# Automated Machine Learning for Remaining Useful Life Predictions
## Online Appendix

## A Search Space Description

The search space $\Lambda$ of AUTORUL allows the creation of 624 unique pipelines, which can be further configured by a total of 168 hyperparameters. In the following tables, more details regarding the available algorithms and hyperparameters are provided. For each algorithm, the total number of hyperparameters (#$\lambda$), the number of categorical (cat) and numerical (num) hyperparameters, is given. Numbers in parentheses denote conditional hyperparameters. The total number of hyperparameters does not add up to the reported 168 hyperparameters as hyperparameters of some algorithms, like, for example, TS Fresh or window generation, are included twice (once for tabular regression and once for sequence-to-sequence regression using neural networks). Components highlighted in italic are not directly included in the pipeline but influence the fitting procedure of the pipeline.

Table 1: Preprocessing Algorithms

| Name | #$\lambda$ | cat | num |
|---|---|---|---|
| Imputation | 1 | 1 (0) | 0 (0) |
| Exponential Smoothing | 2 | 0 (0) | 2 (0) |
| Robust Scaler | 2 | 0 (0) | 2 (0) |
| Normalizer | 0 | 0 (0) | 0 (0) |
| Min Max Scaling | 0 | 0 (0) | 0 (0) |
| Standardizer | 0 | 0 (0) | 0 (0) |

Table 2: Feature Engineering Algorithms

| Name | #$\lambda$ | cat | num |
|---|---|---|---|
| Window Generation | 2 | 0 (0) | 2 (0) |
| Flattening | 0 | 0 (0) | 0 (0) |
| TS Fresh | 43 | 43 (0) | 0 (0) |
| PCA | 2 | 1 (0) | 1 (0) |
| Select Percentile | 2 | 1 (0) | 1 (0) |
| Select Rates | 3 | 2 (0) | 1 (0) |

Table 3: Tabular Regression Algorithms

| Name | #$\lambda$ | cat | num |
|---|---|---|---|
| Extra Trees | 5 | 2 (0) | 3 (0) |
| Gradient Boosting | 6 | 0 (0) | 6 (0) |
| MLP | 6 | 2 (0) | 4 (1) |
| Passive Aggressive | 4 | 2 (0) | 2 (0) |
| Random Forest | 3 | 1 (0) | 2 (0) |
| SGD | 6 | 2 (0) | 4 (1) |

Table 4: Sequence Regression Algorithms

| Name | #$\lambda$ | cat | num |
|---|---|---|---|
| *Optimizer* | 4 | 0 (0) | 4 (0) |
| *Trainer* | 2 | 0 (0) | 2 (0) |
| CNN | 5 | 1 (0) | 4 (1) |
| GRU | 4 | 1 (0) | 3 (1) |
| LSTM | 4 | 1 (0) | 3 (1) |
| TCN | 5 | 1 (0) | 4 (1) |

# B   Detailed Experimental Results

Besides the results reported in the manuscript, we want to provide further insights into the generated pipelines. Figure 1 contains visualizations of the final performances reported in *Table 1* of the manuscript per benchmark dataset.
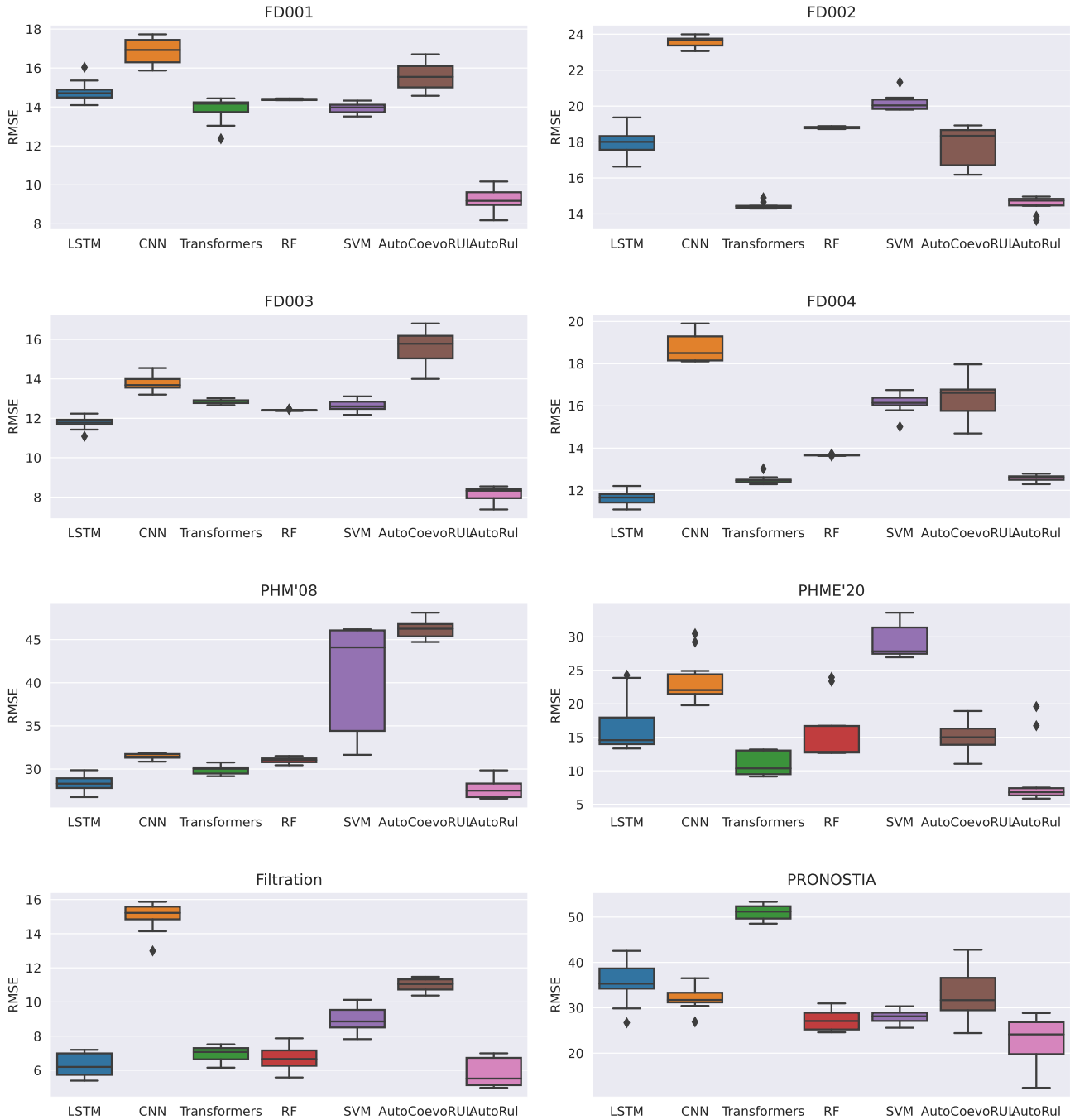


Figure 1: Performance visualizations of all tested RUL methods on all benchmark datasets.

Next, we take a closer look at the behaviour of AUTORUL during the optimization. Table 5 contains statistics about the number of evaluated configurations for each dataset. Reported are the total number of evaluated configurations, the number of successful evaluations, the number of failed evaluations (for example due to an exception during model fitting), the number of configurations where fitting was canceled after five minutes, and the number of configurations that were trained on the complete training data (no multi-fidelity approximation). In general, the vast majority of the evaluated configurations was fitted successfully. Less than 3% of the evaluated configurations were aborted after the configured timeout of five minutes, indicating that the limit was not selected to aggressively. Even though a large number of different configurations was evaluated for each dataset, only roughly 5% of the configurations were evaluated on the full budget, i.e., fully fitted until convergence.

Table 5: Statistics about the evaluated configurations. Results are averaged over ten repetitions.

| Dataset | # Configurations | # Success | # Failed | # Timeout | # Full Budget |
|---|---|---|---|---|---|
| FD001 | $761.8 \pm 144.48$ | $739.1 \pm 135.70$ | $8.5 \pm 9.92$ | $14.2 \pm 6.71$ | $41.3 \pm 7.93$ |
| FD002 | $932.5 \pm 382.89$ | $897.2 \pm 391.24$ | $9.1 \pm 7.71$ | $26.2 \pm 14.28$ | $50.2 \pm 21.61$ |
| FD003 | $648.8 \pm 62.75$ | $631.9 \pm 65.17$ | $3.3 \pm 1.55$ | $13.6 \pm 5.64$ | $35.4 \pm 3.32$ |
| FD004 | $991.9 \pm 462.97$ | $923.0 \pm 356.89$ | $46.7 \pm 5.83$ | $22.2 \pm 5.83$ | $53.2 \pm 26.46$ |
| PHM'08 | $355.0 \pm 26.65$ | $323.1 \pm 25.04$ | $28.5 \pm 0.92$ | $28.5 \pm 4.34$ | $18.7 \pm 1.55$ |
| PHME'20 | $886.9 \pm 125.11$ | $875.0 \pm 122.06$ | $8.0 \pm 2.47$ | $8.0 \pm 4.90$ | $48.4 \pm 6.48$ |
| Filtration | $818.8 \pm 134.86$ | $799.0 \pm 137.68$ | $5.7 \pm 3.29$ | $14.1 \pm 10.09$ | $44.7 \pm 7.87$ |
| PRONOSTIA | $471.5 \pm 79.13$ | $354.7 \pm 76.85$ | $86.8 \pm 12.84$ | $30.0 \pm 9.15$ | $25.3 \pm 4.50$ |

Finally, we take a closer look at the pipelines constructed by AUTORUL. Table 6 provides an overview of the constructed ensembles and the pipelines in them for each of the benchmark datasets. The *ensemble size* column shows the average number of pipelines in the constructed ensemble. It is apparent that the option to build ensembles is used for all datasets but the maximum ensemble size of ten pipelines is not consistently reached. Regarding the used template, both sequence-to-sequence regression and tabular regression are used. Depending on the dataset, either one of the two options can be basically used exclusively or both options can be mixed together. Similarly, for feature engineering all three available algorithms are extensively used with TS Fresh being used the most. Finally, for most datasets a specific type of regression algorithm is used significantly more often than others. Only for the PHM'08 dataset four different algorithm types are used roughly equally often. Random forests and TCNs are by far the most used regression algorithms. In contrast, SGD and extra trees are not included in any of the final ensembles and can probably pruned from the search space.

Table 6: Overview of constructed pipelines for each benchmark dataset. Numbers in parentheses represent the fraction of pipelines using the according component. The sum of fractions within one cell can be unequal to 1.00 due to rounding errors.

| Dataset | Ensemble Size | Template | Feat. Eng. | Regressor | |
|---|---|---|---|---|---|
| FD001 | 8.6 ± 0.92 | seq2seq (1.0) | Flatten (0.35) | CNN (0.20) | GRU (0.03) |
| | | | Limited (0.01) | LSTM (0.02) | TCN (0.74) |
| | | | TS Fresh (0.64) | | |
| FD002 | 6.0 ± 1.84 | seq2seq (0.15) | Flatten (0.12) | Gradient Boosting (0.13) | |
| | | Tabular (0.85) | Limited (0.67) | Random Forest (0.72) | |
| | | | TS Fresh (0.22) | GRU (0.03) | LSTM (0.02) |
| | | | | TCN (0.10) | |
| FD003 | 7.8 ± 1.17 | seq2seq (0.95) | Flatten (0.31) | Gradient Boosting (0.01) | |
| | | Tabular (0.05) | Limited (0.08) | Random Forest (0.04) | |
| | | | TS Fresh (0.61) | CNN (0.10) | GRU (0.01) |
| | | | | LSTM (0.03) | TCN (0.81) |
| FD004 | 5.9 ± 1.22 | seq2seq (0.15) | Flatten (0.16) | Gradient Boosting (0.03) | |
| | | Tabular (0.85) | Limited (0.36) | Random Forest (0.81) | |
| | | | TS Fresh (0.47) | CNN (0.07) | |
| | | | | TCN (0.08) | |
| PHM'08 | 5.8 ± 1.83 | seq2seq (0.74) | Flatten (0.18) | Passive Aggressive (0.02) | |
| | | Tabular (0.26) | Limited (0.54) | Random Forest (0.24) | |
| | | | TS Fresh (0.30) | CNN (0.33) | GRU (0.16) |
| | | | | LSTM (0.05) | TCN (0.21) |
| PHM'20 | 7.2 ± 1.17 | seq2seq (0.99) | Flatten (0.05) | Passive Aggressive (0.01) | |
| | | Tabular (0.01) | Limited (0.04) | CNN (0.01) | GRU (0.36) |
| | | | TS Fresh (0.90) | LSTM (0.33) | TCN (0.28) |
| Filtration | 5.8 ± 1.83 | seq2seq (0.14) | Flatten (0.10) | Gradient Boosting (0.03) | |
| | | Tabular (0.86) | TS Fresh (0.90) | Random Forest (0.83) | |
| | | | | GRU (0.07) | LSTM (0.03) |
| | | | | TCN (0.03) | |
| PRONOSTIA | 4.5 ± 1.12 | Tabular (1.00) | Flatten (0.24) | Gradient Boosting (0.16) | |
| | | | Limited (0.27) | MLP (0.18) | |
| | | | TS Fresh (0.49) | Passive Aggressive (0.18) | |
| | | | | Random Forest (0.49) | |