

ROBERT v 1.2.1 2025/08/05 18:56:41

How to cite: Dalmau, D.; Alegre Requena, J. V. WIREs Comput Mol Sci. 2024, DOI: 10.1002/WCMS.1733



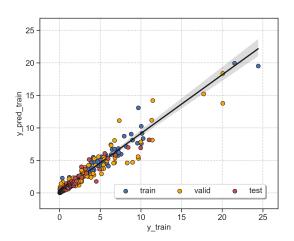
# Section A. ROBERT Score

This score is designed to evaluate the models using different metrics.

### No PFI (standard descriptor filter):

Model = RF · Train:Validation:Test = 54:36:10 Points(train+valid.):descriptors = 1008:6 Score = 10 / 10

## **STRONG**



Train:  $R^2 = 0.97$ , MAE = 0.15, RMSE = 0.39 Valid. :  $R^2 = 0.92$ , MAE = 0.31, RMSE = 0.74 Test:  $R^2 = 0.92$ , MAE = 0.27, RMSE = 0.59

## Severe warnings

No severe warnings detected

#### **Moderate warnings**

- Uneven y distribution (Section C)
- Potential "faulty" outliers (Section E)

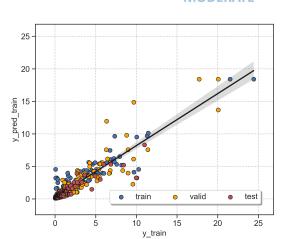
#### Overall assessment

The model seems reliable

## PFI (only most important descriptors):

Model = RF · Train:Validation:Test = 54:36:10 Points(train+valid.):descriptors = 1008:4 Score = 8 / 10

#### **MODERATE**



Train:  $R^2 = 0.89$ , MAE = 0.34, RMSE = 0.7 Valid.:  $R^2 = 0.77$ , MAE = 0.47, RMSE = 1.2 Test:  $R^2 = 0.84$ , MAE = 0.39, RMSE = 0.92

#### Severe warnings

No severe warnings detected

#### **Moderate warnings**

- Uneven y distribution (Section C)
- Potential "faulty" outliers (Section E)

#### **Overall assessment**

Decent model, but it has limitations

ROBERT v 1.2.1 Page 1 of 8

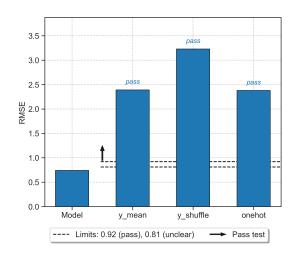


## Section B. Advanced Score Analysis

This section explains each component that comprises the ROBERT score.

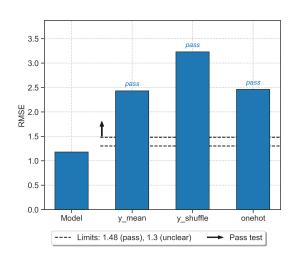
#### 1. Model vs "flawed" models (3 / 3

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 (2 / 2

Good predictive ability with  $R^2$  (test) = 0.92.  $R^2$  0.70-0.85: +1,  $R^2$  >0.85: +2.

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

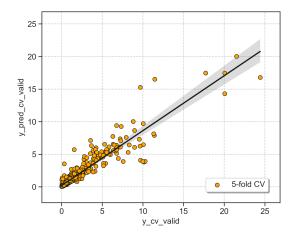
Moderate predictive ability with  $R^2$  (test) = 0.84.  $R^2$  0.70-0.85: +1,  $R^2$  >0.85: +2.

#### 3. Cross-validation (5-fold CV) of the model

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

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

Good predictive ability with  $R^2$  (5-fold CV) = 0.89.  $R^2$  0.70-0.85: +1,  $R^2$  >0.85: +2.

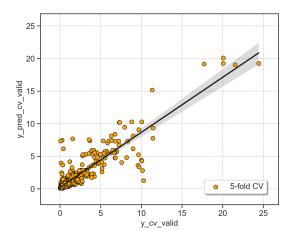


#### 3. Cross-validation (5-fold CV) of the model

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

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

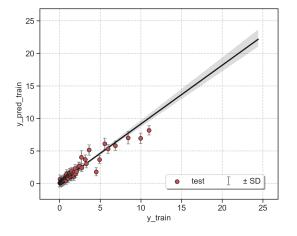
Moderate predictive ability with  $R^2$  (5-fold CV) = 0.82.  $R^2$  0.70-0.85: +1,  $R^2$  >0.85: +2.



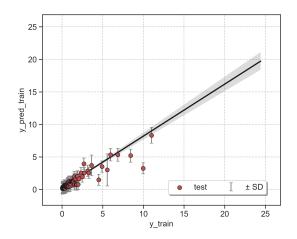
ROBERT v 1.2.1 Page 2 of 8

3b. Avg. standard deviation (SD) (2 / 2 Low variation, 4\*SD (test) = 2.6 (10% y-range). 4\*SD 25-50% y-range: +1, 4\*SD < 25% y-range: +2.

Details here.



3b. Avg. standard deviation (SD) (2 / 2 Low variation, 4\*SD (test) = 3.4 (14% y-range). 4\*SD 25-50% y-range: +1, 4\*SD < 25% y-range: +2. Details here.



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

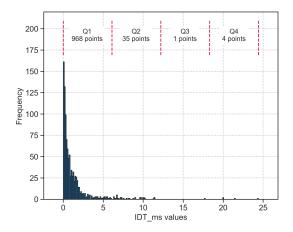
Decent number of descps. (ratio 1008:6). 5 or more points per descriptor: +1.

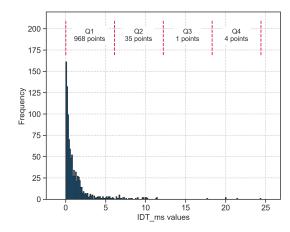
4. Points(train+valid.):descriptors (1 / 1 ===)

Decent number of descps. (ratio 1008:4). 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 not uniform (Q3 has 1 points while Q1 has 968)

y distribution analysis

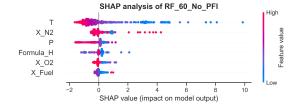
x WARNING! Your data is not uniform (Q3 has 1 points while Q1 has 968)

ROBERT v 1.2.1 Page 3 of 8



## Section D. Feature Importances

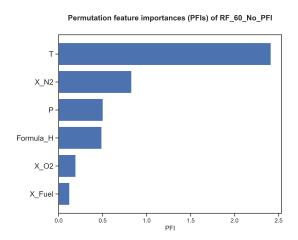
This section presents feature importances measured using the validation set.

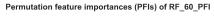


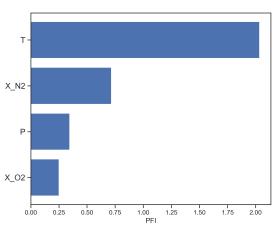
SHAP analysis of RF\_60\_PFI

T
P
X\_N2
X\_O2

-2
0
2
4
6
8
SHAP value (impact on model output)

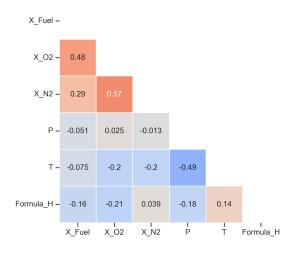


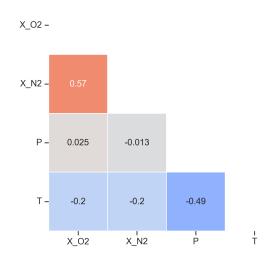




Pearson's r heatmap\_No\_PFI







#### **Correlation analysis**

o Correlations between variables are acceptable

## **Correlation analysis**

o Correlations between variables are acceptable

ROBERT v 1.2.1 Page 4 of 8



### 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: 20 outliers out of 604 datapoints (3.3%)

- 320 (4.9 SDs)
- 322 (2.3 SDs)
- 323 (4.8 SDs)
- 325 (6.4 SDs)
- 361 (4.0 SDs)
- 403 (3.6 SDs)
- 403 (3.0 303)
- 423 (2.0 SDs)
- 426 (2.4 SDs)442 (3.3 SDs)
- 586 (3.0 SDs)

Validation: 36 outliers out of 404 datapoints (8.9%)

- 11 (2.8 SDs)
- 203 (2.4 SDs)
- 204 (3.1 SDs)
- 206 (3.0 SDs)
- 319 (1.2e+01 SDs)
- 321 (5.5 SDs)
- 402 (1.1e+01 SDs)
- 404 (2.6 SDs)
- 406 (2.5 SDs)
- 424 (3.4 SDs)

Test: 8 outliers out of 112 datapoints (7.1%)

- 560 (3.3 SDs)
- 1118 (3.7 SDs)
- 683 (7.1 SDs)
- 378 (3.0 SDs)
- 362 (2.5 SDs)
- 710 (3.5 SDs)
- 324 (8.0 SDs)
- 654 (7.4 SDs)

### PFI (only most important descriptors):

#### Outliers (max. 10 shown)

Train: 18 outliers out of 604 datapoints (3.0%)

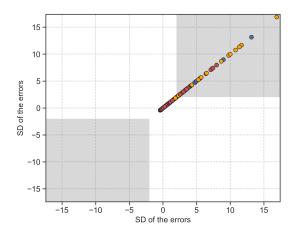
- 103 (3.0 SDs)
- 294 (4.4 SDs)
- 295 (4.5 SDs)
- 296 (4.5 SDs)
- 303 (6.7 SDs)
- 319 (8.5 SDs)
- 323 (8.7 SDs)
- 361 (3.1 SDs)
- 442 (2.4 SDs)
- 561 (2.6 SDs)

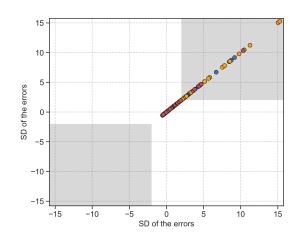
Validation: 24 outliers out of 404 datapoints (5.9%)

- 27 (2.1 SDs)
- 45 (2.7 SDs)
- 320 (1.1e+01 SDs)
- 321 (1.5e+01 SDs)
- 322 (1.5e+01 SDs)
- 325 (1e+01 SDs)
- 402 (5.1 SDs)
- 440 (7.5 SDs)
- 441 (3.9 SDs)
- 522 (5.8 SDs)

Test: 5 outliers out of 112 datapoints (4.5%)

- 683 (4.4 SDs)
- 991 (3.6 SDs)
- 710 (4.7 SDs)
- 324 (1e+01 SDs)
- 654 (3.8 SDs)

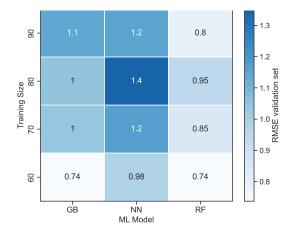


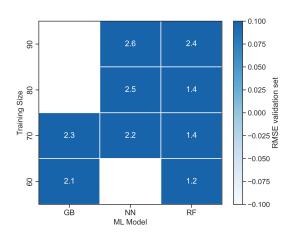


ROBERT v 1.2.1 Page 5 of 8

#### 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 (JetFuel\_Ignition.csv)

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

- Install ROBERT and its dependencies: conda install -y -c conda-forge robert
- Adjust ROBERT version: pip install robert==1.2.1
- Install scikit-learn-intelex: pip install scikit-learn-intelex==2024.7.0

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

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

python -m robert --names "Point" --y "IDT ms" --model "[RF,GB,NN]" --csv name "JetFuel Ignition.csv"

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

Originally run in Python 3.12.2 using Linux #1 SMP Fri Apr 20 16:44:24 UTC 2018

Total execution time: 493.79 seconds (the number of processors should be specified by the user)

ROBERT v 1.2.1 Page 6 of 8



# Section H. Transparency

min\_samples\_leaf: 1

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): PFI (only most important descriptors):

sklearn model: RandomForestRegressor sklearn model: RandomForestRegressor

random state: 233 random state: 43 names: Point names: Point n estimators: 60 n estimators: 5 max depth: 60 max depth: 5 max features: 0.75 max features: 1.0 min samples split: 2 min samples split: 2

min\_weight\_fraction\_leaf: 0 min\_weight\_fraction\_leaf: 0

ccp\_alpha: 0 ccp\_alpha: 0 oob\_score: False oob\_score: True max\_samples: 0.75 max\_samples: 0.75

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

min\_samples\_leaf: 1

#### No PFI (standard descriptor filter): PFI (only most important descriptors):

split: KN split: KN type: reg type: reg

error\_type: rmse error\_type: rmse



## Section I. Abbreviations

GB: gradient boosting

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 RMSE: root mean square error MCC: Matthew's correl. coefficient

**CLAS:** classification ML: machine learning RND: random

CV: cross-validation MVL: multivariate lineal models SHAP: Shapley additive explanations

PFI: permutation feature importance

F1 score: balanced F-score NN: neural network VR: voting regressor

R2: coefficient of determination GP: gaussian process

ROBERT v 1.2.1 Page 7 of 8

#### 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.

ROBERT v 1.2.1 Page 8 of 8