# forester: A Novel Approach to Accessible and Interpretable AutoML for Tree-Based Modeling

Hubert Ruczyński, Anna Kozak
Warsaw University of Technology

# forester: A Novel Approach to Accessible and Interpretable AutoML for Tree-Based Modeling
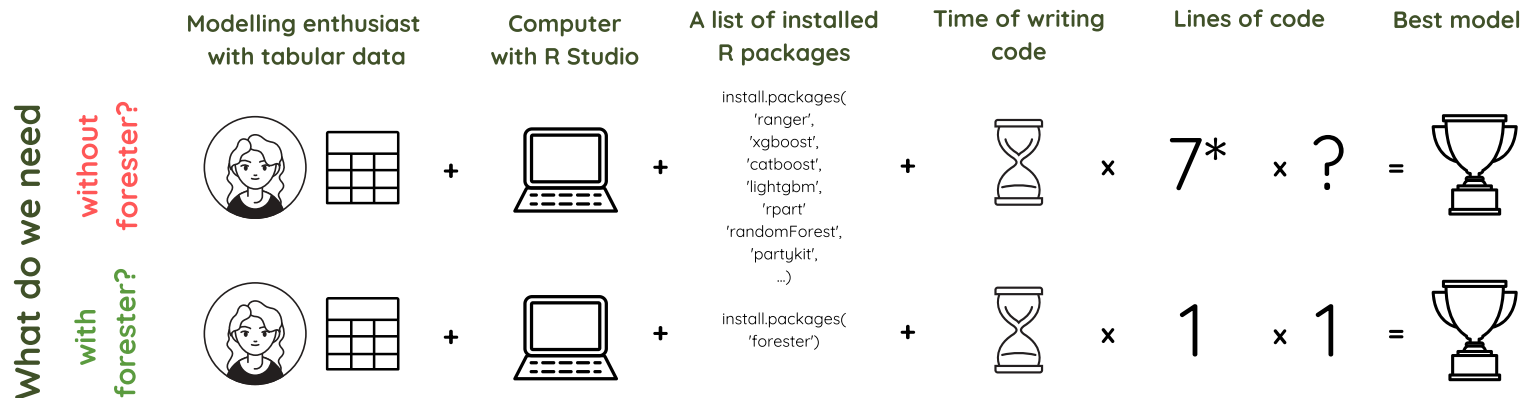
**Anna Kozak**[1]  **Hubert Ruczyński**[1]

[1]Warsaw University of Technology

**Abstract**   The majority of AutoML solutions are developed in Python. However, a large percentage of data scientists are associated with the R language. Unfortunately, there are limited R solutions available with high entry level which means they are not accessible to everyone. To fill this gap, we present the *forester* package, which offers ease of use regardless of the user's proficiency in the area of machine learning.

The *forester* package is an open-source AutoML package implemented in R designed for training high-quality tree-based models on tabular data. It supports regression and binary classification tasks. A single line of code allows the use of unprocessed datasets, informs about potential issues concerning them, and handles feature engineering automatically. Moreover, hyperparameter tuning is performed by Bayesian optimization, which provides high-quality outcomes. The results are later served as a ranked list of models. Finally, the *forester* package offers a vast training report, including the ranked list, a comparison of trained models, and explanations for the best one.

# How to build models in R?

| What do we need | | Modelling enthusiast with tabular data | Computer with R Studio | A list of installed R packages | Time of writing code | Lines of code | Best model |
|---|---|---|---|---|---|---|---|
| | without forester? | | + | + install.packages( 'ranger', 'xgboost', 'catboost', 'lightgbm', 'rpart' 'randomForest', 'partykit', ...) | + ⏳ | × 7* × ? | = 🏆 |
| | with forester? | | + | + install.packages( 'forester') | + ⏳ | × 1 × 1 | = 🏆 |

* dependent on the number of packages used

## How to use it?

```
library(forester)
data(`lisbon`)
train_output <- train(lisbon,`Price`)
```
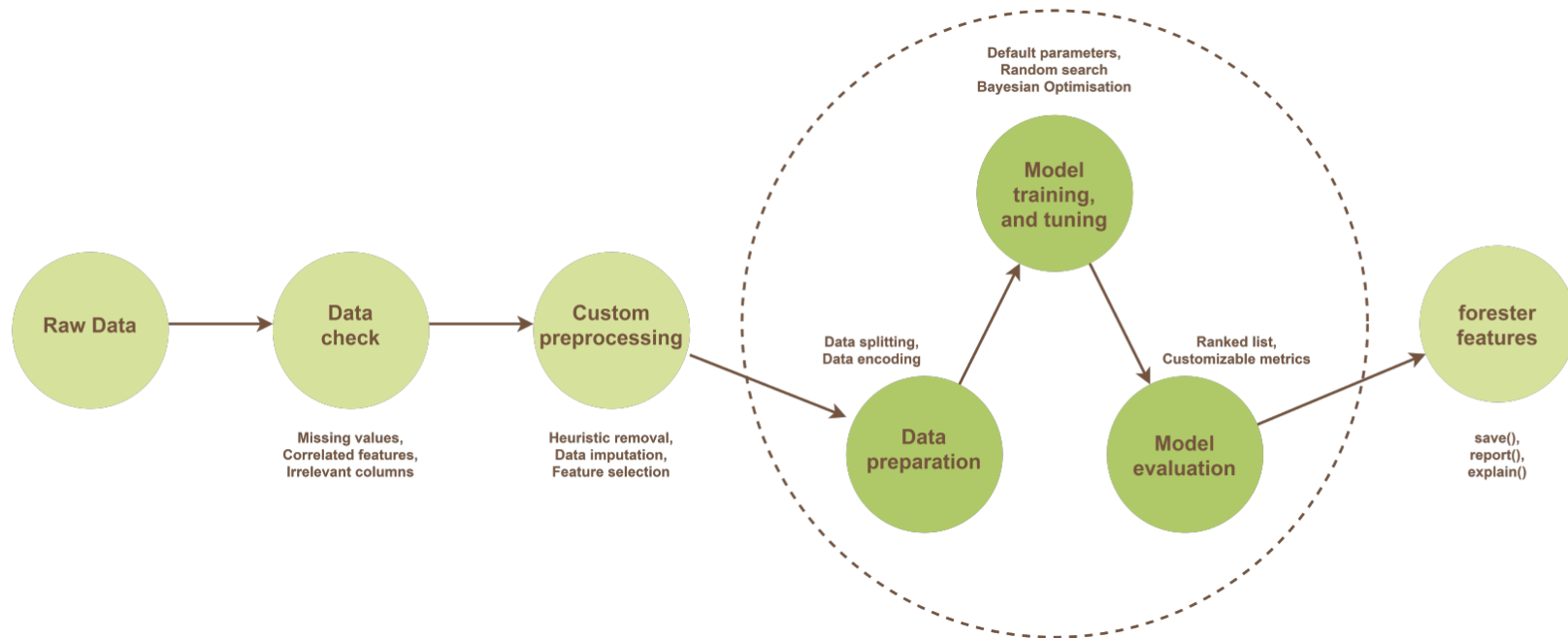
# What is the forester?

The forester is an AutoML tool in R for **tabular data regression** and **binary classification tasks\***, that wraps up all machine learning processes into a single `train()` function, which includes:

1. rendering a brief **data check report**,
2. **preprocessing** the initial dataset enough for models to be trained,
3. **training 5 tree-based models** with default parameters, random search and Bayesian optimization,
4. evaluating them and providing a **ranked list**.

However, that's not everything that the forester has to offer. Via additional functions, the user can easily explain created models with the usage of *DALEX* or generate one of the predefined **reports** including:

1. information about the dataset,
2. in-depth parameters of trained models,
3. visualizations comparing the best models,
4. explanations of the aforementioned models.

# forester pipeline



**Raw Data**

**Data check**
Missing values, Correlated features, Irrelevant columns

**Custom preprocessing**
Heuristic removal, Data imputation, Feature selection

Data splitting, Data encoding

**Data preparation**

**Model training, and tuning**
Default parameters, Random search Bayesian Optimisation

**Model evaluation**
Ranked list, Customizable metrics

**forester features**
save(), report(), explain()

# Automatic reports and XAI

# Evaluation

| Name | Number of columns | Number of rows |
|------|-------------------|----------------|
| kr-vs-kp | 37 | 3196 |
| breast-w | 10 | 699 |
| credit-approval | 16 | 690 |
| credit-g | 21 | 1000 |
| diabetes | 9 | 768 |
| phoneme | 6 | 5404 |
| banknote-authentication | 5 | 1372 |
| blood-transfusion-service-center | 5 | 748 |

## Performance comparison of forester and H2O
### for the binary classification task

# Evaluation

| Name | Number of columns | Number of rows |
|---|---|---|
| bank32nh | 33 | 8192 |
| wine_quality | 12 | 6497 |
| Mercedes_Benz_Greener_Manufacturing | 378 | 4209 |
| kin8nm | 9 | 8192 |
| pol | 49 | 15000 |
| 2dplanes | 11 | 40768 |
| elevators | 19 | 16599 |



Performance comparison of forester and H2O for the regression task

# Evaluation

Table 11: The comparison of mean execution times in seconds for the *forester* and *H2O* for binary classification experiments.

| task_name | forester | H2O | difference | relative difference |
|---|---|---|---|---|
| banknote-authentication | 818.33 | 2521.33 | -1703 | 0.28 |
| blood-transfusion-service-center | 155.67 | 555.67 | -400 | 0.26 |
| breast-w | 451.33 | 797.33 | -346 | 0.57 |
| credit-approval | 805 | 1513 | -708 | 0.53 |
| credit-g | 2453 | 4234 | -1781 | 0.58 |
| diabetes | 1645.67 | 2643.67 | -998 | 0.62 |
| kr-vs-kp | 451.33 | 806.67 | -355.33 | 0.57 |
| phoneme | 2748.33 | 3695.33 | -947 | 0.67 |

Table 12: The comparison of mean execution times in seconds for the *forester* and *H2O* for regression experiments.

| task_name | forester | H2O | difference | relative difference |
|---|---|---|---|---|
| 2dplanes | 401 | 1050.67 | -649.67 | 0.38 |
| bank32nh | 708.67 | 1214.67 | -506 | 0.58 |
| elevators | 720.33 | 1435.33 | -715 | 0.5 |
| kin8nm | 544.67 | 1564 | -1019.33 | 0.35 |
| Mercedes_Benz_Greener_Manufacturing | 848 | 1371.67 | -523.67 | 0.61 |
| pol | 756 | 1548.33 | -792.33 | 0.49 |
| wine_quality | 1317.33 | 2130 | -812.67 | 0.63 |

forester  Public

Edit Pins ▾   👁 Unwatch  11 ▾   Fork  15 ▾   ⭐ Starred  103 ▾

main ▾   |   10 branches   ⊙ 0 tags

Go to file   Add file ▾   <> Code ▾

HubertR21 Merge pull request #120 from ModelOriented/report  ···   ✓ 090bfe0 · last week   ⏱ 226 commits

| | | |
|---|---|---|
| 📁 R | fixes | last week |
| 📁 catboost_info | new dev | last year |
| 📁 data | Version 1.0.0 part 1 | last year |
| 📁 docs | update forester page | 7 months ago |
| 📁 inst/rmd | fixes | last week |
| 📁 man | fixes | last week |
| 📁 misc | fix | last week |
| 📁 pkgdown | Version 1.0.0 part 2 | last year |
| 📁 tests | forester 1.4.1 - enhanced reports | last week |
| 📁 vignettes | knowledge check fixes | 10 months ago |
| 📄 .Rbuildignore | new dev | last year |
| 📄 .gitignore | new dev | last year |
| 📄 DESCRIPTION | forester 1.4.1 - enhanced reports | last week |
| 📄 NAMESPACE | report and plots update | last week |
| 📄 NEWS.md | forester 1.4.1 - enhanced reports | last week |
| 📄 README.md | update forester page | 7 months ago |
| 📄 forester.Rproj | Version 1.0.0 part 2 | last year |

## About

Trees are all you need

🔗 modeloriented.github.io/forester/

📖 Readme
〰 Activity
⭐ 103 stars
👁 11 watching
⑂ 15 forks

Report repository

## Releases

No releases published
Create a new release

## Packages

No packages published
Publish your first package

## Contributors 9

## Deployments 74

🟢 github-pages last week

+ 73 deployments

## Languages

- ● HTML 98.5%
- ● R 1.5%

📄 README.md

# forester: Quick and Simple Tools for Training and Testing of Tree-based Models 🔗

A significant amount of time is spent on building models with high performance. Selecting the appropriate model structures, optimizing hyperparameters and explainability are only part of the process of creating a machine learning-based solution. Despite the wide range of structures considered, tree-based models are champions in competitions or

# Current directions

# The impact of data preparation on the quality of tree-based models created with AutoML forester package

Hubert Ruczyński[1]

[1]Warsaw University of Technology

**Abstract**   Automated Machine Learning (AutoML) solutions are increasingly popular, as they allow data scientists to train high-quality machine learning models with minimal effort. However, the majority of AutoML solutions are developed in Python, leaving R users with limited, and overly complex tools. In this paper, we introduce *forester*, an open-source AutoML package implemented in R that is designed for training high-quality tree-based models on tabular data. The *forester* supports regression, binary classification, and newly implemented survival analysis tasks. The focus put on a single model family, gives us an opportunity to derive conclusions about its behaviour.

Additionally, we introduce a custom preprocessing module, which creates an opportunity to validate a common belief that tree-based models do not require any data preprocessing. We answer this question by conducting a thorough ablation study, of the *forester* package, where we evaluate the impact of dozens of preprocessing strategies. Obtained results let us believe that in the case of the tree-based models family some methods, prove to be more efficient than, others. Finally, we cannot fully agree with the presented thesis, and we provide the reasons supporting our belief.

# Ablation study

| Data set | Rows | Columns | Static | Duplicate pairs | Missing fields | Dimensional issues | Correlation pairs | Imbalance | ID-like |
|---|---|---|---|---|---|---|---|---|---|
| banknote-authentication | 1372 | 5 | 0 | 0 | 0 | No | 1 | No | No |
| blood-transfusion-service-center | 748 | 5 | 0 | 0 | 0 | No | 1 | Yes | No |
| breast-w | 699 | 10 | 0 | 0 | 16 | No | 9 | Yes | No |
| credit-approval | 690 | 16 | 0 | 0 | 37 | No | 1 | No | No |
| credit-g | 1000 | 21 | 0 | 0 | 0 | No | 0 | Yes | No |
| diabetes | 768 | 9 | 0 | 0 | 0 | No | 0 | Yes | No |
| kr-vs-kp | 3196 | 37 | 4 | 0 | 0 | Yes | 0 | No | No |
| phoneme | 5403 | 6 | 0 | 0 | 0 | No | 0 | Yes | No |

| Data set | Rows | Columns | Static | Duplicate pairs | Missing fields | Dimensional issues | Correlation pairs | Imbalance | ID-like |
|---|---|---|---|---|---|---|---|---|---|
| 2dplanes | 40768 | 11 | 0 | 0 | 0 | No | 0 | No | No |
| bank32nh | 8192 | 33 | 0 | 0 | 0 | Yes | 0 | Yes | No |
| elevators | 16599 | 19 | 2 | 0 | 0 | No | 11 | Yes | No |
| kin8nm | 8192 | 9 | 0 | 0 | 0 | No | 0 | No | No |
| Mercedes_Benz_Greener_Manufacturing | 4209 | 378 | 145 | 134 | 0 | Yes | 522 | No | Yes |
| pol | 15000 | 49 | 22 | 156 | 0 | Yes | 2 | Yes | No |
| wine_quality | 6497 | 12 | 0 | 0 | 0 | No | 1 | Yes | No |

# Ablation study

# Master's degree

1. Implementation of multiclass classification (MC).

2. Adding a few more datasets for current tasks.

3. Conducting preprocessing methods ablations study for MC.

4. Evaluation of various preprocessing strategies on new data.

5. Uploading forester to CRAN Repository.

Thank you for attention!

# forester: an R package for automated building of tree-based machine learning models

Anna Kozak, Adrianna Grudzień, Hubert Ruczyński, Patryk Słowakiewicz

MI2.AI Group, Faculty of Mathematics and Information Science, Warsaw University of Technology

## Introduction

A significant amount of time is spent on building models with high performance. Selecting the appropriate model structures, optimising hyperparameters and explainability are only part of the process of creating a machine learning-based solution. Despite the wide range of structures considered, tree-based models are champions in competitions or hackathons. So, aren't tree-based models enough? They are, and that's why **we want to fully automate the process of training tree-based models so that even the newcomers can easily build**, train and understand these powerful prediction tools. At the same time, **the experienced users gain a powerful tool for making high-quality baseline models** for new tasks, they start working with.
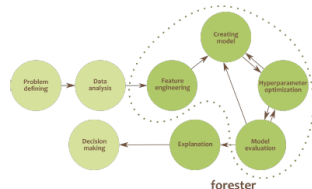
## What is the *forester*?

The *forester* is **an autoML tool in R** that wraps up all machine learning processes into a single `train()` function, which includes:

1. rendering a brief data check report,
2. preprocessing initial dataset enough for models to be trained,
3. training 5 tree-based models with default parameters, random search and Bayesian optimisation,
4. evaluating them and providing a ranked list.

However, that's not everything that the *forester* has to offer. Via additional functions, the user can easily explain created models with the usage of *DALEX* or generate one of the predefined reports including:

1. information about the dataset,
2. in-depth parameters of trained models,
3. visualisations comparing the best models,
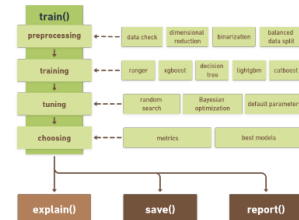4. explanations of the aforementioned models.

## Why tree-based models?

Tree-based models, especially **XGBoost are extremely popular amongst winners in Kaggle competitions** and they firmly show their superiority with tabular data, not only in terms of fast computations. Moreover, the researchers also prove that **tree-based models are superior to deep learning neural networks** because they don't suffer from uninformative columns presence and are not biased toward overly smoothed solutions.

## Package structure

With functions in *forester* package users can create a well-tuned tree-based model with a unified, simple formula. With the usage of only two required parameters: the raw, not preprocessed dataset and target column name, the user is able to achieve satisfying results. The *forester* automatically handles the "ugly" part for you.



forester

## For whom is this package created?

The *forester* is designed for beginners in data science, but also for more experienced users. They get an easy-to-use tool that can be used to prepare high-quality baseline models for comparison with more advanced methods or a set of output parameters for more thorough optimisations. **Tree-based models are created in just one line of code.** The package differentiates itself in this aspect from powerful autoML frameworks like *mlr3* and *H2O*.

| | forester | mlr3 | H2O |
|---|---|---|---|
| easy to use | ✔ | | |
| preprocessing | ✔ | ✔ | ✔ |
| autoML | ✔ | ✔ | ✔ |
| feature selection | 🕐 | ✔ | ✔ |
| model tuning | ✔ | ✔ | ✔ |
| vizualization | ✔ | ✔ | ✔ |
| explanation | ✔ | ✔ | ✔ |
| report | ✔ | | |

## Contact info

✉ anna.kozak@pw.edu.pl
✉ grudziena@outlook.com
✉ hruczynski21@interia.pl
✉ slowakiewiczpatryk@outlook.com
○ https://github.com/ModelOriented/forester

---

# *forester*: growing transparent tree-based models for everyone

Anna Kozak[1], Adrianna Grudzień[1], Hubert Ruczyński[1],

Patryk Słowakiewicz[1], Przemysław Biecek[1,2]

[1]MI2.AI, Warsaw University of Technology  [2]MI2.AI, University of Warsaw

## Let's talk about AutoML, tree-based models, explainable AI (XAI), exploratory data analysis (EDA)!



### How to build tree-based models in R?

## What is *forester*?

- full automation of the process of training tree-based models
- no demand for ML expertise
- powerful tool for making high-quality baseline models for experienced users

The *forester* package is **an AutoML tool in R** that wraps up all machine learning processes into a single `train()` function, which includes:

1. rendering a brief **data check** report,
2. **preprocessing** initial dataset enough for models to be trained,
3. **training** 5 tree-based models with default parameters, random search and Bayesian optimisation,
4. **evaluating** them and providing a ranked list.



forester::train()

## How to use it?

```
library(forester)
data('lisbon')
train_output <- train(lisbon, 'Price')
```

## For whom is this package created?

The *forester* package is designed for beginners in data science, but also for more experienced users. They get an easy-to-use tool that can be used to prepare high-quality baseline models for comparison with more advanced methods or a set of output parameters for more thorough optimisations.

## Contact info

✉ anna.kozak@pw.pl
○ https://github.com/ModelOriented/forester

## References

P. Biecek. DALEX: Explainers for Complex Predictive Models in R. *Journal of Machine Learning Research*, 19(84):1–5, 2018. URL https://jmlr.org/papers/v19/18-416.html.

A. Kozak, H. Ruczyński, P. Słowakiewicz, A. Grudzień, and P. Biecek. *forester: Quick and Simple Tools for Training and Testing of Tree-based Models*, 2022. URL https://github.com/ModelOriented/forester. R package version 1.0.0.

## Prepare meaningful report less than in 60 seconds!

As data scientists, we are fully aware that there are some time expensive processes in out work. One of them is creating a report with meaningful results. That's why one of the most powerful forester feature, which makes it a efficient tool for both experienced users and the newcomers, is a `report()` function. This single-line command is designed to **provide a holistic view on the outcomes of the ML process happening inside the forester.**

### See the report yourself!

# *forester*: A Novel Approach to Accessible and Interpretable AutoML for Tree-Based Modeling

Anna Kozak [1], Hubert Ruczyński[1]

[1]Warsaw University of Technology

**Let's talk about AutoML, tree-based models, explainable AI (XAI), and exploratory data analysis (EDA)!**

## How to build tree-based models in R?

| | Modelling enthusiast with tabular data | Computer with R Studio | A list of installed R packages | Time of writing code | Lines of code | Best model |
|---|---|---|---|---|---|---|

**What do we need without forester?**

install.packages(
'ranger',
'xgboost',
'catboost',
'lightgbm',
'rpart'
'randomForest',
'partykit',
...)

+ + + × 7* × ? =

**What do we need with forester?**

install.packages(
'forester')

+ + + × 1 × 1 =

\* dependent on the number of packages used

## What is *forester*?

- Full automation of the process of training tree-based models,
- No demand for ML expertise,
- Powerful tool for making high-quality baseline models for experienced users.

The *forester* package is **an AutoML tool in R** that wraps up all machine learning processes into a single `train()` function, which includes:

1. Rendering a brief **data check** report,
2. **Preprocessing** initial dataset enough for models to be trained,
3. **Training** 5 tree-based models with default parameters, random search and Bayesian optimisation,
4. **Evaluating** them and providing a ranked list.

## For whom is this package created?

The *forester* package is designed for beginners in data science, but also for more experienced users. They get an easy-to-use tool that can be used to prepare high-quality baseline models for comparison with more advanced methods or a set of output parameters for more thorough optimisations.

## How to use it?

```
library(forester)
data('lisbon')
train_output <- train(lisbon, 'Price')
```

Raw data → Data check (1) → Data preparation (2) → Model training and tuning (3) → Model evaluation (4) → forester features

Missing values, Correlated features, Irrelevant columns

Data splitting, Preprocessing, Data imputation

Default parameters, Random search, Bayesian Optimization

Ranked list, Customizable metrics

save(), report(), explain()

## Prepare meaningful report less than in 60 seconds!

As data scientists, we are fully aware that there are some time expensive processes in out work. One of them is creating a report with meaningful results. That's why one of the most powerful *forester* feature, which makes it a efficient tool for both experienced users and the newcomers, is a `report()` function. This single-line command is designed to **provide a holistic view on the outcomes of the ML process** happening inside of the *forester*.
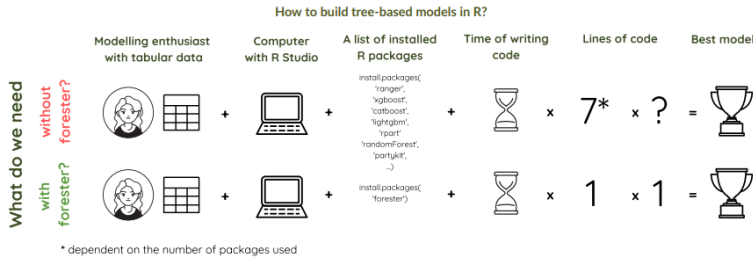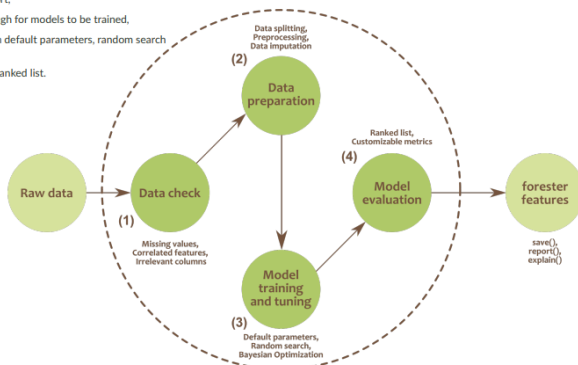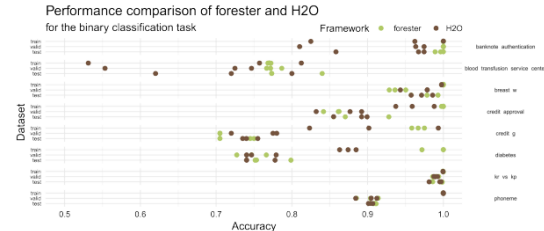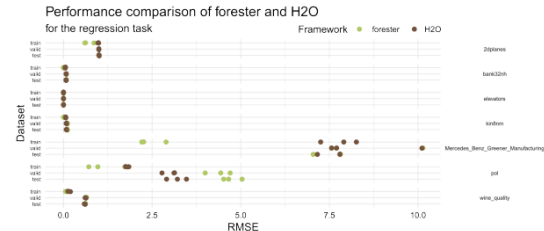
## Simple doesn't mean worse!

According to our experiments, the forester package achieves competitive results in much shorter time in comparison to well-known H2O AutoML tool. We have compared their performance on 8 binary classification, and 7 regression tasks, and the calculations were repeated 3 times for each dataset and framework. The forester outperformed H2O most of the times, even though the latter package's training lasted 2 times longer on average.

### Performance comparison of forester and H2O
for the regression task

### Performance comparison of forester and H2O
for the binary classification task

## Contact info

✉ anna.kozak@pw.edu.pl
✉ hruczynski21@interia.pl
⌂ https: //github.com/ModelOriented/forester

## Paper

## GitHub

## References

L. Grinsztajn, E. Oyallon, and G. Varoquaux. Why do tree-based models still outperform deep learning on typical tabular data? In *Thirty-sixth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2022. URL https://openreview.net/forum?id=Fp7__phQszn.

A. Kozak, H. Ruczyński, P. Słowakiewicz, A. Grudzień, and P. Biecek. *forester: Quick and Simple Tools for Training and Testing of Tree-based Models*, 2023. URL https://github.com/ModelOriented/forester. R package version 1.1.4.

# Investigating the Efficiency of Tree-based Models for Tabular Data with *forester* Package
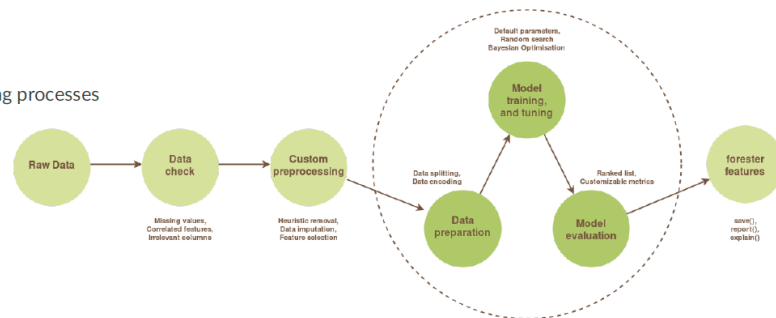
Hubert Ruczyński, Anna Kozak

Warsaw University of Technology
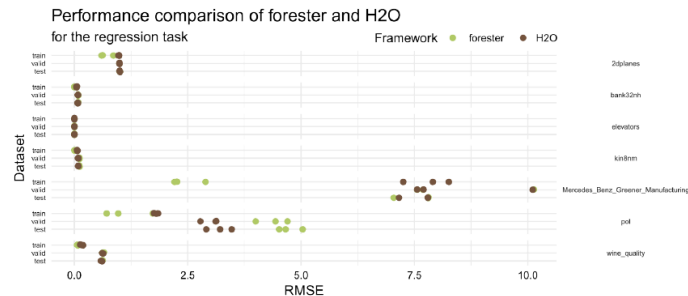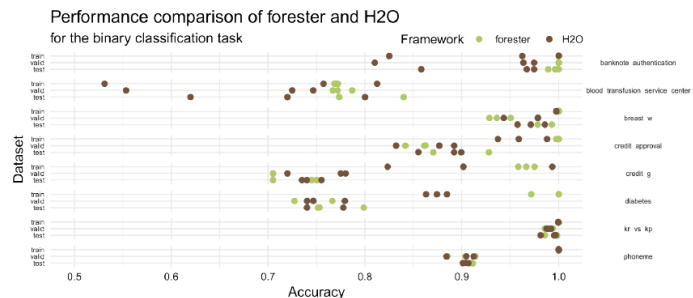
## What is *forester*?

The *forester* package is **an AutoML tool in R** that wraps up all machine learning processes into a single `train()` function, which includes:

1. rendering a brief **data check** report,

2. **preprocessing** initial dataset enough for models to be trained,

3. **training** 5 tree-based models with default parameters, random search and Bayesian optimization,

4. **evaluating** them and providing a ranked list.



## Simple doesn't mean worse!

According to our experiments, the forester package achieves competitive results in much shorter time in comparison to well-known H2O AutoML tool. We have compared their performance on 8 binary classification, and 7 regression tasks, and the calculations were repeated 3 times for each dataset and framework. The forester outperformed H2O most of the times, even though the latter package's training lasted 2 times longer on average.

Preprocessing time comparison
if any feature selection method was used or not

Feature Selection: none, yes



Preprocessing time comparison
depending on feature selection method
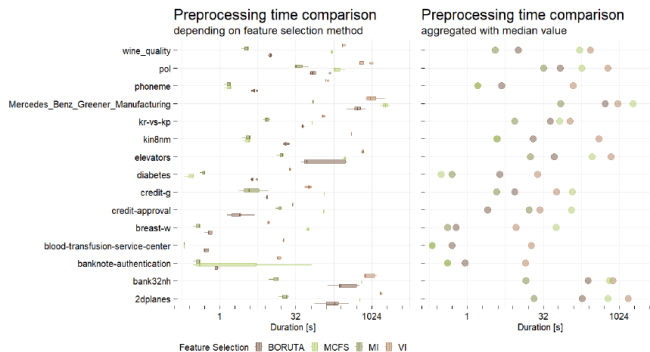
Preprocessing time comparison
aggregated with median value

Feature Selection: BORUTA, MCFS, MI, VI
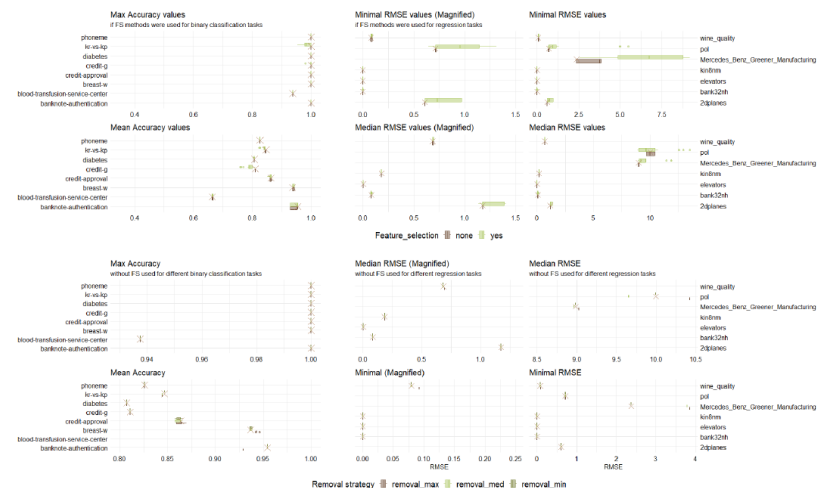
## Time complexity

The most important findings consider feature selection (FS) methods. They are the most expensive part of preprocessing, and their execution times differ significantly between the methods. The right plot shows us that Mutual Information (MI) based selection method, and BORUTA are relatively fast, whereas Monte Carlo Feature Selection (MCFS), and Variable Importance (VI) are rather slow.



Max Accuracy values
if FS methods were used for binary classification tasks

Minimal RMSE values (Magnified)
if FS methods were used for regression tasks

Minimal RMSE values

Mean Accuracy values

Median RMSE values (Magnified)

Median RMSE values

Feature_selection: none, yes

Max Accuracy
without FS used for different binary classification tasks

Median RMSE (Magnified)
without FS used for different regression tasks

Median RMSE
without FS used for different regression tasks

Mean Accuracy

Minimal (Magnified)

Minimal RMSE

Removal strategy: removal_max, removal_med, removal_min

## Feature selection impact on performance

FS methods are responsible for unstable results, and in most cases, its usage leads to worse results than for baseline methods marked with **X**. In some cases however, with FS methods we can obtain better results.

When we consider preprocessing strategies based on heuristic removals, the results are less significant, but in most cases lead to enhancements of the results.

## Contact info

✉ hruczynski21@interia.pl

⌂ https://github.com/ModelOriented/forester

## References

A. Kozak and H. Ruczyński. forester: A Novel Approach to Accessible and Interpretable AutoML for Tree-Based Modeling. In *AutoML Conference 2023 (ABCD Track)*, 2023. URL https://openreview.net/forum?id=Q3DWpGoX7PD.