

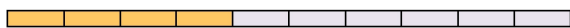
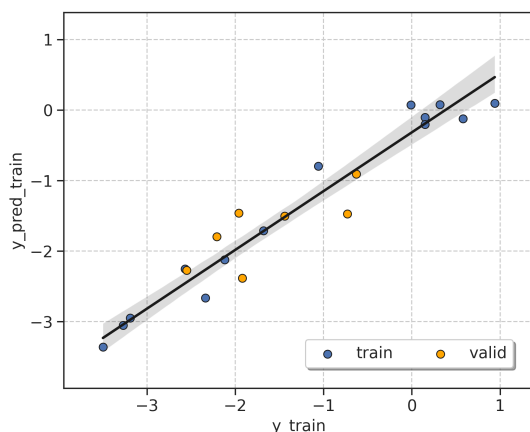


ROBERT v 1.2.0 2024/09/13 20:25:04

How to cite: Dalmau, D.; Alegre Requena, J. V. ChemRxiv, 2023, DOI: 10.26434/chemrxiv-2023-k994h**Section A. ROBERT Score***This score is designed to evaluate the models using different metrics.***No PFI (standard descriptor filter):**

Model = NN · Train:Validation = 67:33

Points(train+valid.):descriptors = 21:6

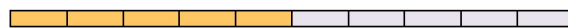
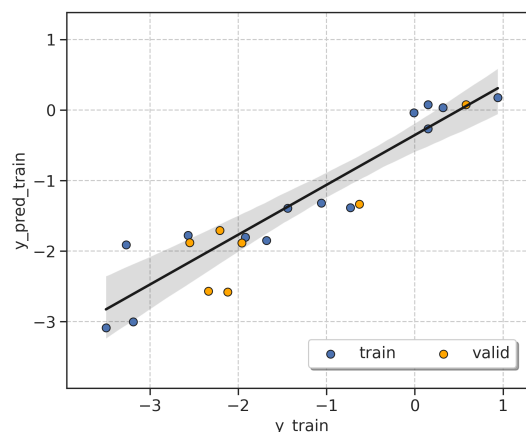
**WEAK**

Train : $R^2 = 0.97$, MAE = 0.29, RMSE = 0.37
Valid. : $R^2 = 0.59$, MAE = 0.39, RMSE = 0.44

PFI (only most important descriptors):

Model = RF · Train:Validation = 67:33

Points(train+valid.):descriptors = 21:2

**WEAK**

Train : $R^2 = 0.9$, MAE = 0.4, RMSE = 0.54
Valid. : $R^2 = 0.8$, MAE = 0.45, RMSE = 0.5

Severe warnings☒ No severe warnings detected**Moderate warnings**

- ☒ Imprecise predictions (Section B.3b)
- ☒ Slightly uneven y distribution (Section C)

Overall assessment☒ The model is unreliable**Severe warnings**☒ No severe warnings detected**Moderate warnings**

- ☒ Imprecise predictions (Section B.3b)
- ☒ Slightly uneven y distribution (Section C)

Overall assessment☒ The model is unreliable



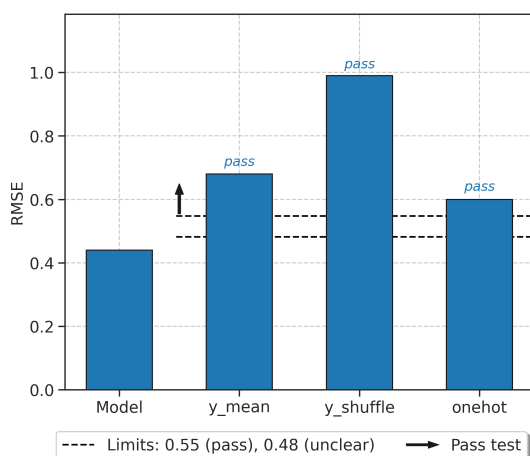
Section B. Advanced Score Analysis

This section explains each component that comprises the ROBERT score.

1. Model vs "flawed" models (3 / 3)

The model predicts right for the right reasons.

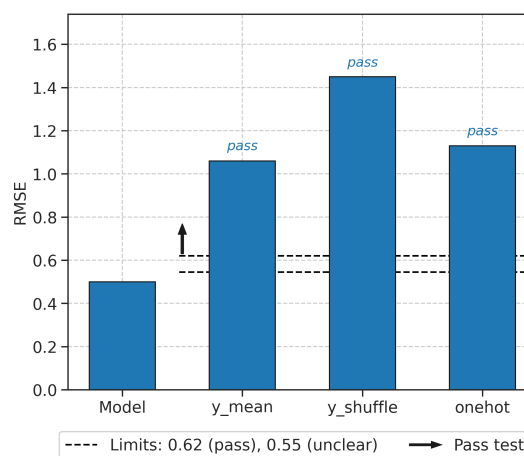
Pass: +1, Unclear: 0, Fail: -1. [Details here.](#)



1. Model vs "flawed" models (3 / 3)

The model predicts right for the right reasons.

Pass: +1, Unclear: 0, Fail: -1. [Details here.](#)



2. Predictive ability of the model (0 / 2)

Low predictive ability with R^2 (valid.) = 0.59.

R^2 0.70-0.85: +1, R^2 >0.85: +2.

2. Predictive ability of the model (1 / 2)

Moderate predictive ability with R^2 (valid.) = 0.8.

R^2 0.70-0.85: +1, R^2 >0.85: +2.

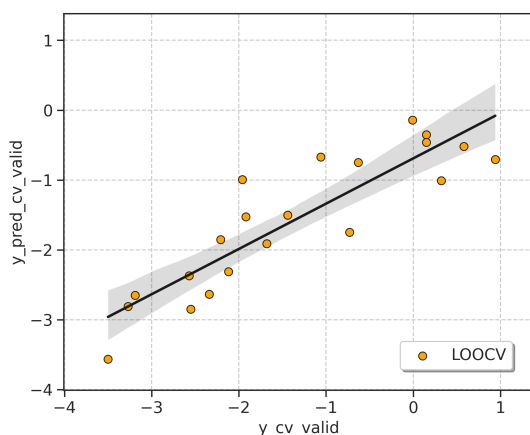
3. Cross-validation (LOOCV) of the model

Overfitting analysis on the model with 3a and 3b:

3a. CV predictions train + valid. (1 / 2)

Moderate predictive ability with R^2 (LOOCV) = 0.79.

R^2 0.70-0.85: +1, R^2 >0.85: +2.



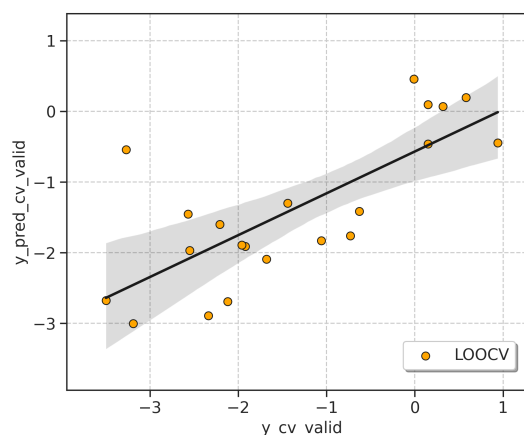
3. Cross-validation (LOOCV) of the model

Overfitting analysis on the model with 3a and 3b:

3a. CV predictions train + valid. (0 / 2)

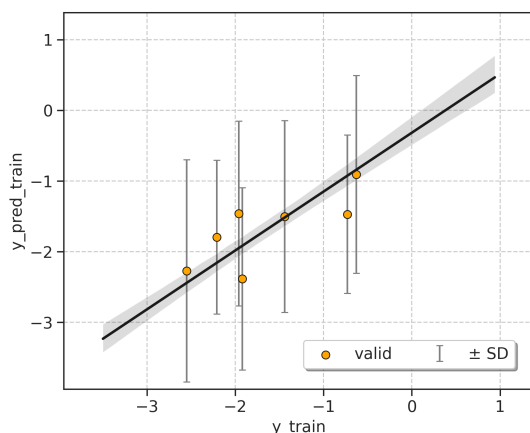
Low predictive ability with R^2 (LOOCV) = 0.57.

R^2 0.70-0.85: +1, R^2 >0.85: +2.

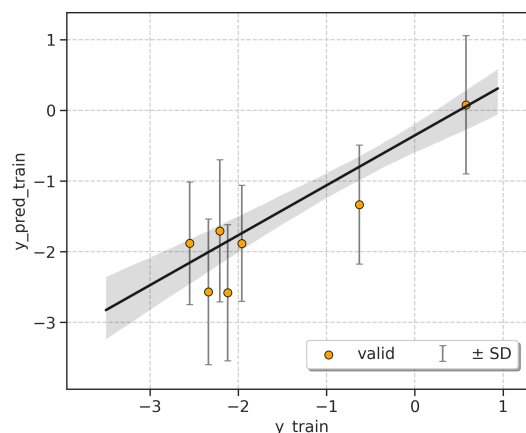


3b. Avg. standard deviation (SD) (0 / 2 ☐)

High variation, $4 \times \text{SD (valid.)} = 5.2$ (118% y-range).
 $4 \times \text{SD 25-50\% y-range: } +1$, $4 \times \text{SD} < 25\% \text{ y-range: } +2$.
[Details here.](#)

**3b. Avg. standard deviation (SD)** (0 / 2 ☐)

High variation, $4 \times \text{SD (valid.)} = 3.7$ (84% y-range).
 $4 \times \text{SD 25-50\% y-range: } +1$, $4 \times \text{SD} < 25\% \text{ y-range: } +2$.
[Details here.](#)

**4. Points(train+valid.):descriptors** (0 / 1 ☐)

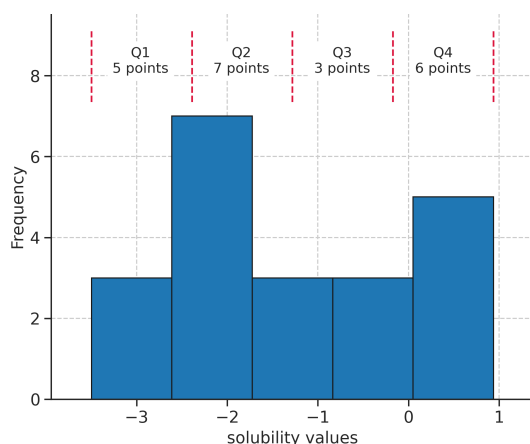
Number of descps. could be lower (ratio 21:6).
 5 or more points per descriptor: +1.

4. Points(train+valid.):descriptors (1 / 1 ☒)

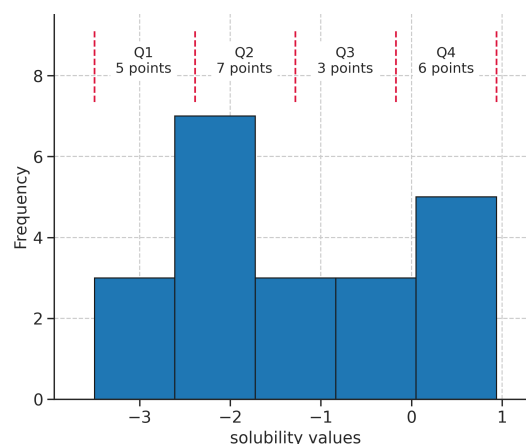
Decent number of descps. (ratio 21:2).
 5 or more points per descriptor: +1.

**Section C. Distribution of y Values**

This section shows the distribution of y values within the training and validation sets.

**y distribution analysis**

x **WARNING!** Your data is slightly not uniform (Q3 has 3 points while Q2 has 7)

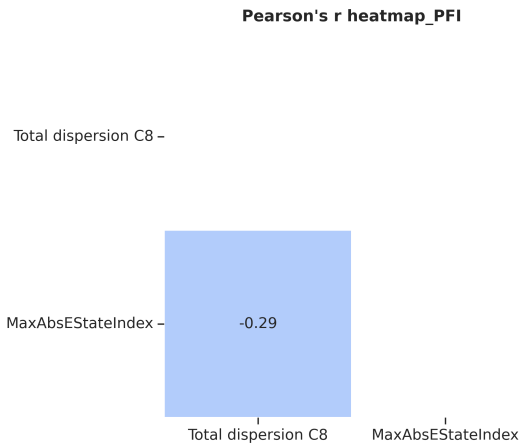
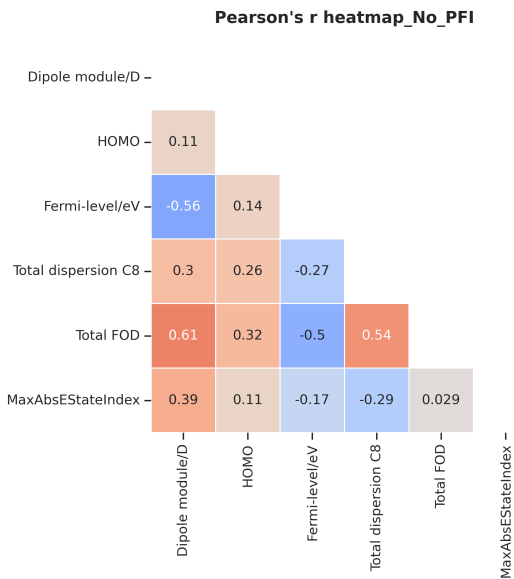
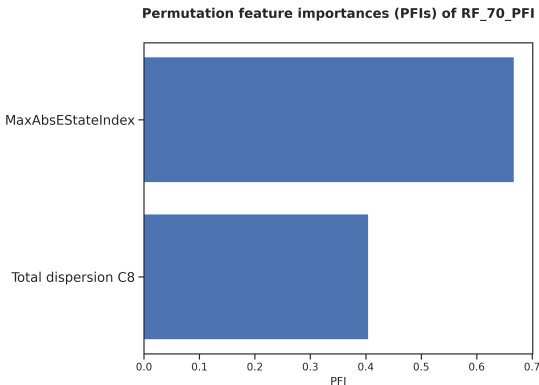
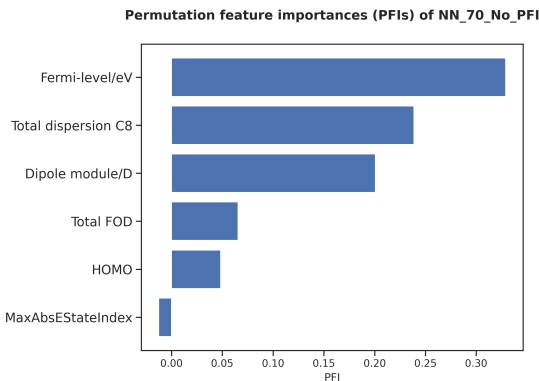
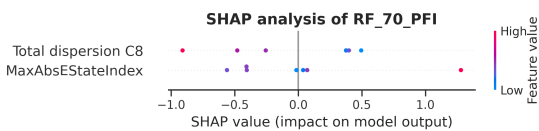
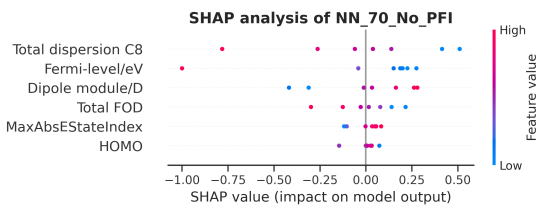
**y distribution analysis**

x **WARNING!** Your data is slightly not uniform (Q3 has 3 points while Q2 has 7)



Section D. Feature Importances

This section presents feature importances measured using the validation set.



Correlation analysis

o Correlations between variables are acceptable

Correlation analysis

o Correlations between variables are acceptable



Section E. Outlier Analysis

This section detects outliers using the standard deviation (SD) of errors from the training set.

No PFI (standard descriptor filter):

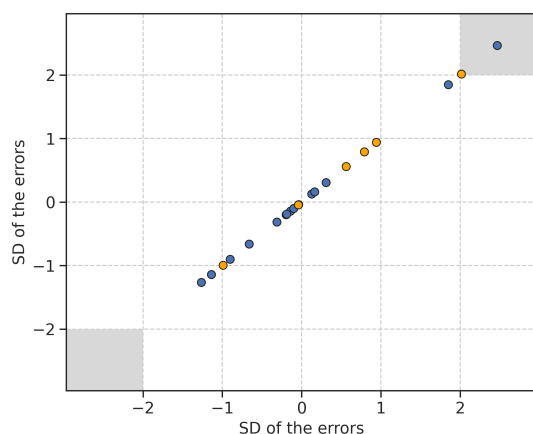
Outliers (max. 10 shown)

Train: 1 outliers out of 14 datapoints (7.1%)

- mol_1083 (2.5 SDs)

Validation: 1 outliers out of 7 datapoints (14.3%)

- mol_1000 (2.0 SDs)



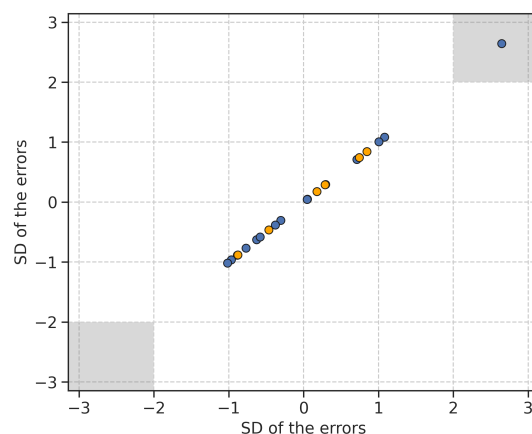
PFI (only most important descriptors):

Outliers (max. 10 shown)

Train: 1 outliers out of 14 datapoints (7.1%)

- mol_100 (2.6 SDs)

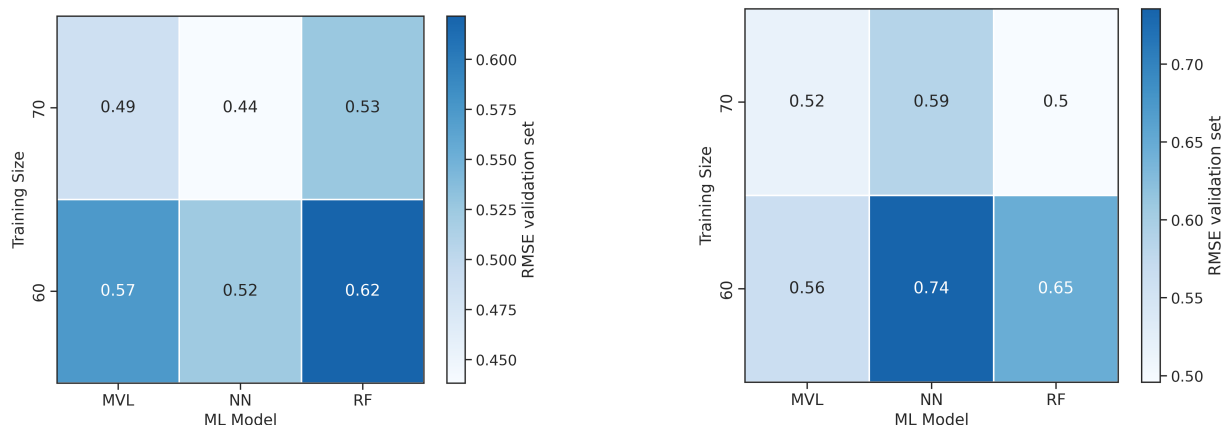
Validation: 0 outliers out of 7 datapoints (0.0%)





Section F. Model Screening

This section compares different combinations of hyperoptimized algorithms and partition sizes.



Section G. 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!*):

- 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==1.2.0`
- Install scikit-learn-intelex: `pip install scikit-learn-intelex==2024.5.0`

(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==1.6.1`
- 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 3.10.13 using Linux #1 SMP Fri Mar 29 23:14:13 UTC 2024

Total execution time: 78.74 seconds (*the number of processors should be specified by the user*)



Section H. 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 (standard descriptor filter):

sklearn model: MLPRegressor
 random_state: 233
 names: code_name
 batch_size: 4
 hidden_layer_sizes: [16, 16]
 learning_rate_init: 0.01
 max_iter: 200
 validation_fraction: 0.2
 alpha: 0.0001
 shuffle: True
 tol: 0.0001
 early_stopping: False
 beta_1: 0.999
 beta_2: 0.999
 epsilon: 1e-08

PFI (only most important descriptors):

sklearn model: RandomForestRegressor
 random_state: 43
 names: code_name
 n_estimators: 60
 max_depth: 5
 max_features: 1.0
 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 (standard descriptor filter):

split: KN
 type: reg
 error_type: rmse

PFI (only most important descriptors):

split: KN
 type: reg
 error_type: rmse



Section I. Abbreviations

Reference section for the abbreviations used.

| | | |
|------------------------------------|--|--|
| ACC: accuracy | KN: k-nearest neighbors | REG: Regression |
| ADAB: AdaBoost | MAE: root-mean-square error | RF: random forest |
| CSV: comma separated values | MCC: Matthew's correl. coefficient | RMSE: root mean square error |
| CLAS: classification | ML: machine learning | RND: random |
| CV: cross-validation | MVL: multivariate lineal models | SHAP: Shapley additive explanations |
| F1 score: balanced F-score | NN: neural network | VR: voting regressor |
| GB: gradient boosting | PFI: permutation feature importance | |
| GP: gaussian process | R2: coefficient of determination | |

Miscellaneous

General tips to improve the models and instructions to predict new values.

Some general tips to improve the score

1. Adding meaningful datapoints might help to improve the model. Also, using a uniform population of datapoints across the whole range of y values usually helps to obtain reliable predictions across the whole range. More information about the range of y values used is available in Section C.
2. Adding meaningful descriptors or replacing/deleting the least useful descriptors used might help. Feature importances are gathered in Section D.

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 shown at the end of the resulting PDF report and will be stored in the last column of two CSV files called MODEL_SIZE_test(_No)_PFI.csv, which are in the PREDICT folder.
-