86.86 seconds (the number of processors should be specified by the user)



## TRANSPARENCY

*This section contains important parameters used in scikit-learn models and ROBERT.*

### 1. Parameters of the scikit-learn models (same keywords as used in scikit-learn):

#### No PFI (all descriptors):

sklearn model: MLPRegressor  
 random\_state: 19  
 names: code\_name  
 batch\_size: 4  
 hidden\_layer\_sizes: [16, 16]  
 learning\_rate\_init: 0.01  
 max\_iter: 200  
 validation\_fraction: 0.1  
 alpha: 0.0001  
 shuffle: True  
 tol: 0.0001  
 early\_stopping: False  
 beta\_1: 0.999  
 beta\_2: 0.999  
 epsilon: 1e-08

#### PFI (only important descriptors):

sklearn model: RandomForestRegressor  
 random\_state: 43  
 names: code\_name  
 n\_estimators: 100  
 max\_depth: 5  
 max\_features: 0.5  
 min\_samples\_split: 2  
 min\_samples\_leaf: 1  
 min\_weight\_fraction\_leaf: 0  
 ccp\_alpha: 0  
 oob\_score: False  
 max\_samples: 0.75

### 2. ROBERT options for data split (KN or RND), predict type (REG or CLAS) and hyperopt error (RMSE, etc.):

#### No PFI (all descriptors):

split: KN  
 type: reg  
 error\_type: rmse

#### PFI (only important descriptors):

split: KN  
 type: reg  
 error\_type: rmse



## ABBREVIATIONS

*Reference section for the abbreviations used.*

**ACC:** accuracy

**ADAB:** AdaBoost

**CSV:** comma separated values

**CLAS:** classification

**CV:** cross-validation

**F1 score:** balanced F-score

**GB:** gradient boosting

**GP:** gaussian process

**KN:** k-nearest neighbors

**MAE:** root-mean-square error

**MCC:** Matthew's correl. coefficient

**ML:** machine learning

**MVL:** multivariate lineal models

**NN:** neural network

**PFI:** permutation feature importance

**R2:** coefficient of determination

**REG:** Regression

**RF:** random forest

**RMSE:** root mean square error

**RND:** random

**SHAP:** Shapley additive explanations

**VR:** voting regressor

**AQME**

*This module performs RDKit conformer generation from SMILES, followed by the creation of 200+ molecular and atomic descriptors using RDKit, xTB and DBSTEP (saved as AQME-ROBERT\_FILENAME.csv).*

The complete output (AQME\_data.dat) and raw data are stored in the AQME folder.

Time AQME: 25.75 seconds

---

**CURATE**

*This module takes care of data curation, including filters for correlated descriptors, noise, and duplicates, as well as conversion of categorical descriptors.*

The complete output (CURATE\_data.dat) and curated database are stored in the CURATE folder.

Time CURATE: 1.36 seconds

----- Images generated by the CURATE module -----

---

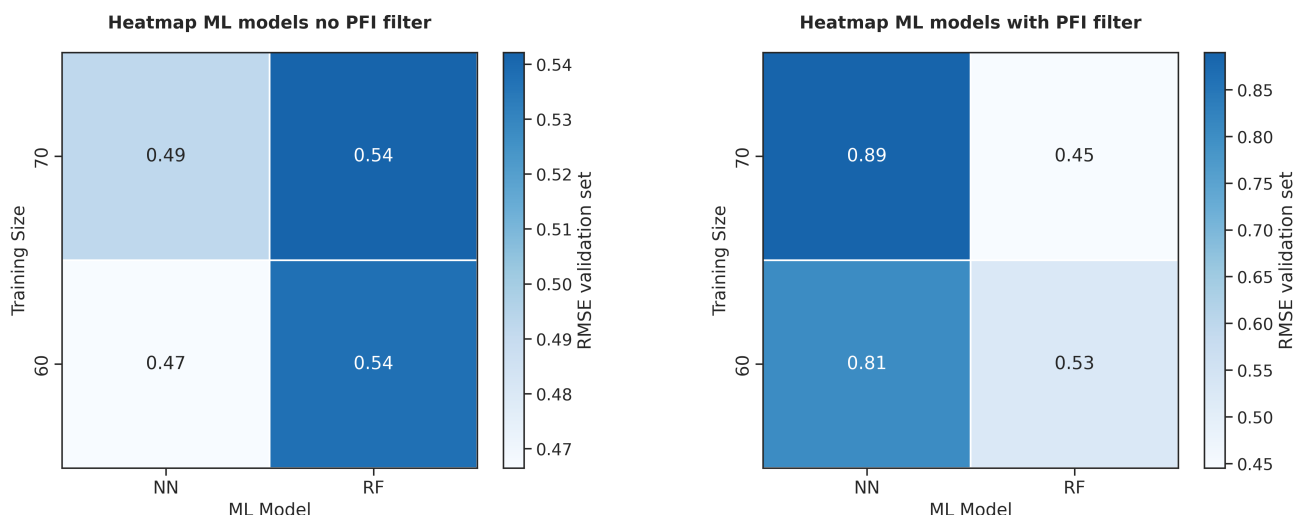
**GENERATE**

*This module carries out a screening of ML models and selects the most accurate one. It includes a comparison of multiple hyperoptimized models and training sizes.*

The complete output (GENERATE\_data.dat) and heatmaps are stored in the GENERATE folder.

Time GENERATE: 50.14 seconds

----- Images generated by the GENERATE module -----



## VERIFY

Determination of predictive ability of models using four tests: 5-fold CV, y-mean (error against the mean y baseline), y-shuffle (predict with shuffled y values), and one-hot (predict using one-hot encoding instead of the X values).

The complete output (VERIFY\_data.dat) and donut plot are stored in the VERIFY folder.

Time VERIFY: 1.41 seconds

### ----- Images and summary generated by the VERIFY module -----

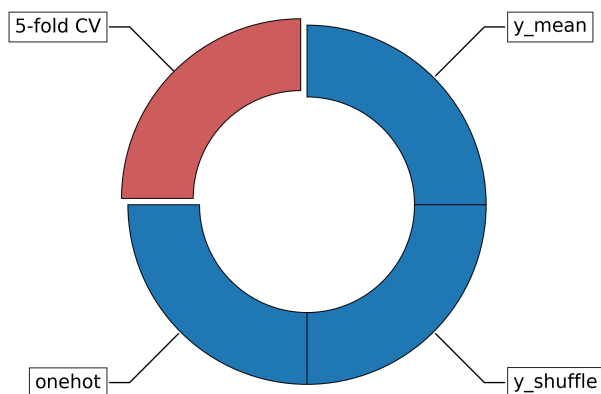
#### No PFI (all descriptors):

Original RMSE (valid. set)  $0.47 + 25\% \text{ thres.} = 0.58$   
 x 5-fold CV: FAILED, RMSE = 0.63, higher than thres.  
 o y\_mean: PASSED, RMSE = 1.3, higher than thres.  
 o y\_shuffle: PASSED, RMSE = 1.7, higher than thres.  
 o onehot: PASSED, RMSE = 0.69, higher than thres.

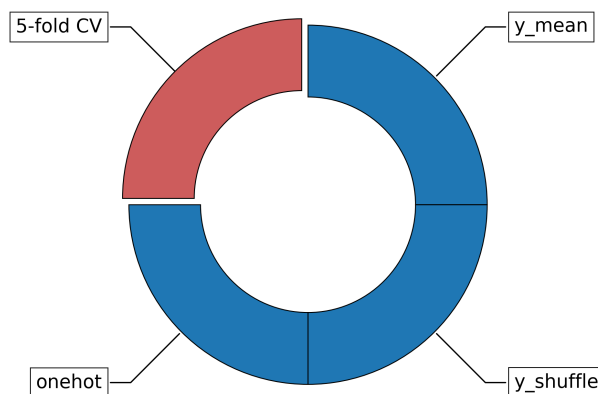
#### PFI (only important descriptors):

Original RMSE (valid. set)  $0.45 + 25\% \text{ thres.} = 0.56$   
 x 5-fold CV: FAILED, RMSE = 0.79, higher than thres.  
 o y\_mean: PASSED, RMSE = 1.2, higher than thres.  
 o y\_shuffle: PASSED, RMSE = 1.3, higher than thres.  
 o onehot: PASSED, RMSE = 0.98, higher than thres.

VERIFY tests of NN\_60\_No\_PFI



VERIFY tests of RF\_70\_PFI



**PREDICT**

This module predicts and plots the results of training and validation sets from GENERATE, as well as from external test sets (if any). Feature importances from SHAP and PFI, and outlier analysis are also represented.

The complete output (PREDICT\_data.dat) and heatmaps are stored in the PREDICT folder.

Time PREDICT: 8.2 seconds

----- Images and summary generated by the PREDICT module -----

**No PFI (all descriptors):**Prediction metrics and descriptors

- Points Train:Validation = 12:9
- Proportion Train:Validation = 57:43
- Number of descriptors = 91
- Proportion (train+valid.) points:descriptors = 21:91
- Train :  $R^2 = 0.96$ , MAE = 0.2, RMSE = 0.28
- Valid. :  $R^2 = 0.89$ , MAE = 0.42, RMSE = 0.47

Outliers (max. 10 shown)

Train: 1 outliers out of 12 datapoints (8.3%)

- mol\_1069 (2.3 SDs)

Validation: 3 outliers out of 9 datapoints (33.3%)

- mol\_100 (2.2 SDs)
- mol\_106 (2.6 SDs)
- mol\_1076 (2.5 SDs)

**PFI (only important descriptors):**Prediction metrics and descriptors

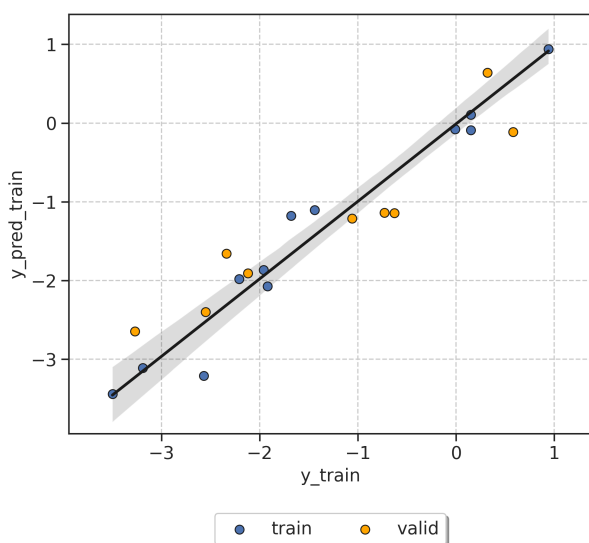
- Points Train:Validation = 14:7
- Proportion Train:Validation = 67:33
- Number of descriptors = 4
- Proportion (train+valid.) points:descriptors = 21:4
- Train :  $R^2 = 0.91$ , MAE = 0.39, RMSE = 0.45
- Valid. :  $R^2 = 0.88$ , MAE = 0.34, RMSE = 0.45

Outliers (max. 10 shown)

Train: 0 outliers out of 14 datapoints (0.0%)

Validation: 0 outliers out of 7 datapoints (0.0%)

Predictions\_train\_valid of NN\_60\_No\_PFI



Predictions\_train\_valid of RF\_70\_PFI

