



UNIVERSITY OF
COPENHAGEN

Correlating FIES radial velocity drifts and telemetry data to improve calibration

Presented by:

Jonathan Jegstrup, Marcus Bredtved, Simone Vejlgard and Runi Sørensen

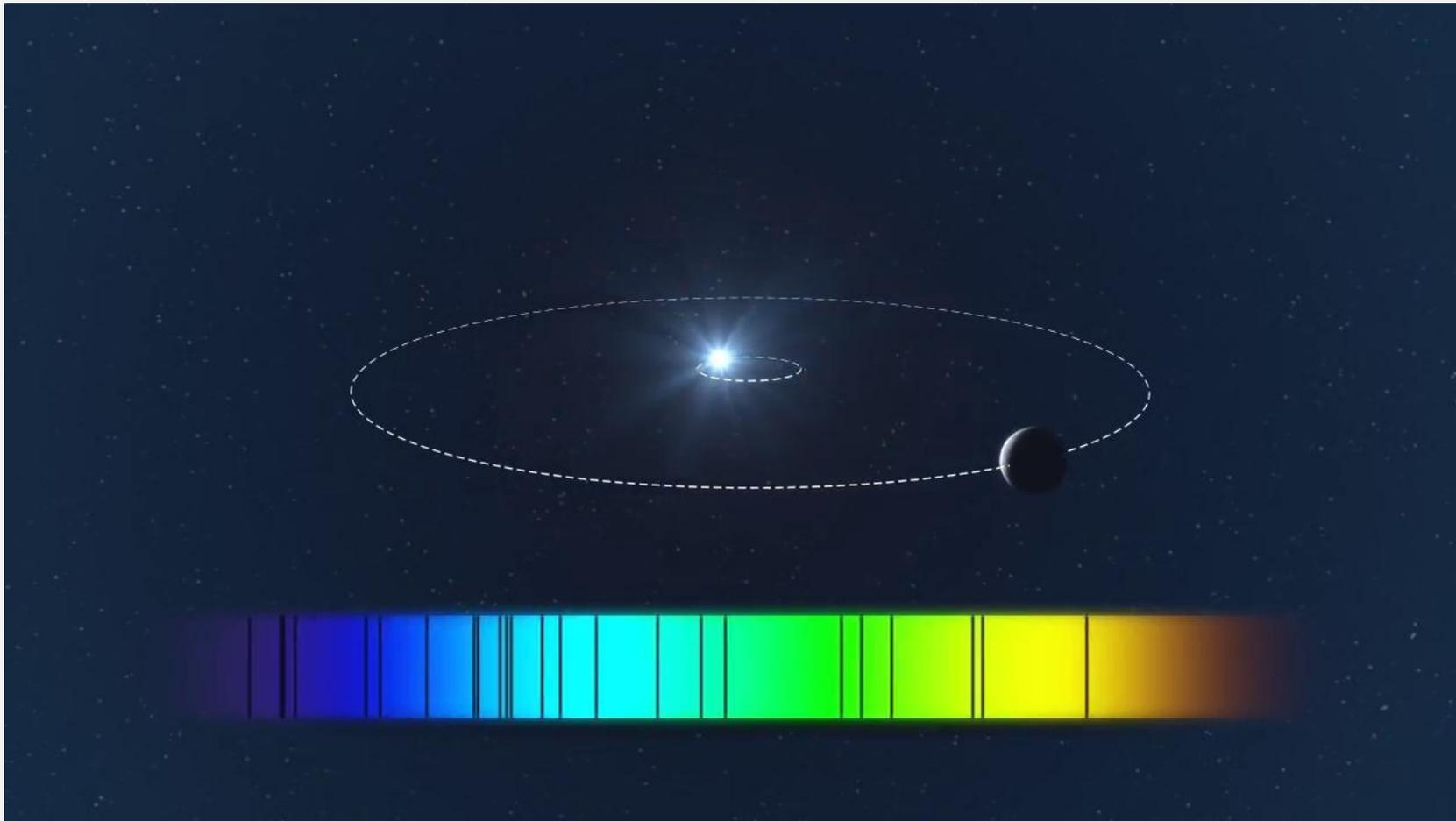
Thanks to

Lars A. Buchhave and the DTU Space Exoplanet Group for kindly providing us data
Troels C. Petersen for continuous feedback

Introduction and motivation

- Data from exoplanet research group at DTU Space
- ‘Going in blind’ - Uncertain ML applicability, but scientific merit!
- Science case: Improved calibration of radial velocity drift due to instrument drift will increase FIES (spectrograph) ability to provide confirmations and mass measurements for exoplanet candidates discovered by TESS (space observatory) around bright, nearby stars - and help the search for Earth-like exoplanets.

Quick intro to radial velocity method and the spectrograph



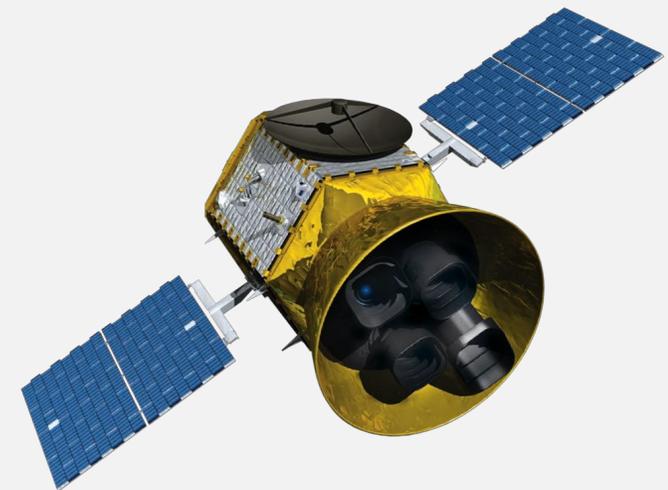
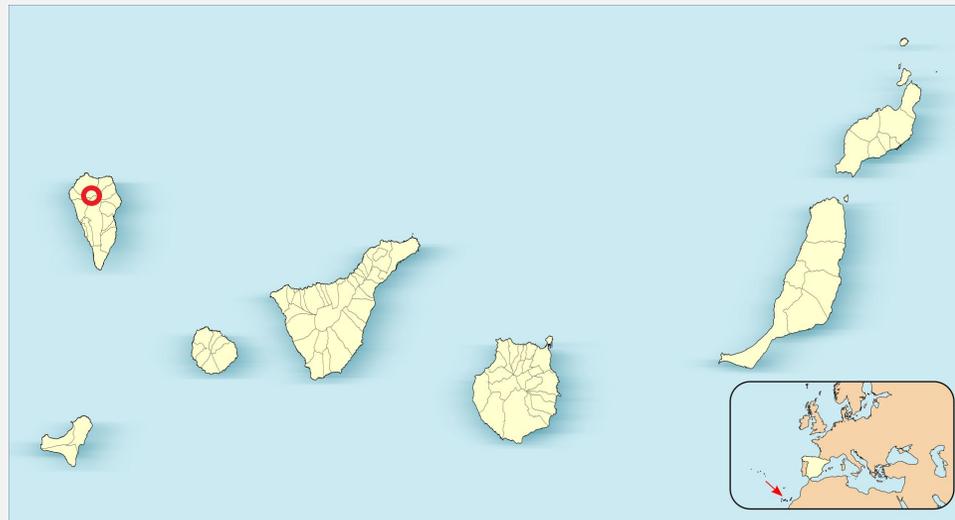
Credit: ESO/L. Calçada

Quick intro to radial velocity method and the spectrograph

Nordic Optic Telescope on the island of La Palma, Canary Islands, Spain



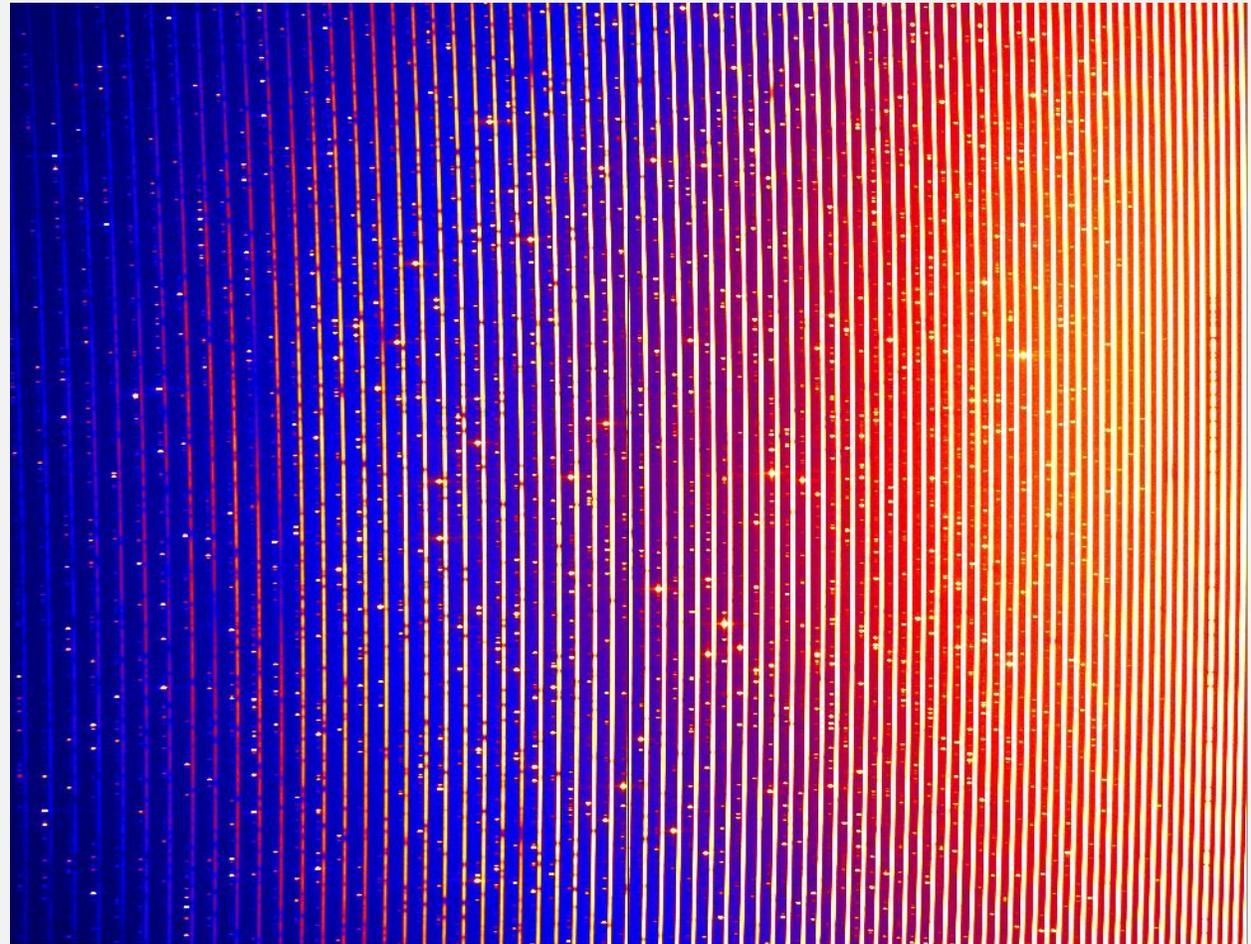
NOT



TESS

Quick intro to radial velocity method and FIES

- Path of light is very complicated - many steps
- Extreme stability is needed for precise RV measurements
 - $1/1000$ of a CCD pixel = 1 m/s
- Drift arises from tiny changes due to mechanical or thermal fluctuations - changes the physical properties of the instrument

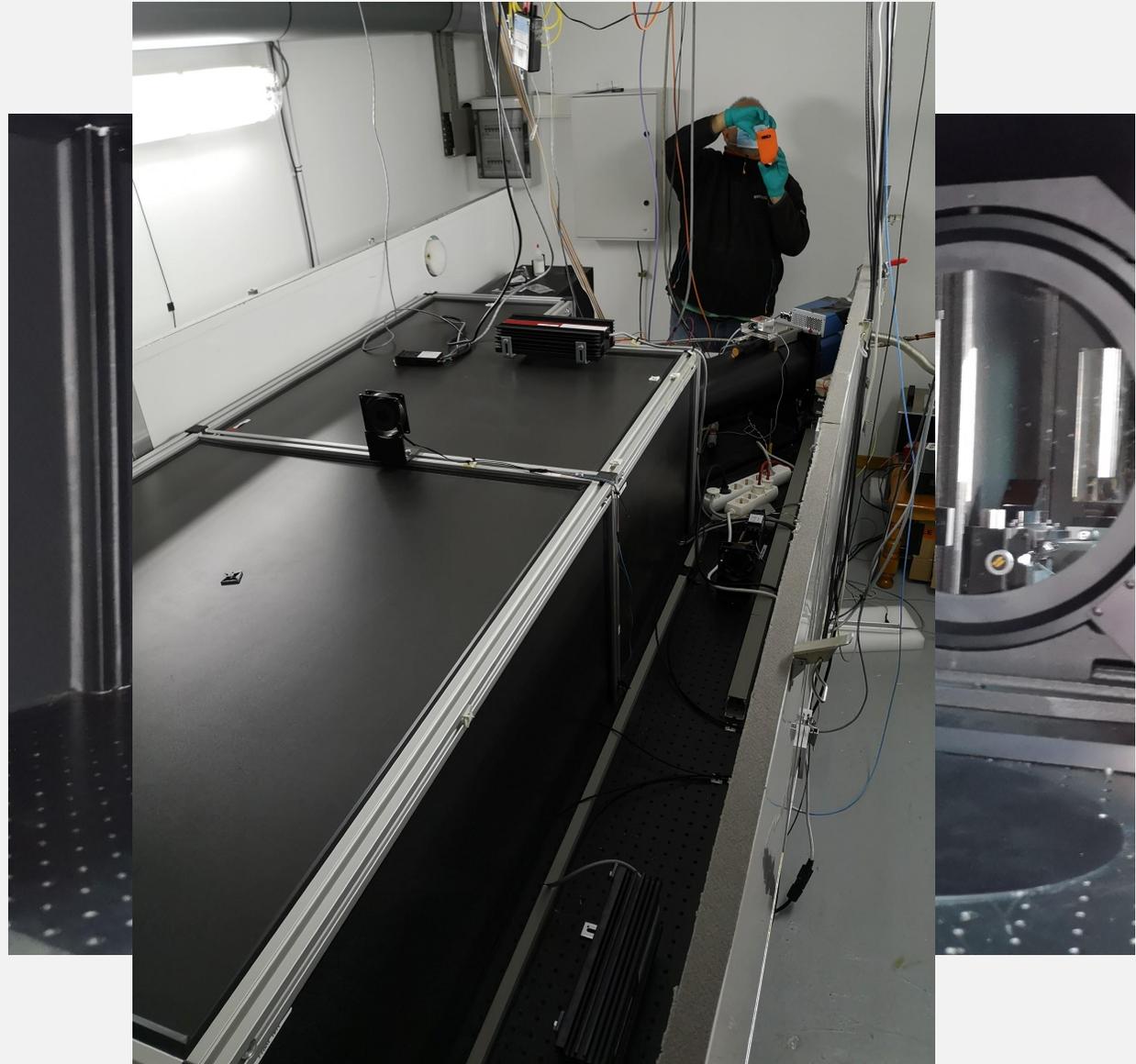
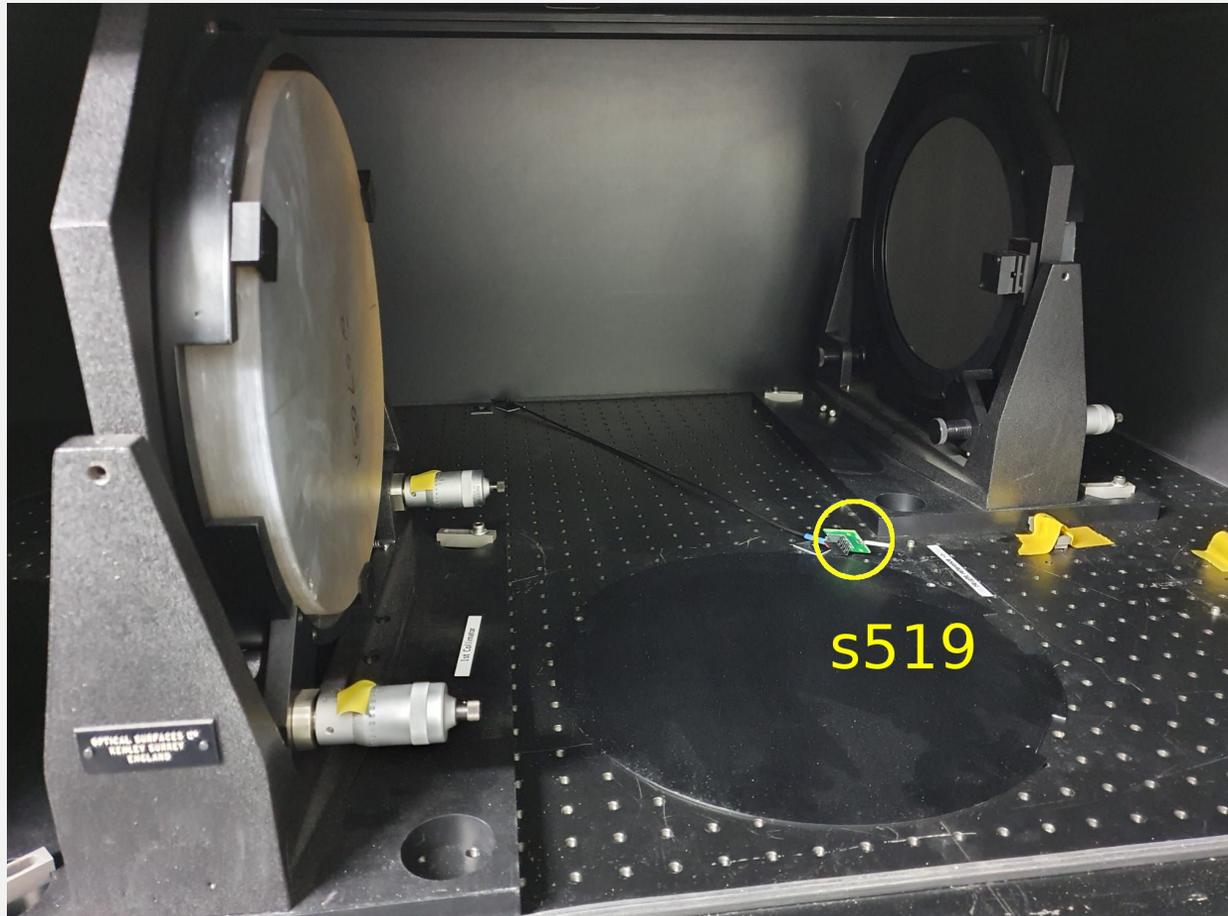


Raw spectrum from FIES

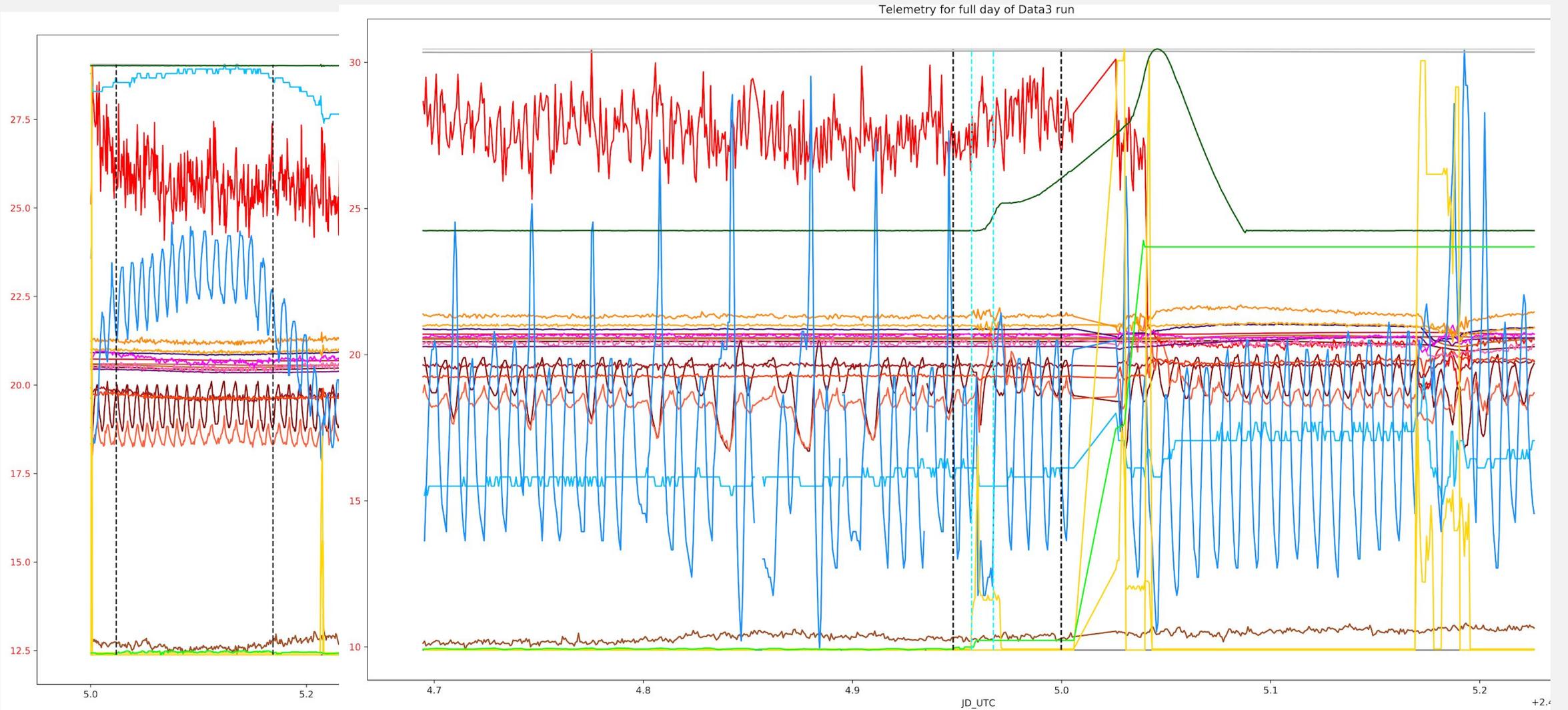
Motivation and goals

- Data of the uncalibrated spectrograph is drifting, likely due to temperature shifts in the spectrograph itself
- Can we understand these drifts based on the available telemetry placed around the instrument?
- Task is to see if a combination of telemetry sensors can train a model that describes the drift - or if we can ascertain that the telemetry itself needs to be improved

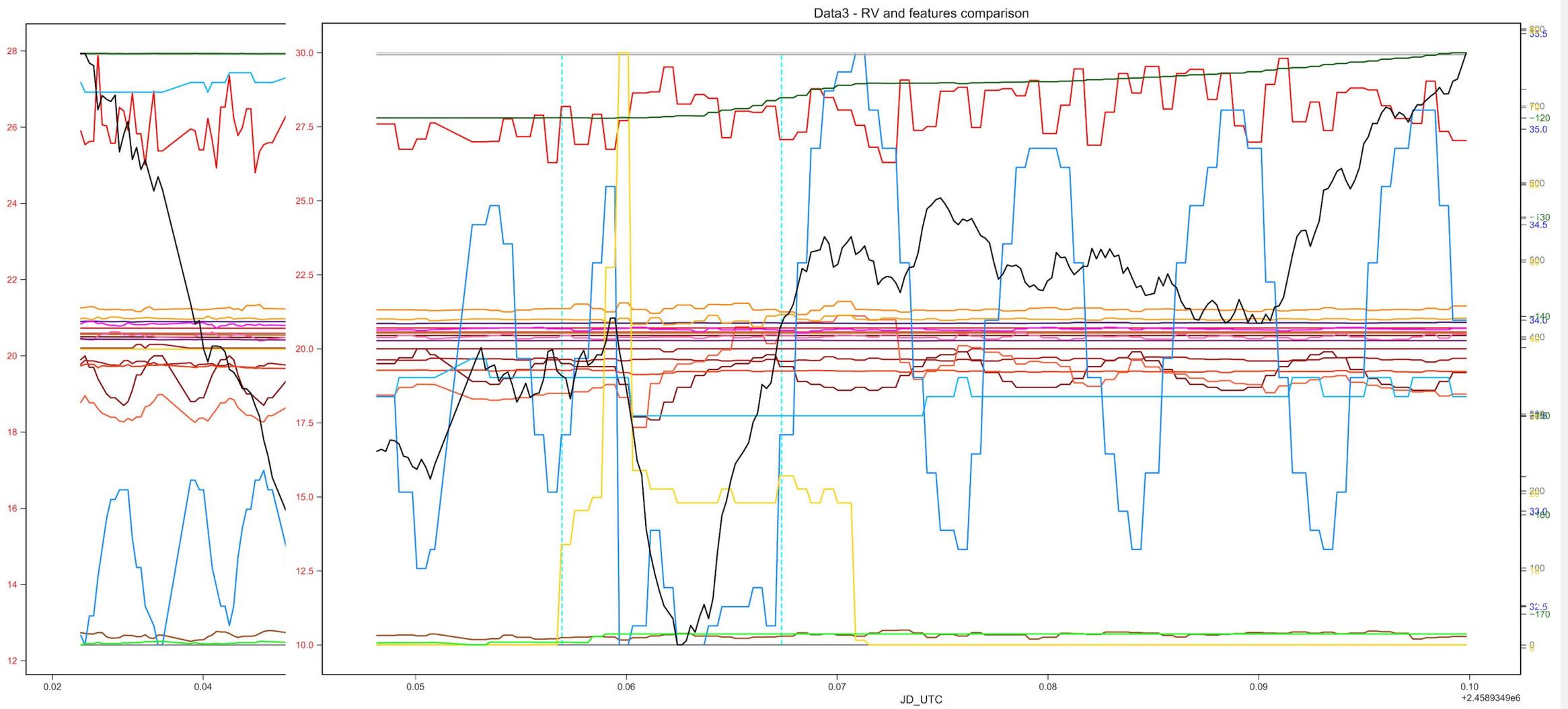
Instrumentation - the insides and sensors



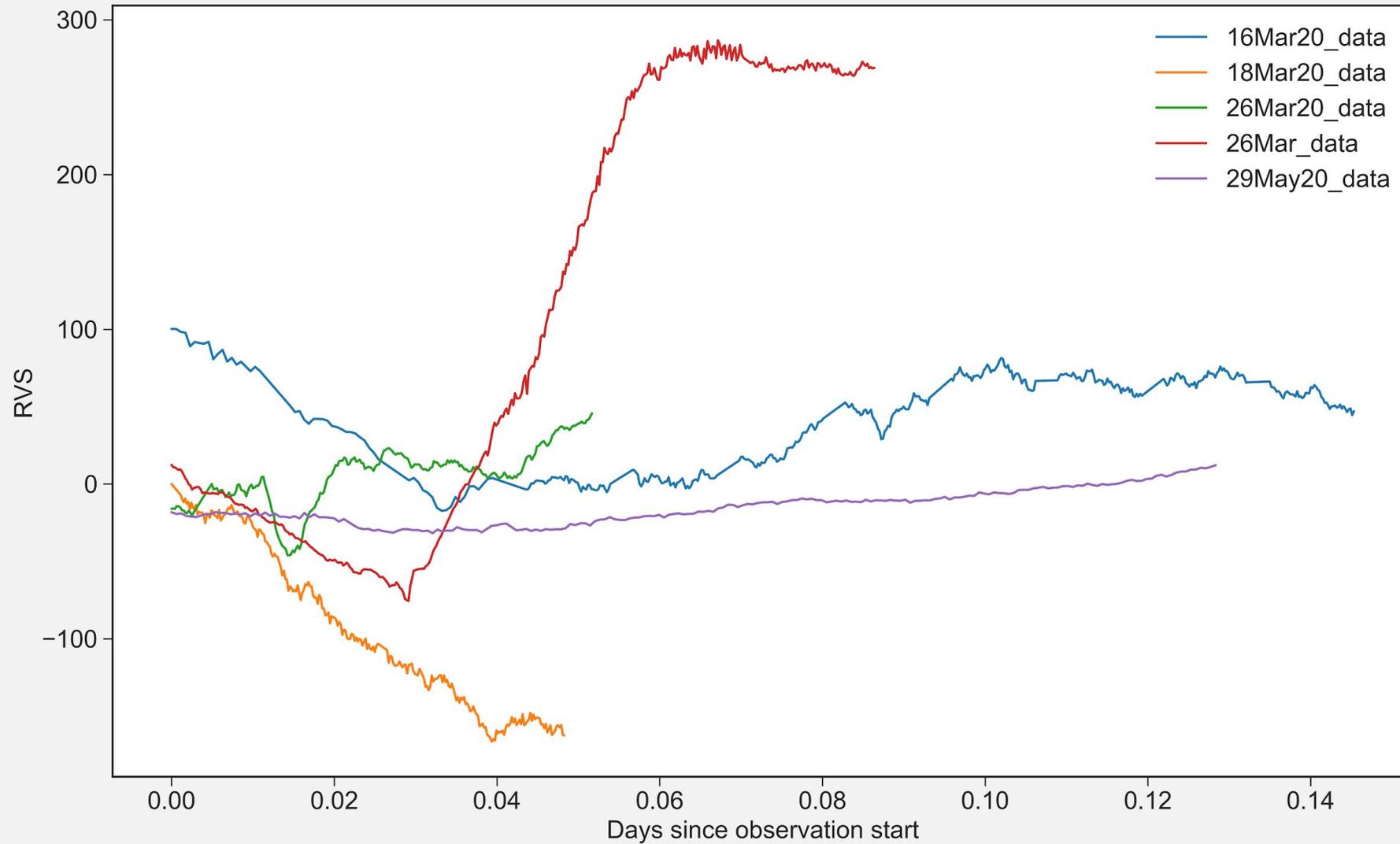
Visualization of the telemetry data



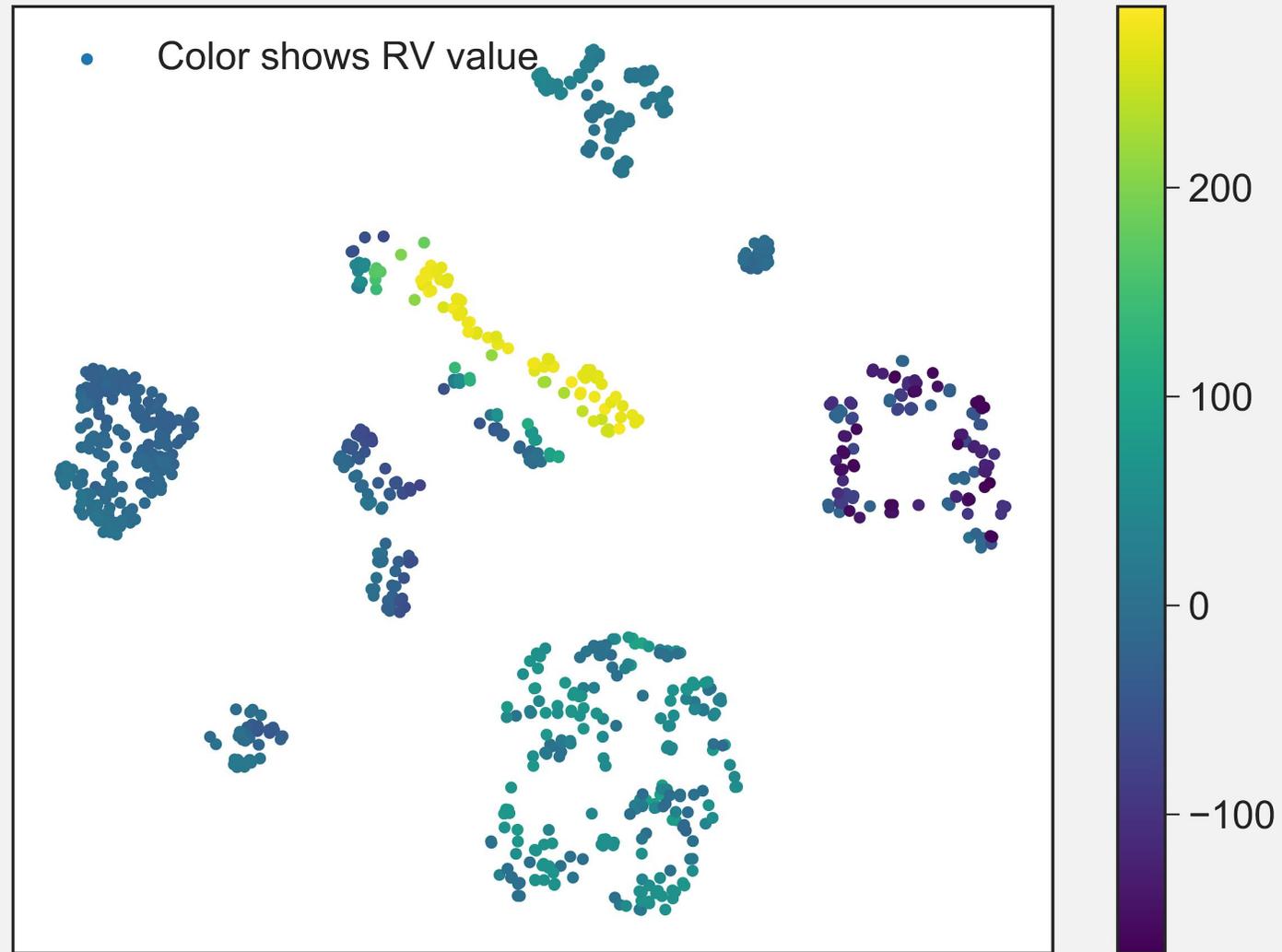
Visualization of the telemetry with RV data



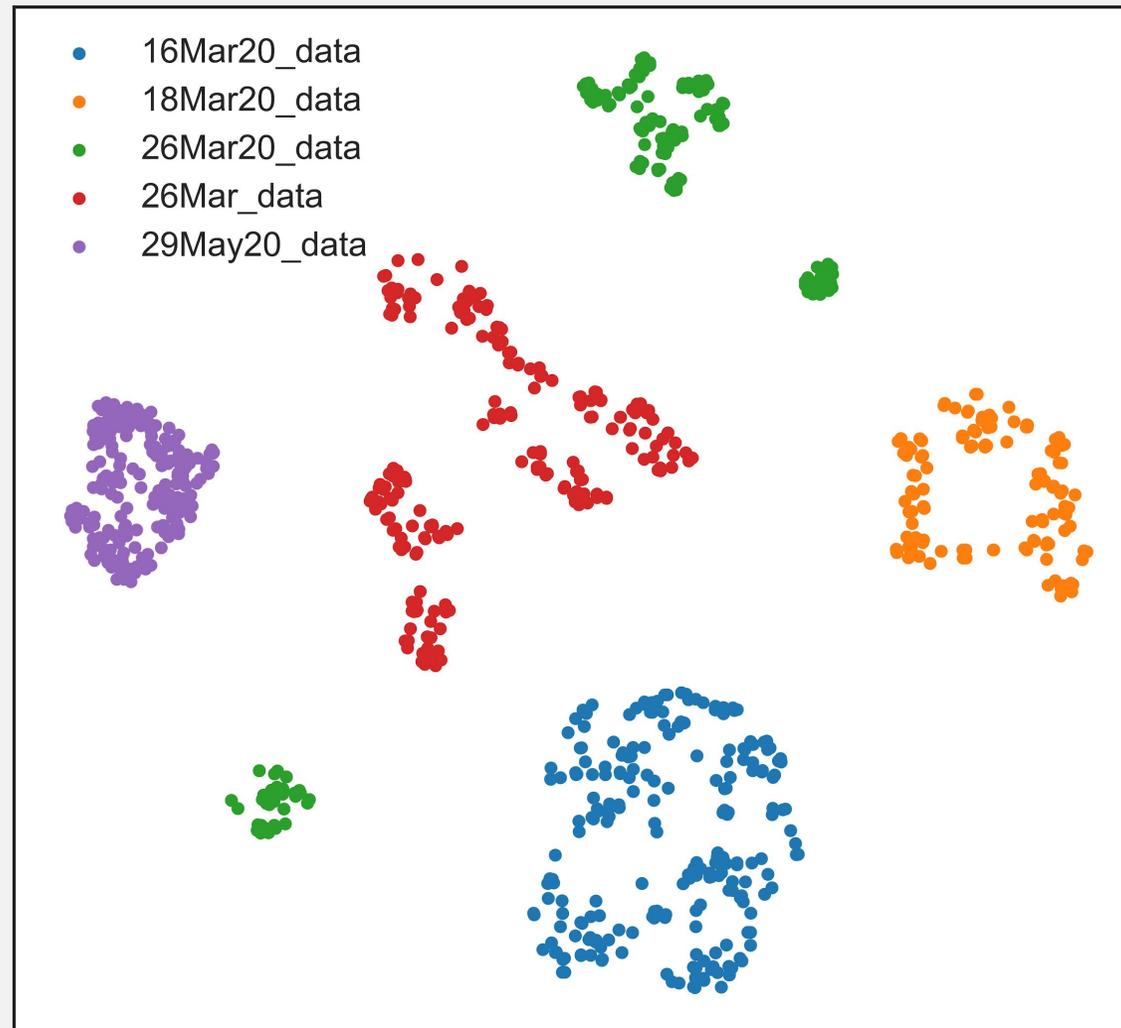
The 5 runs



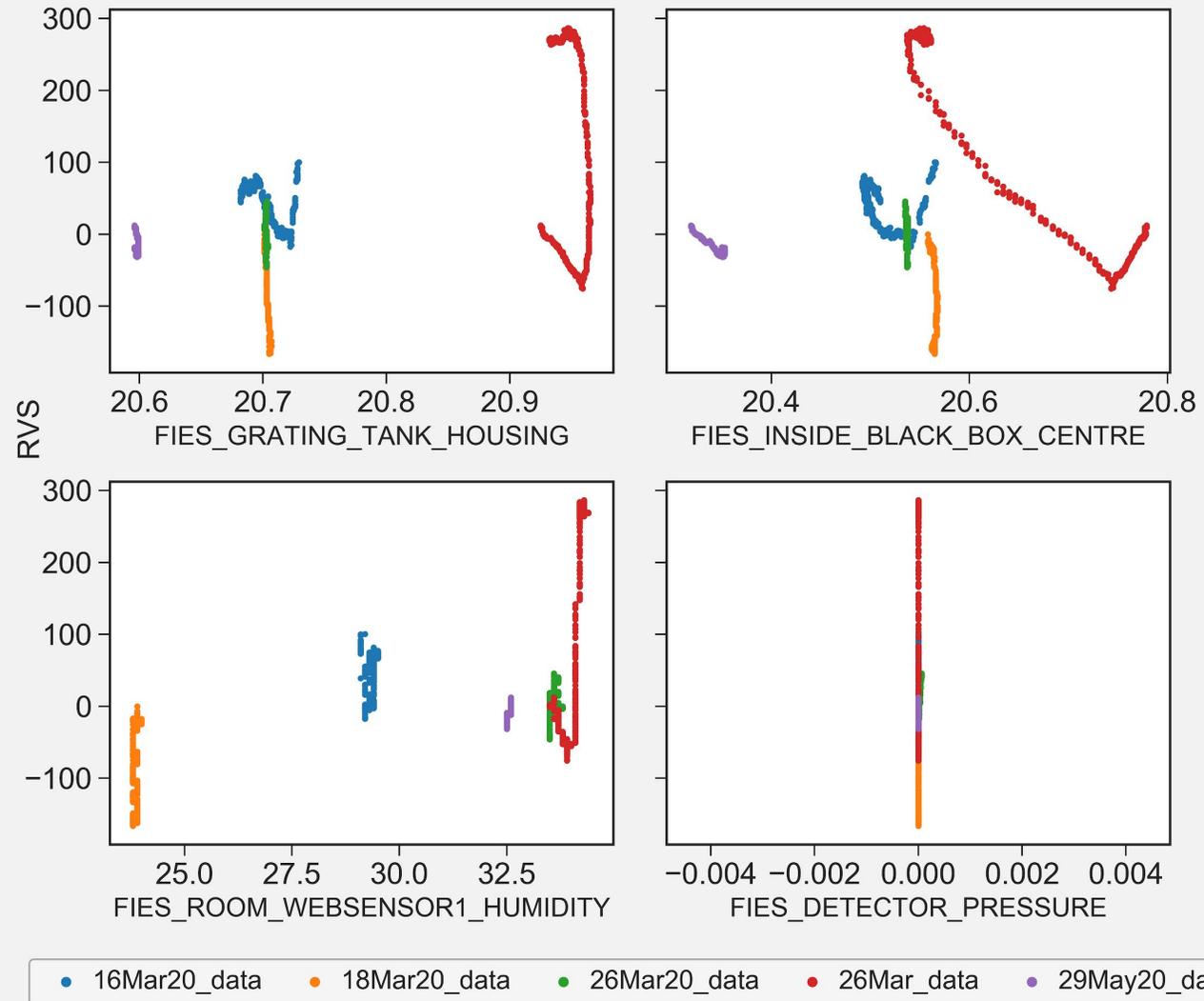
t-SNE dimensionality reduction



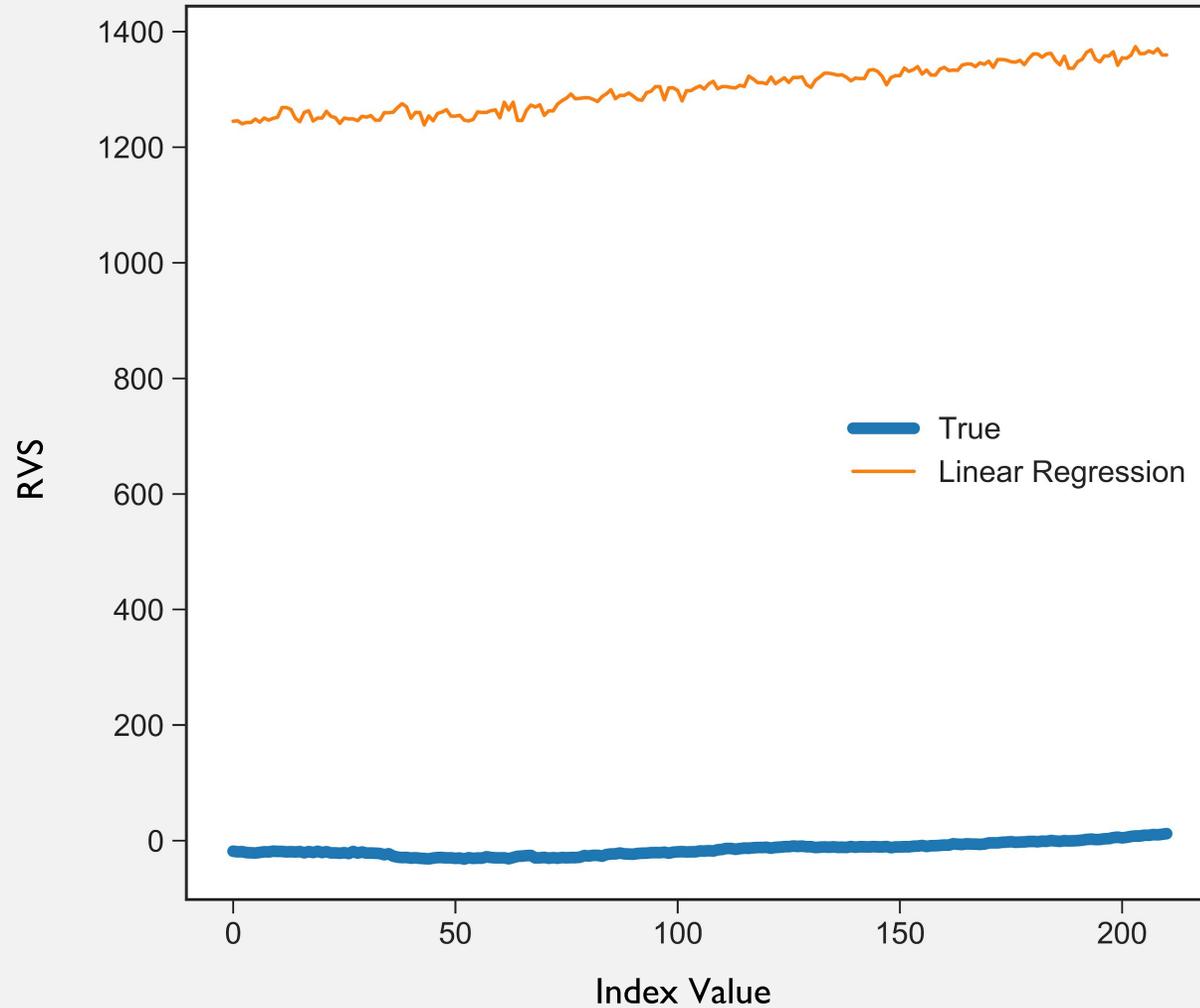
t-SNE dimensionality reduction



Parameter correlation to RVS



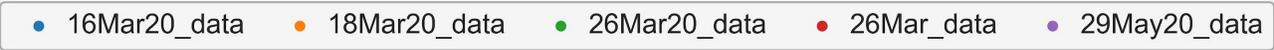
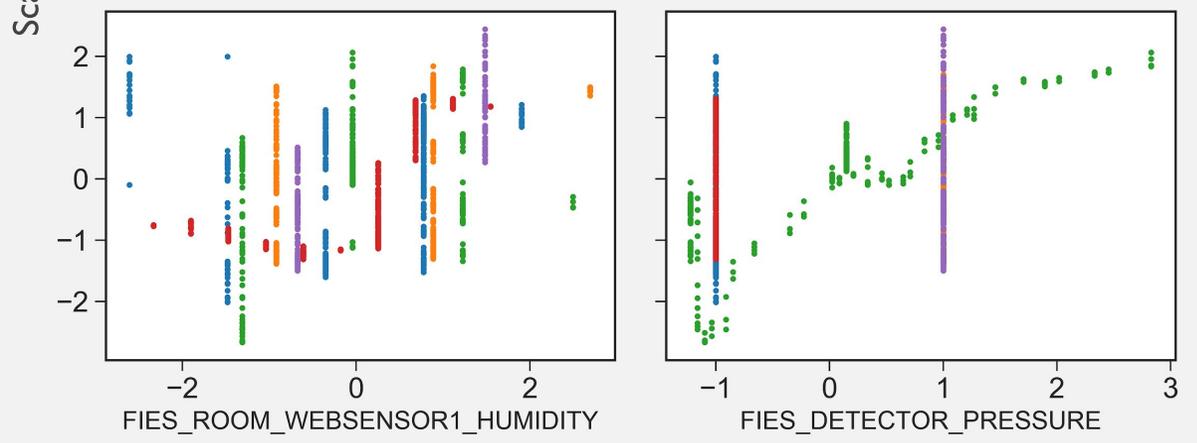
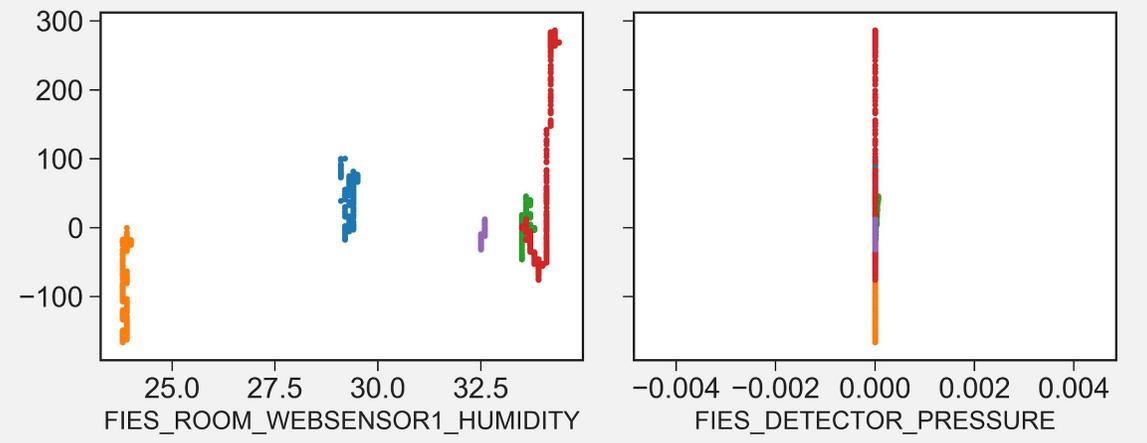
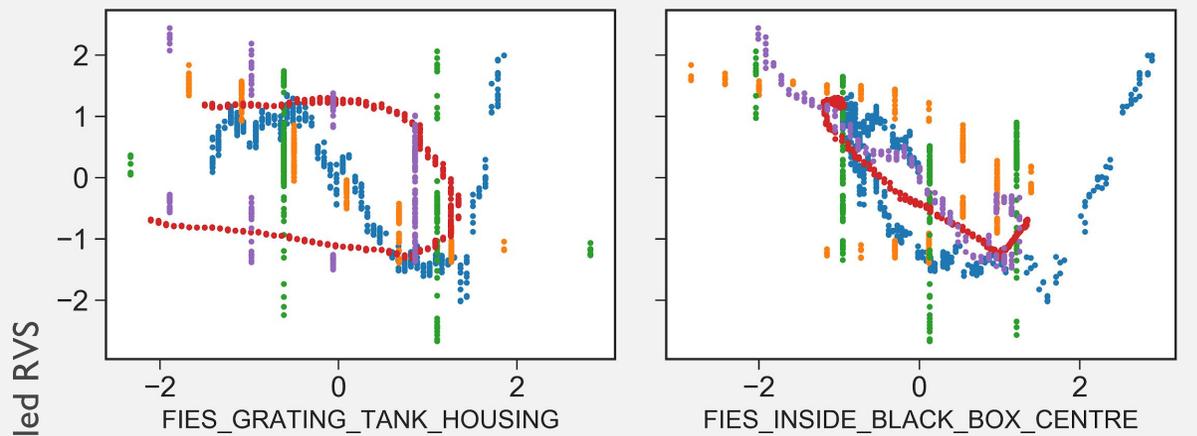
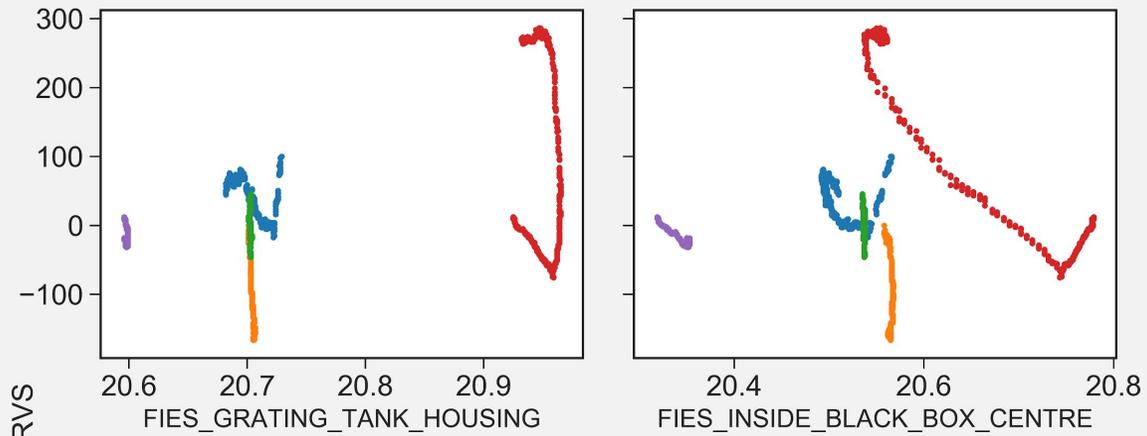
Simple linear regression - Test on timeseries 5



Parameter correlation to RVS

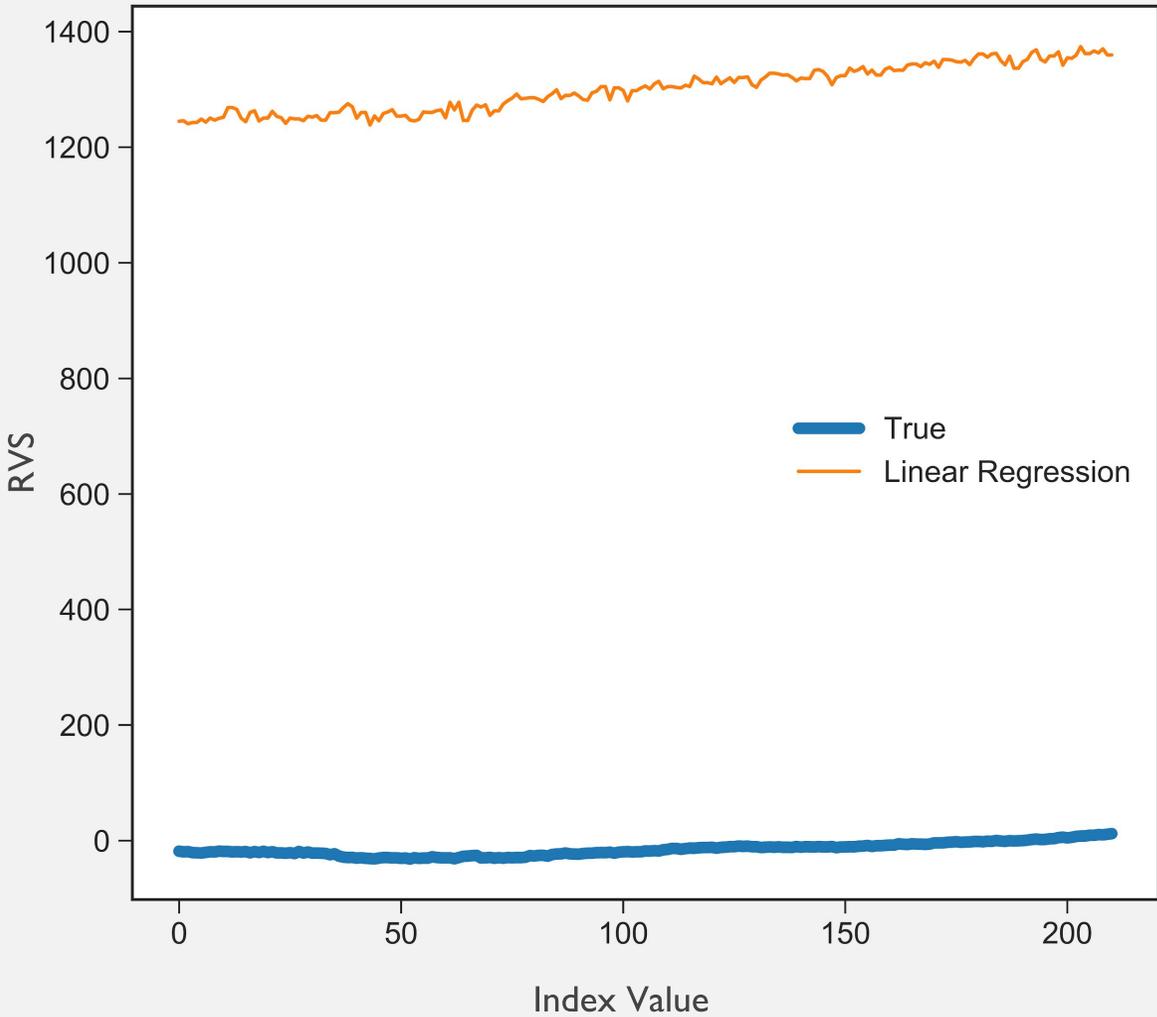
Not scaled

Scaled

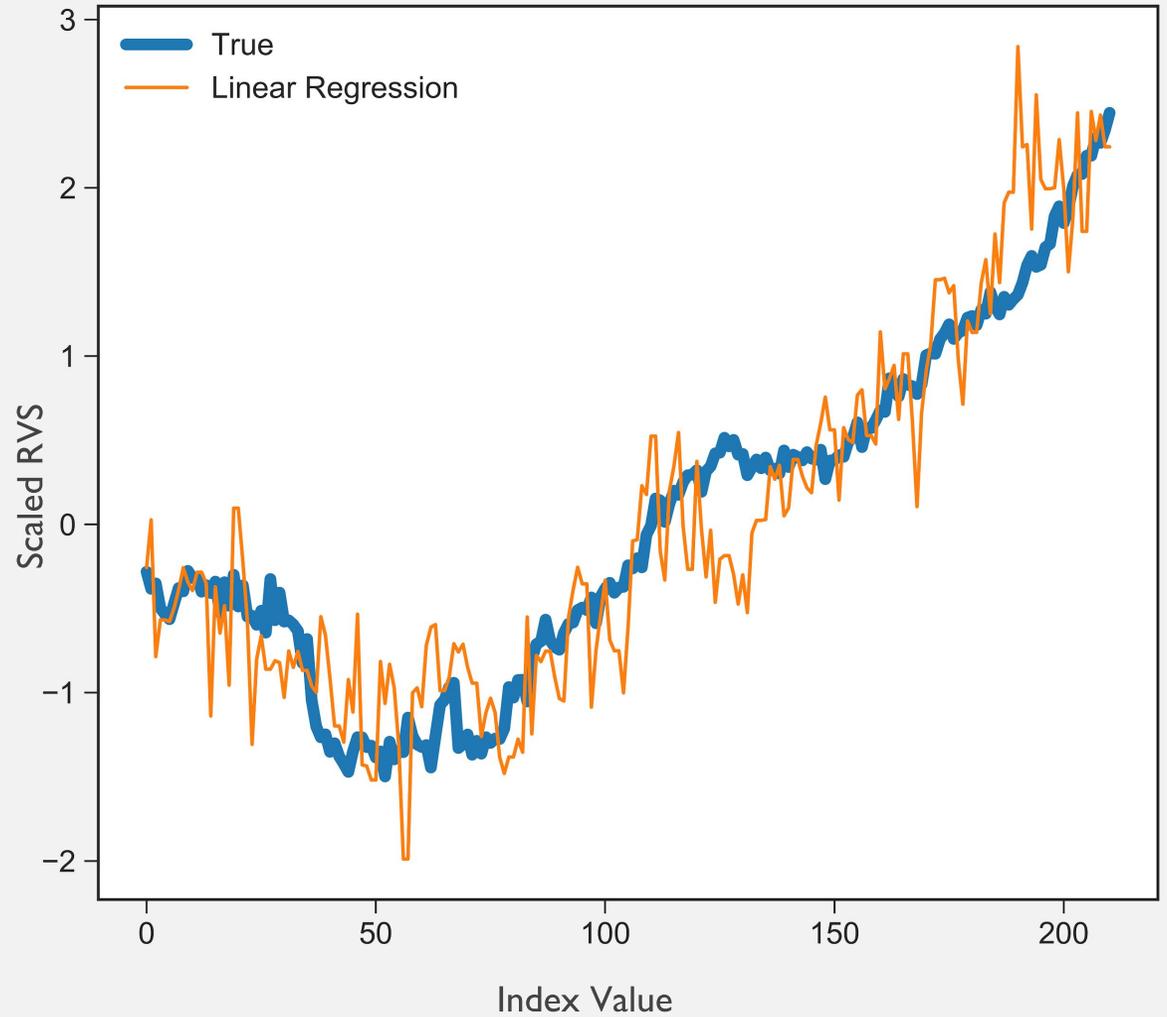


Simple linear regression - Test on timeseries 5

Not scaled

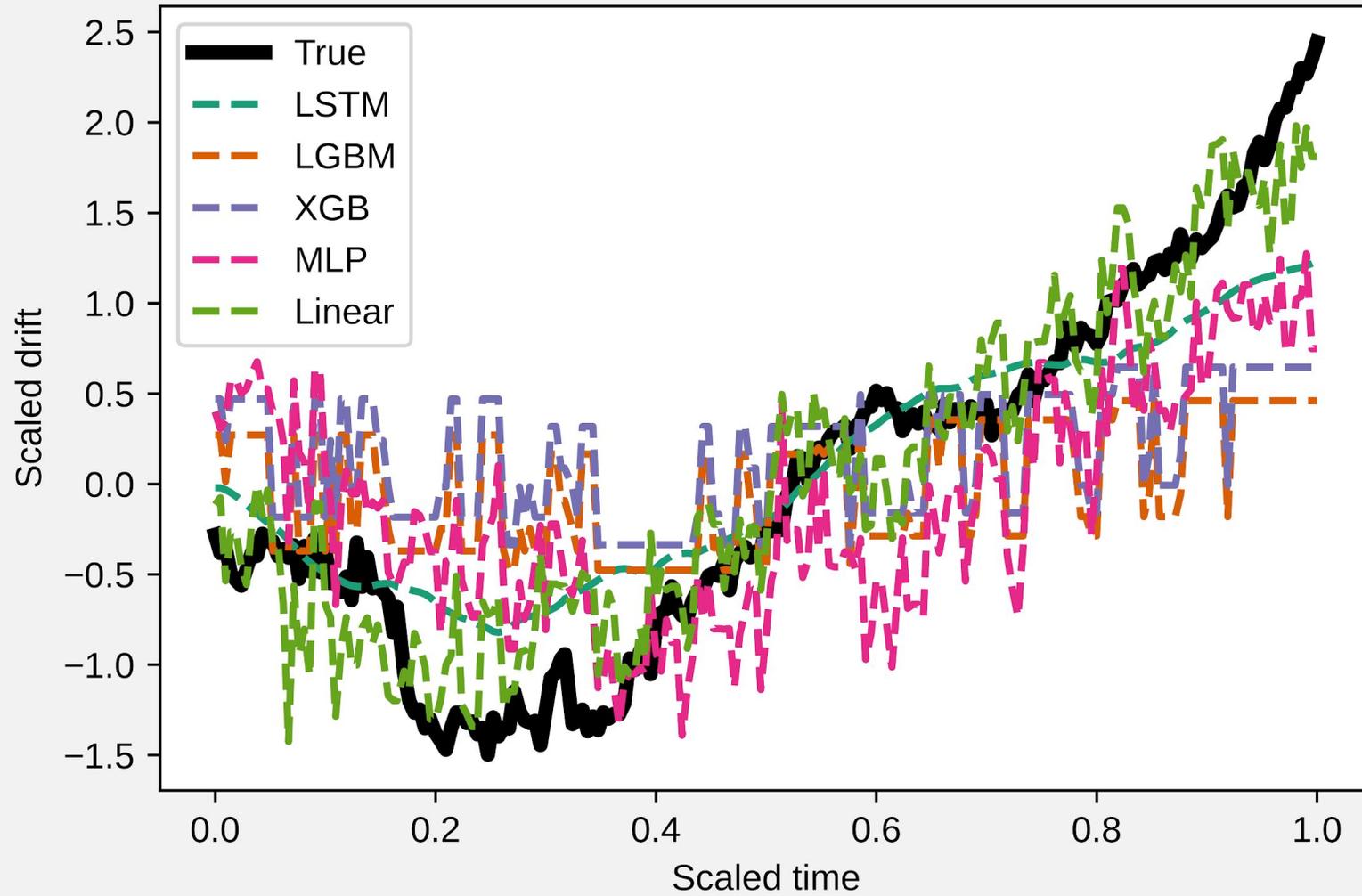


Scaled



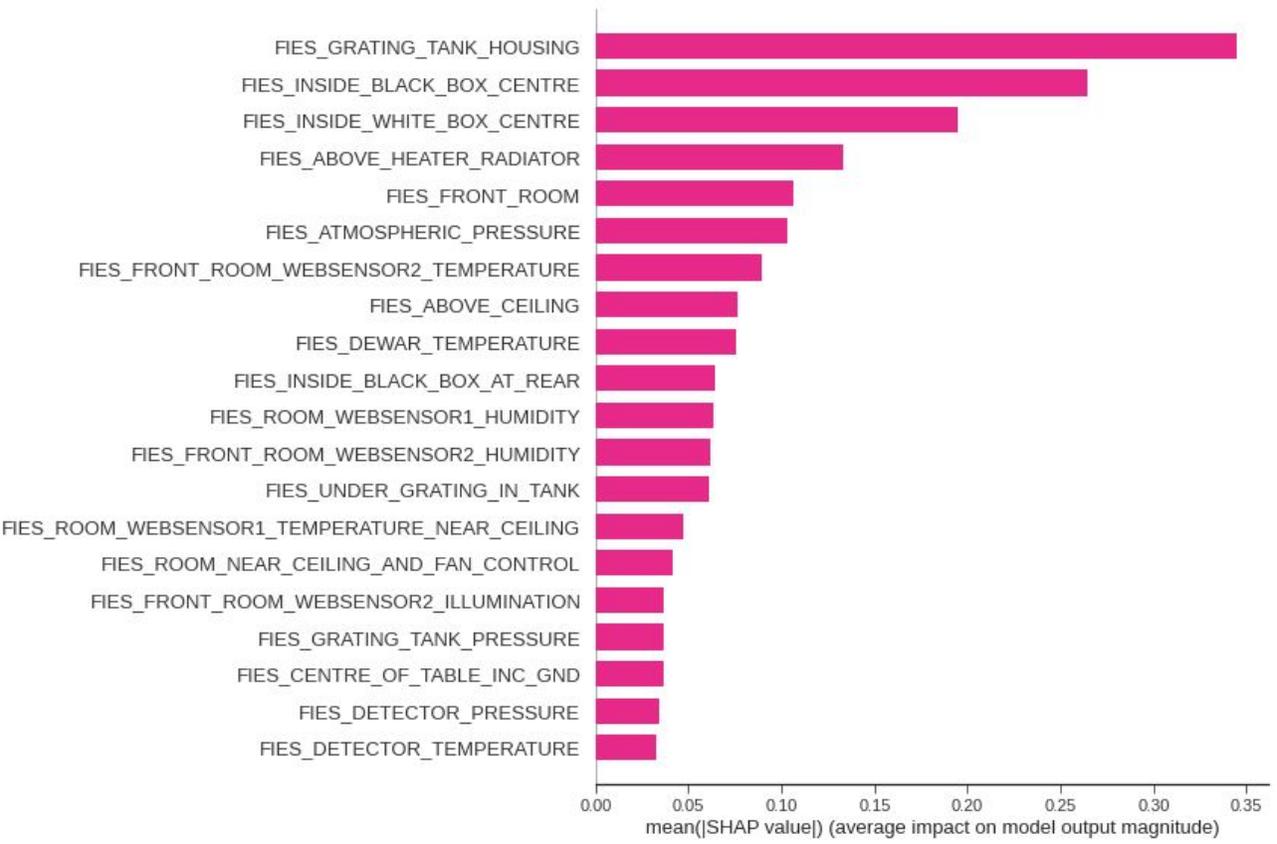
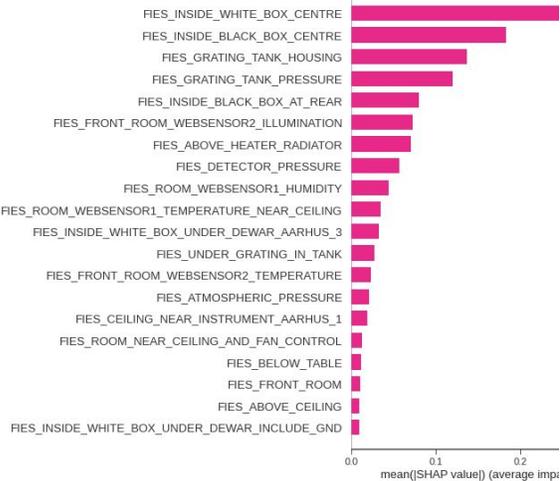
Predictions

Testing on Data5

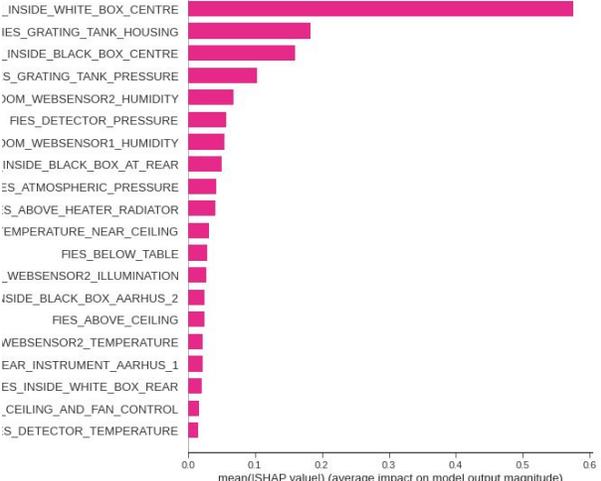


SHAP-values

XG Boost



Light GBM



Coherent with expectations?

“” Thank you for the update!

Our top features:

1. FIES_INSIDE_GRATING_TANK

2. FIES_INSIDE_BLACK_BOX_CENTRE

3. FIES_INSIDE_WHITE_BOX_CENTRE

4. FIES_ABOVE_HEATER_RADIATOR

5. FIES_ATMOSPHERIC_PRESSURE

(...) I can say that we expect that the temperature near the grating, (...), would be a place where we expect large correlation with the drifts (RVs).

FIES_GRATING_TANK_HOUSING is just by the grating, so this makes sense. The other temperatures might also make sense, such as

FIES_INSIDE_BLACK_BOX_CENTRE and FIES_INSIDE_WHITE_BOX_CENTRE.

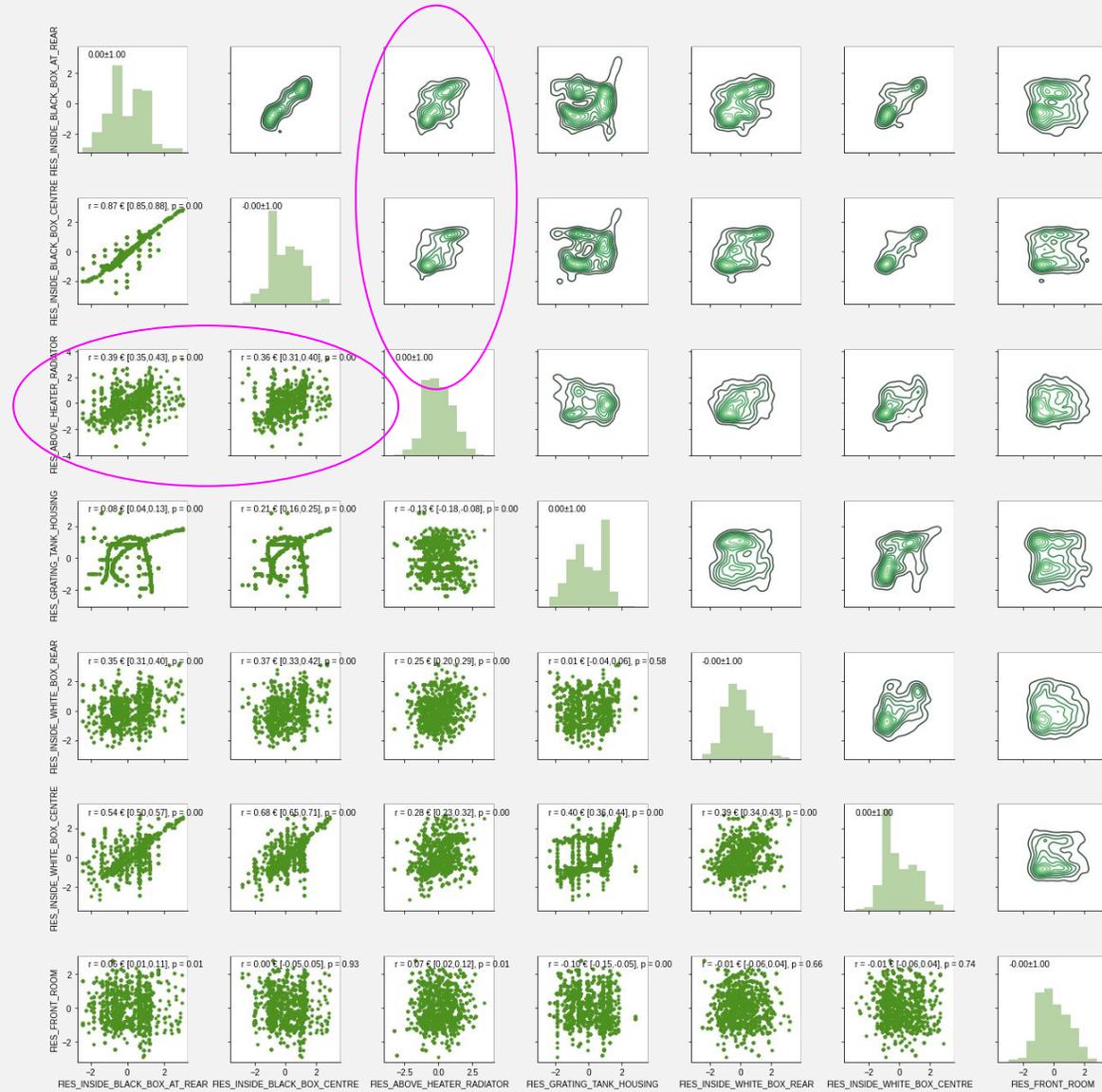
(...) surprised that

FIES_ABOVE_HEATER_RADIATOR has such great impact, given that it is a very noisy instrument.

(...) On the other hand, it surprises me that

FIES_ATMOSPHERIC_PRESSURE has any influence (...) in a pressure chamber (...) as we have a very expensive barometer!“”

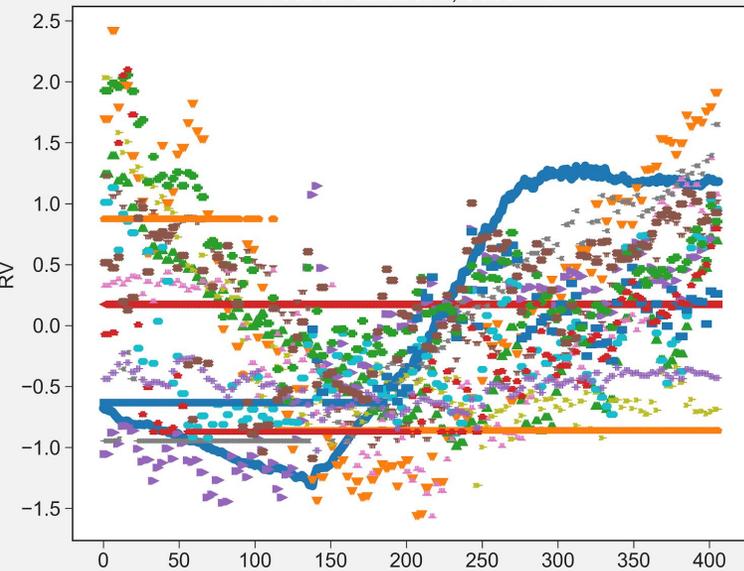
Telemetry data correlations



Importance of data sample size

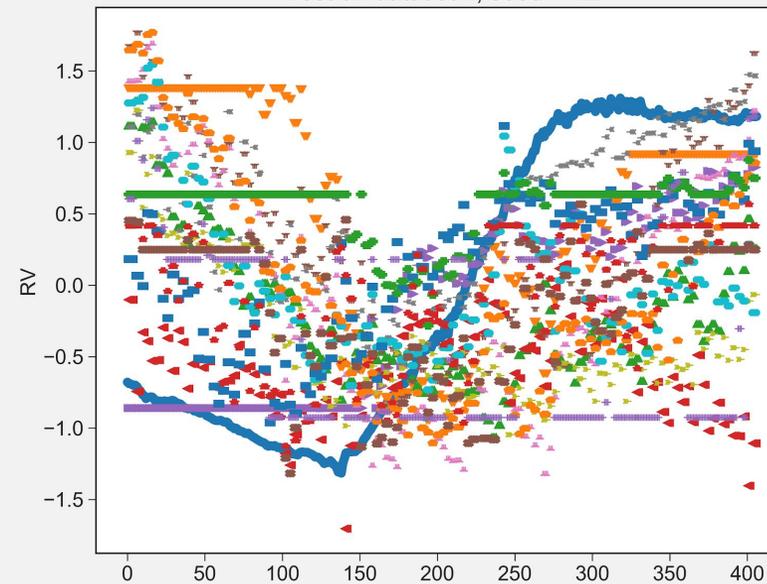
MLPR NN optimized with random search

Test on dataset 4, seed = 0



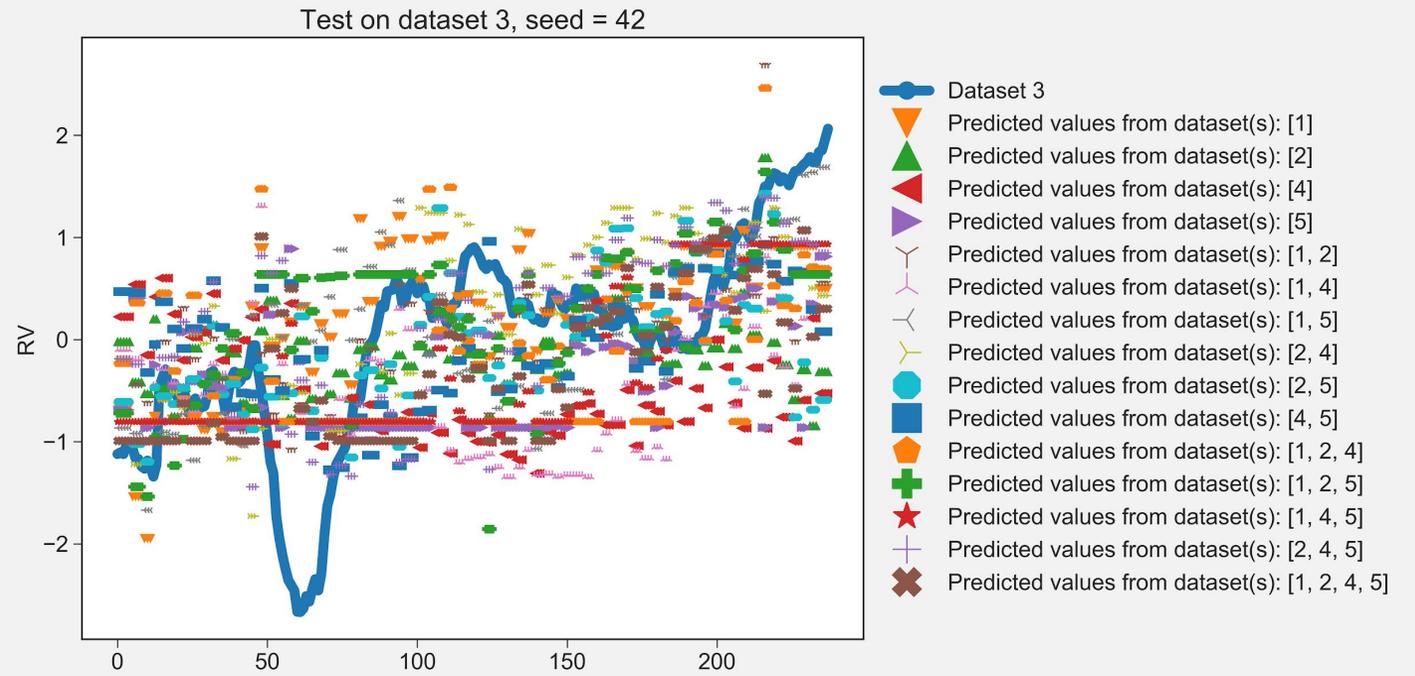
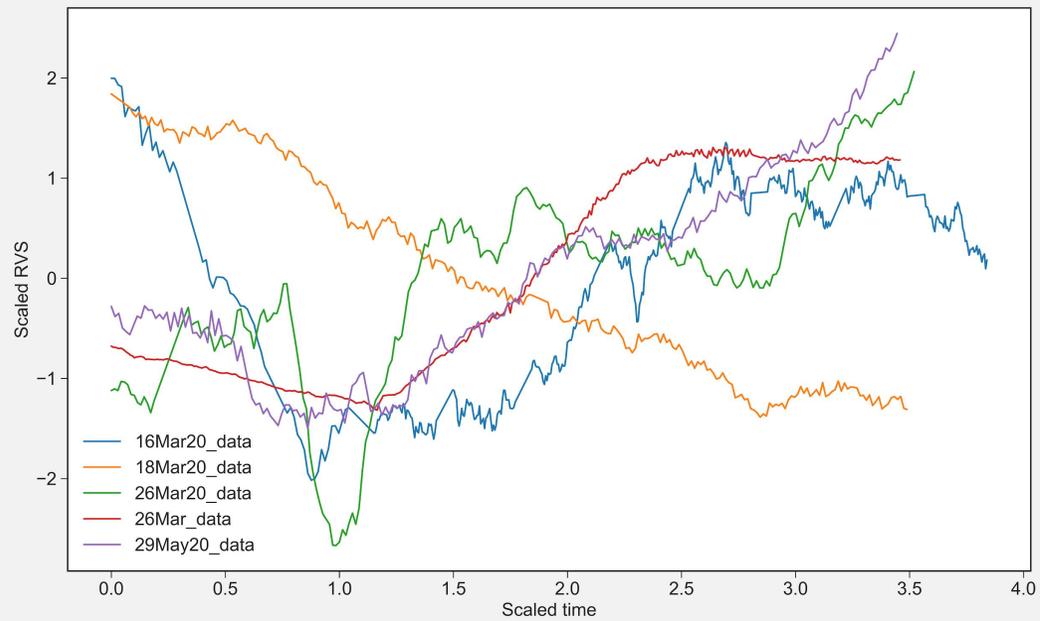
- Dataset 4
- Predicted values from dataset(s): [1]
- Predicted values from dataset(s): [2]
- Predicted values from dataset(s): [3]
- Predicted values from dataset(s): [5]
- Predicted values from dataset(s): [1, 2]
- Predicted values from dataset(s): [1, 3]
- Predicted values from dataset(s): [1, 5]
- Predicted values from dataset(s): [2, 3]
- Predicted values from dataset(s): [2, 5]
- Predicted values from dataset(s): [3, 5]
- Predicted values from dataset(s): [1, 2, 3]
- Predicted values from dataset(s): [1, 2, 5]
- Predicted values from dataset(s): [1, 3, 5]
- Predicted values from dataset(s): [2, 3, 5]
- Predicted values from dataset(s): [1, 2, 3, 5]

Test on dataset 4, seed = 42



- Dataset 4
- Predicted values from dataset(s): [1]
- Predicted values from dataset(s): [2]
- Predicted values from dataset(s): [3]
- Predicted values from dataset(s): [5]
- Predicted values from dataset(s): [1, 2]
- Predicted values from dataset(s): [1, 3]
- Predicted values from dataset(s): [1, 5]
- Predicted values from dataset(s): [2, 3]
- Predicted values from dataset(s): [2, 5]
- Predicted values from dataset(s): [3, 5]
- Predicted values from dataset(s): [1, 2, 3]
- Predicted values from dataset(s): [1, 2, 5]
- Predicted values from dataset(s): [1, 3, 5]
- Predicted values from dataset(s): [2, 3, 5]
- Predicted values from dataset(s): [1, 2, 3, 5]

Representative Data



MLPR NN optimized with random search

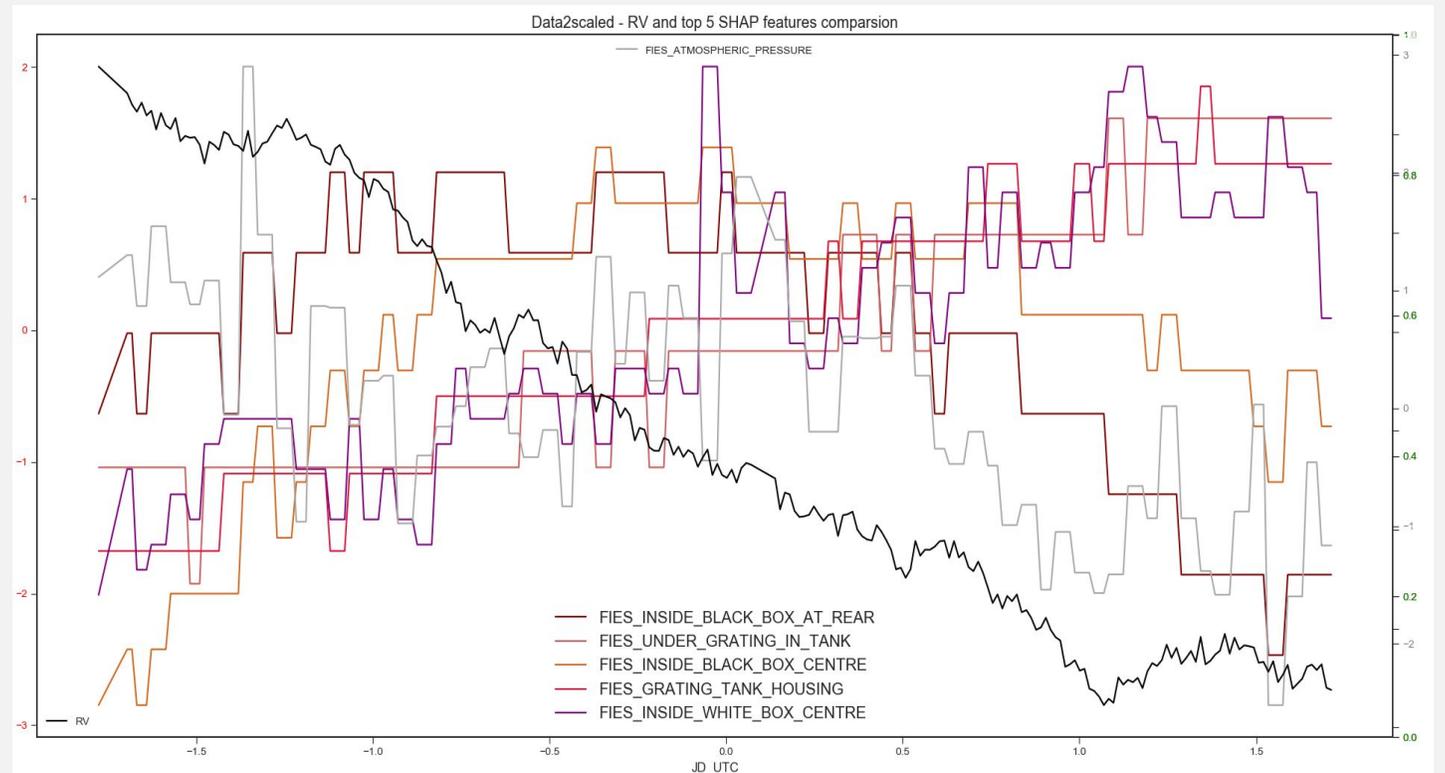
Conclusion

Data sample size makes predictions difficult

Confidence of ranked variables

Expectations compared to results

Further work



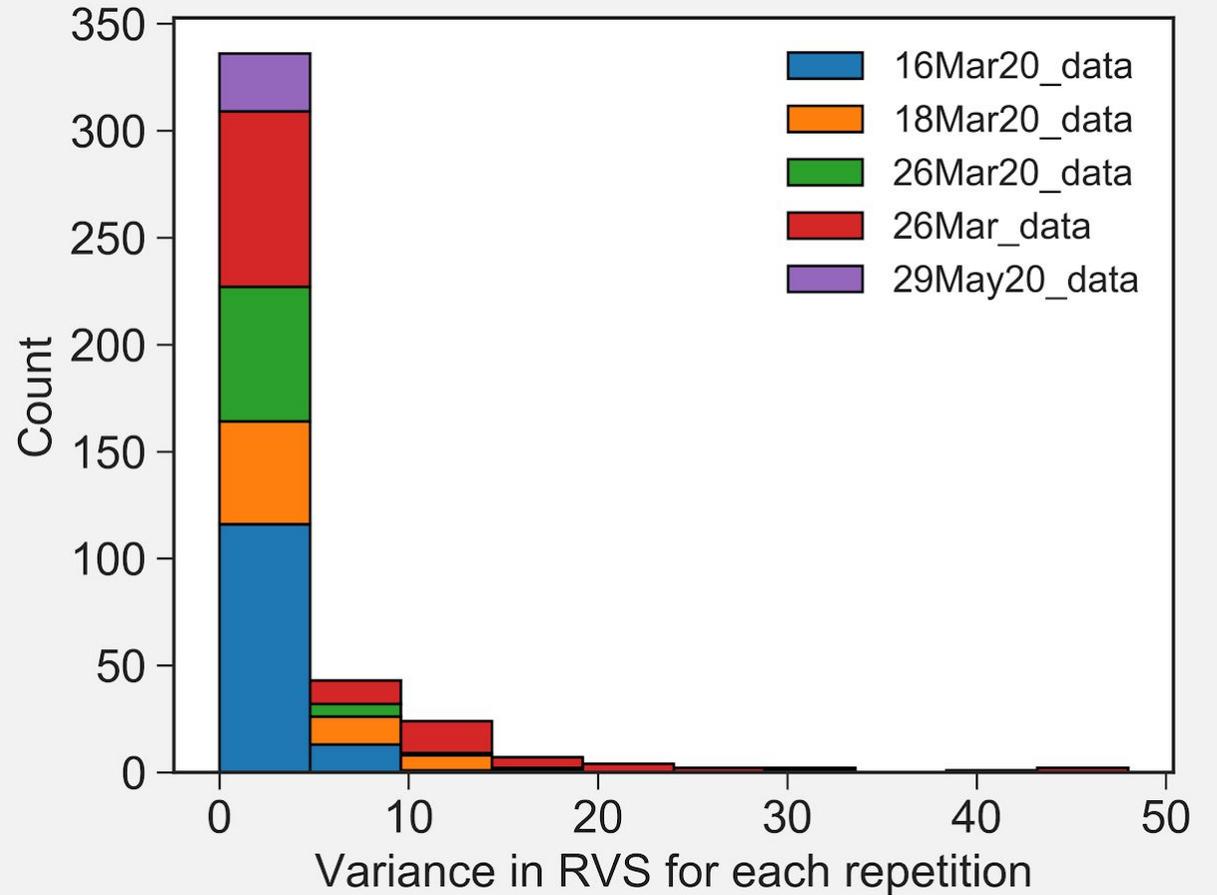
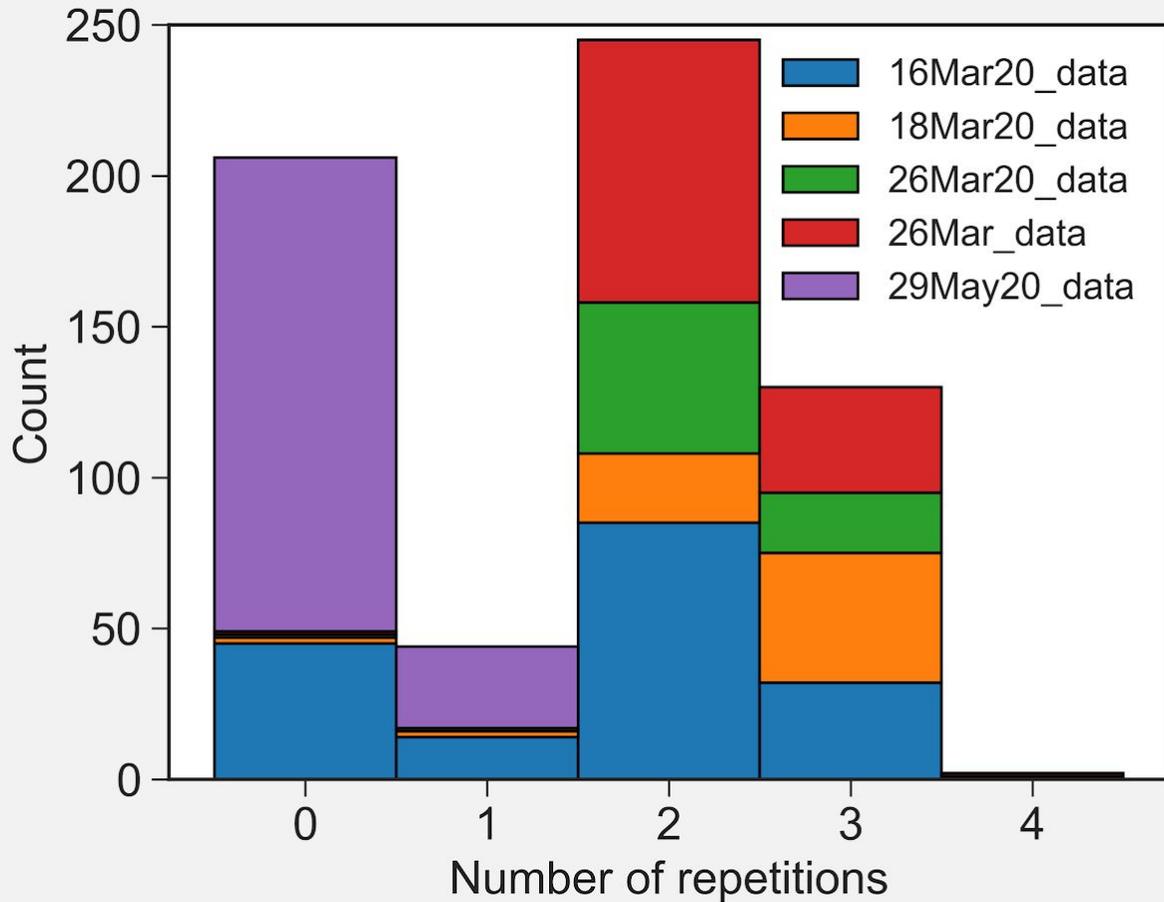
Appendix

The scripts behind the scenes

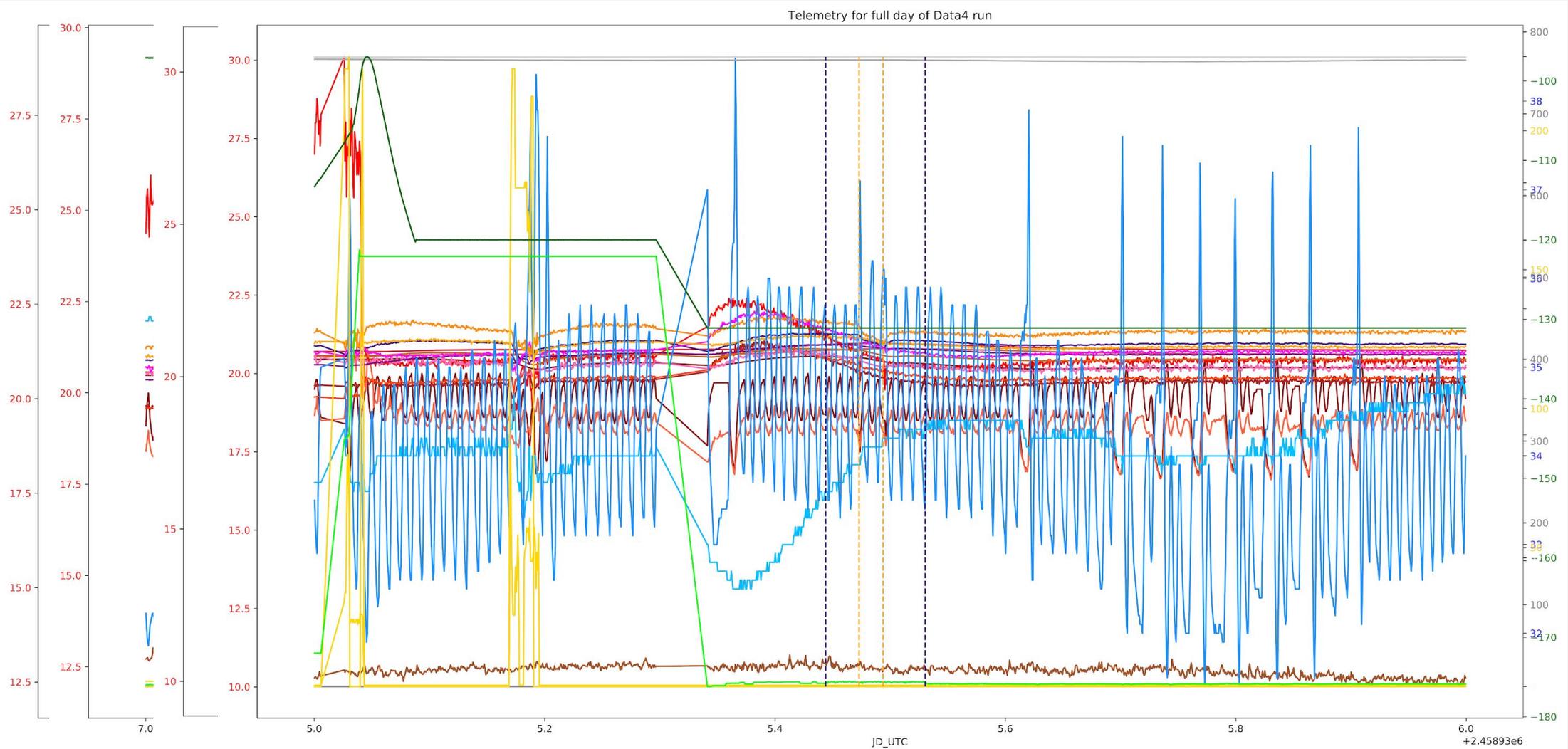
To see all the code we wrote in order to get the presented results (and all the attempts that didn't produce any results), please refer to [our shared GitHub repo](#)

To run the LSTM efficiently, we have parallelized the scripts including this algorithm to a Google Colab GPU. An example is found [here](#)

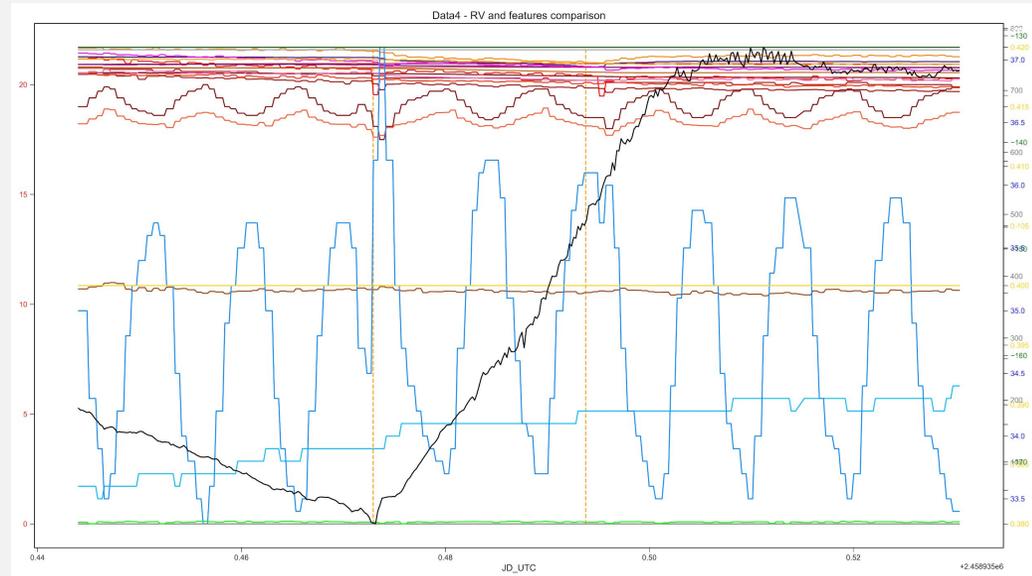
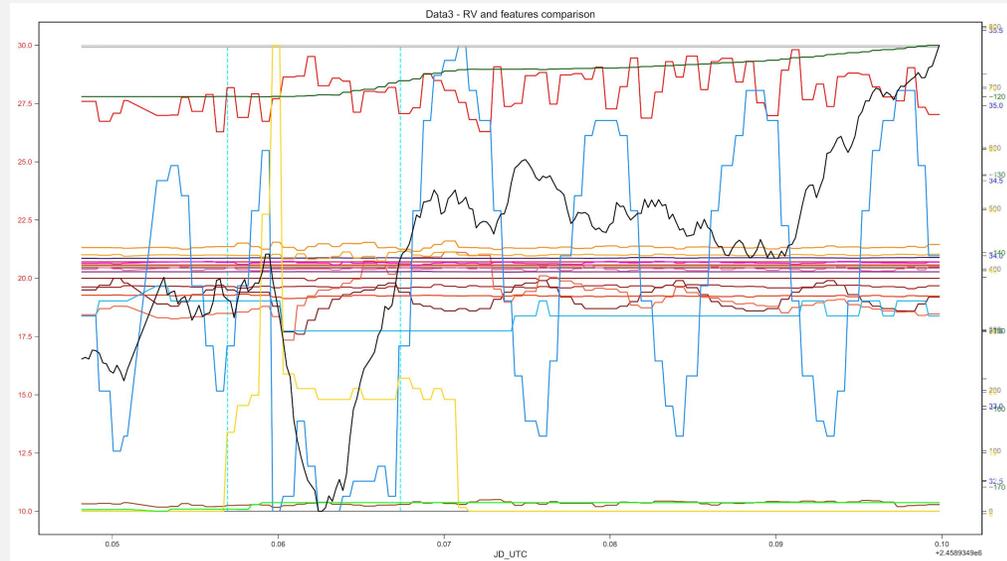
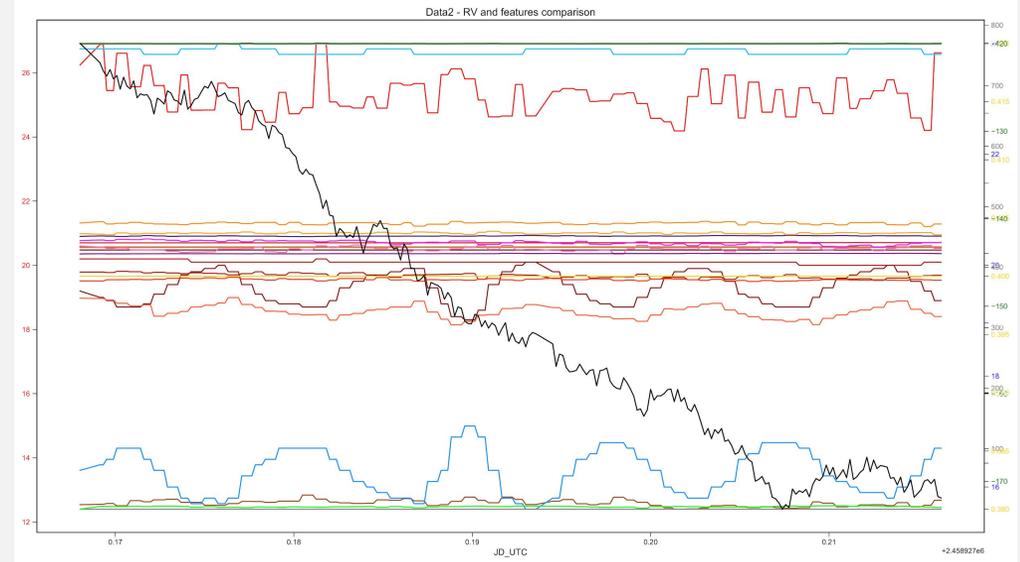
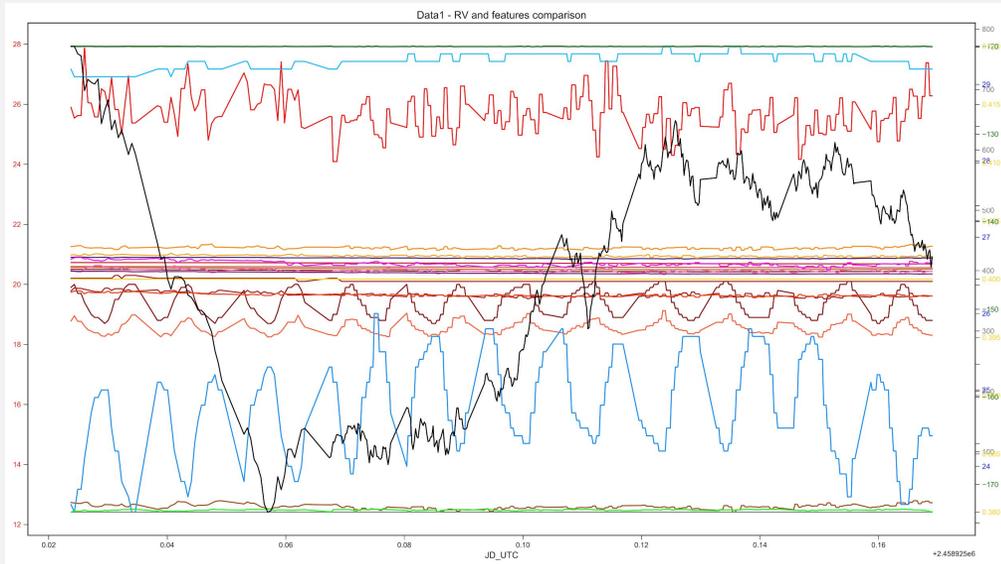
Number of telemetry copies (repetitions)



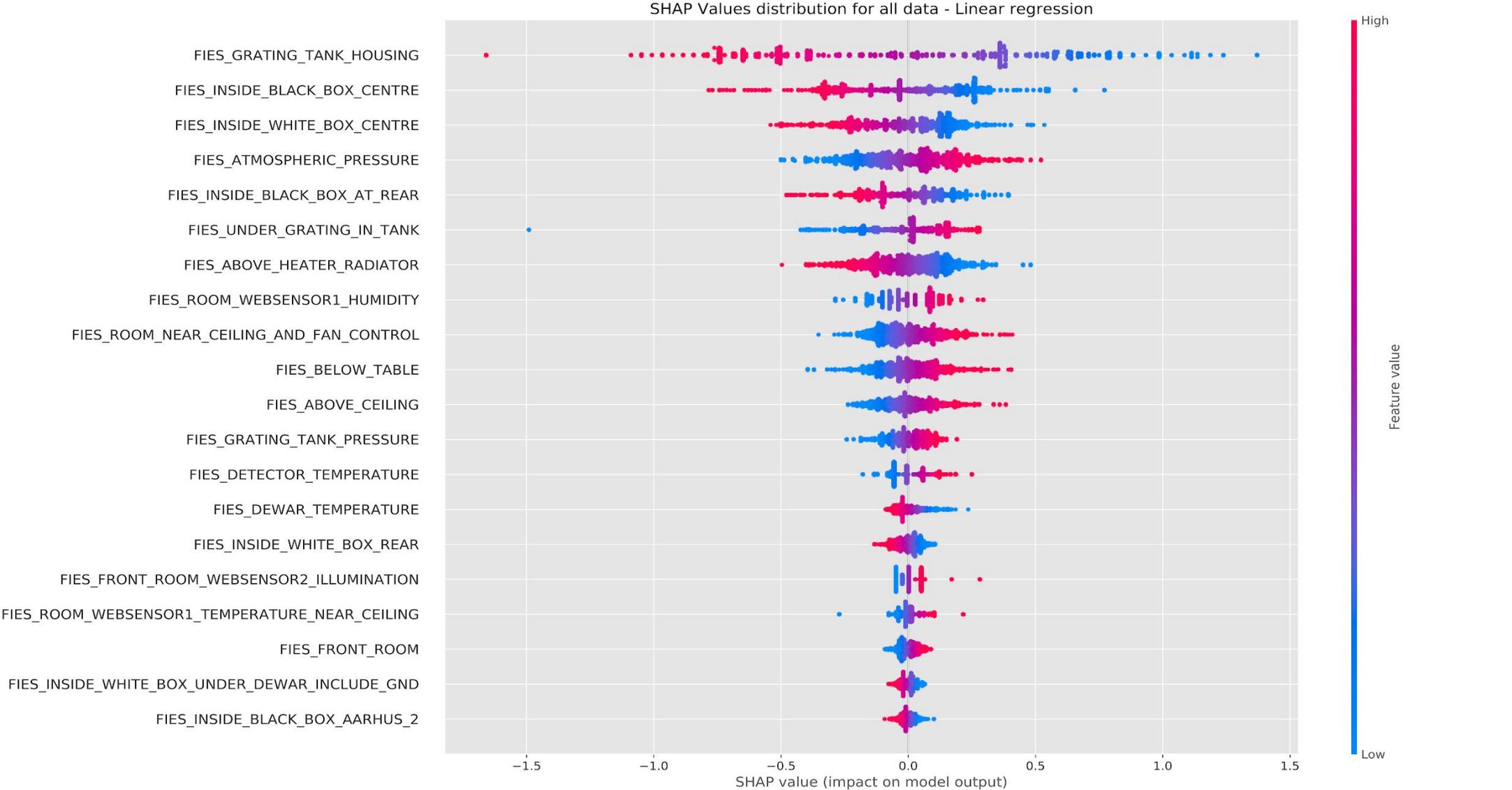
Visualization of the telemetry data



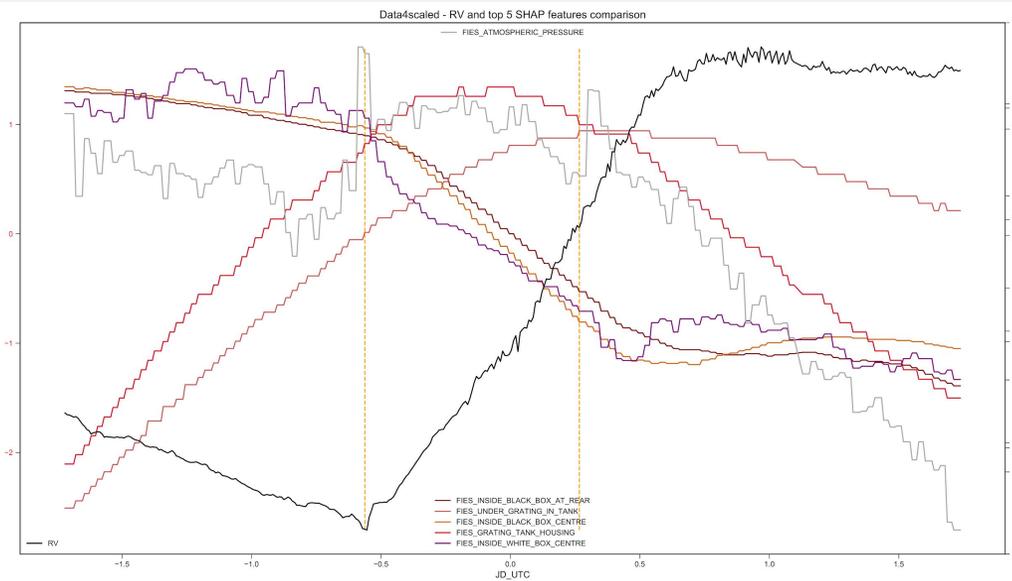
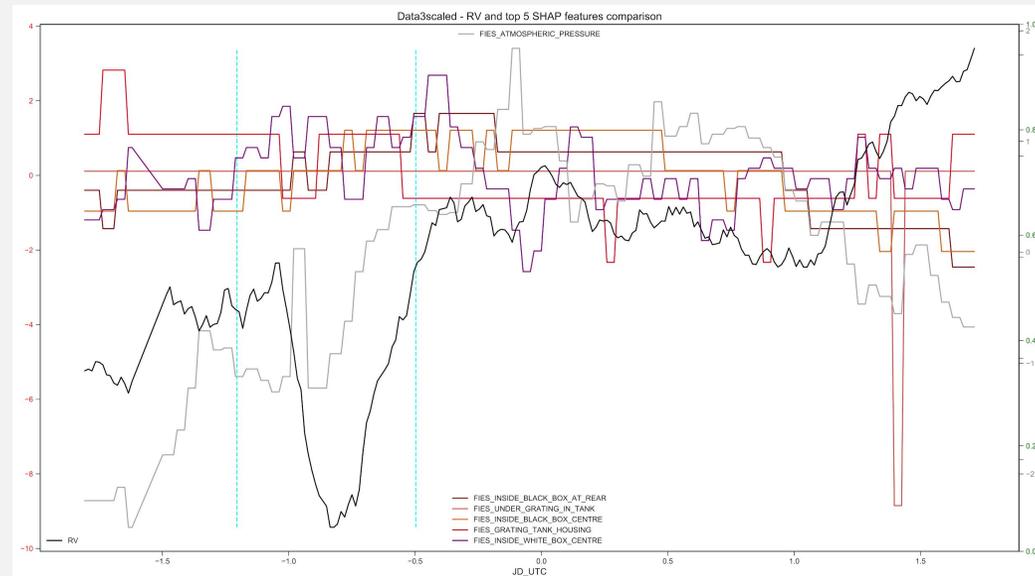
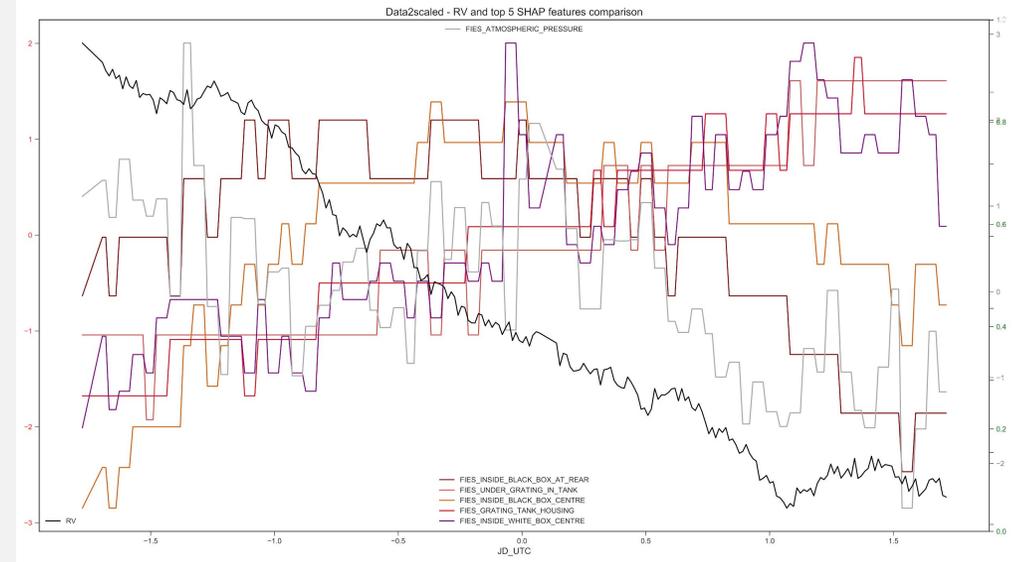
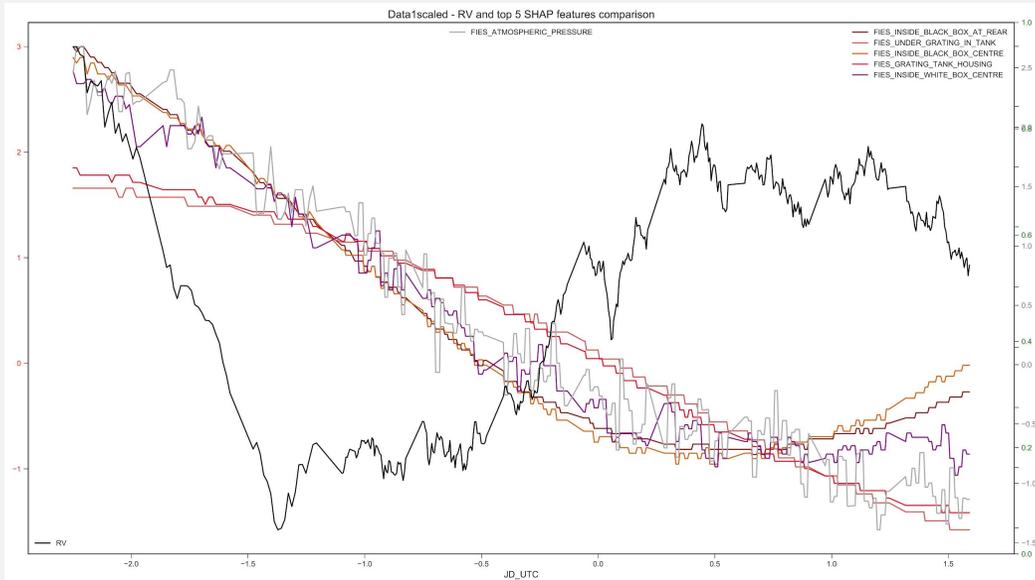
Visualization of the telemetry with RV data



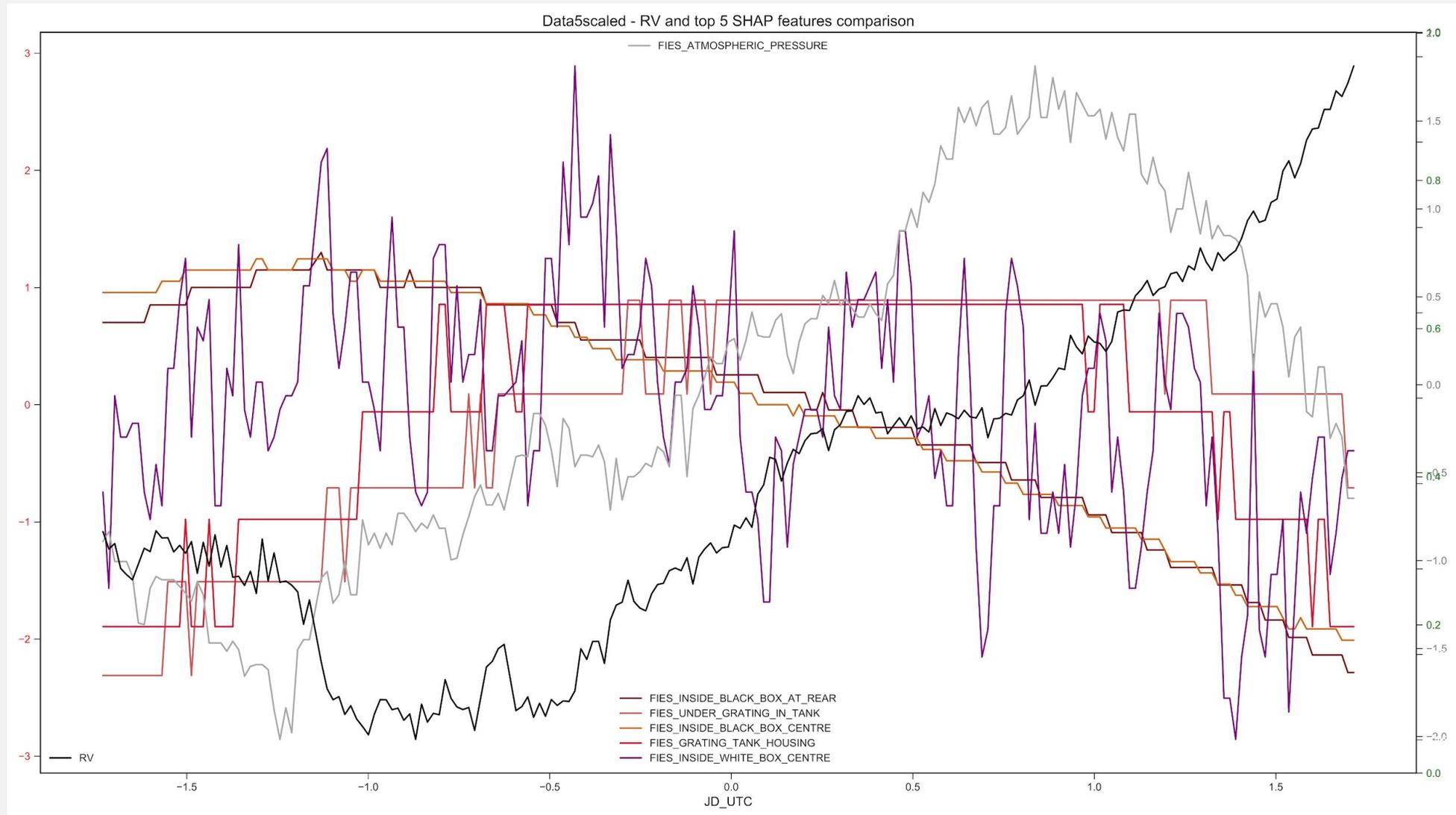
SHAP-values



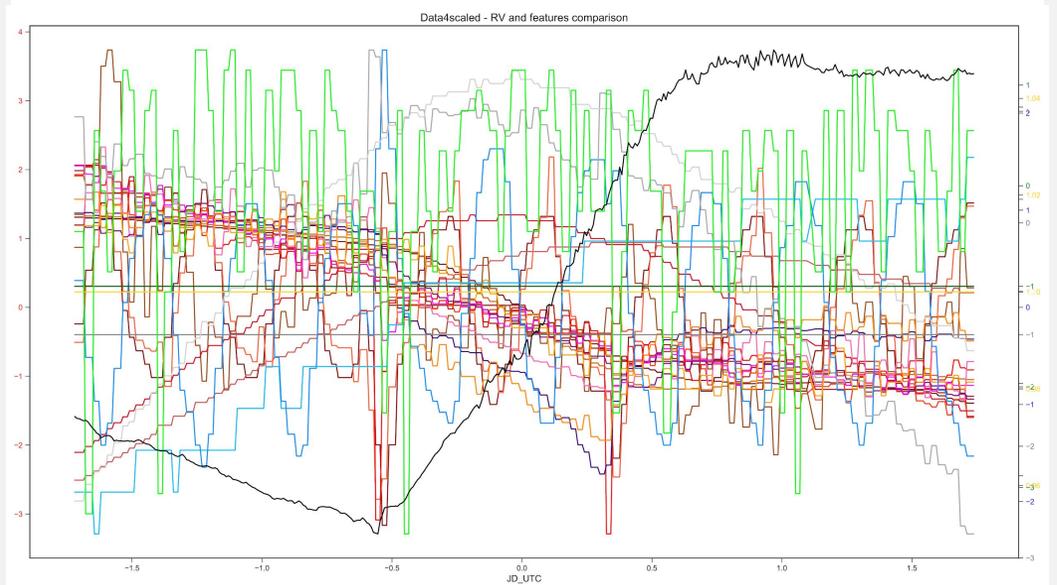
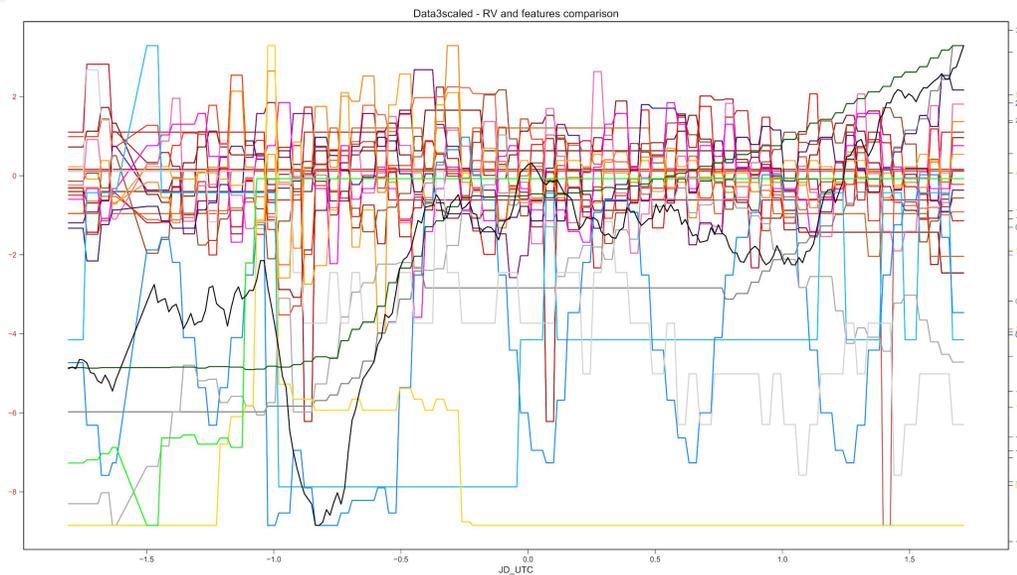
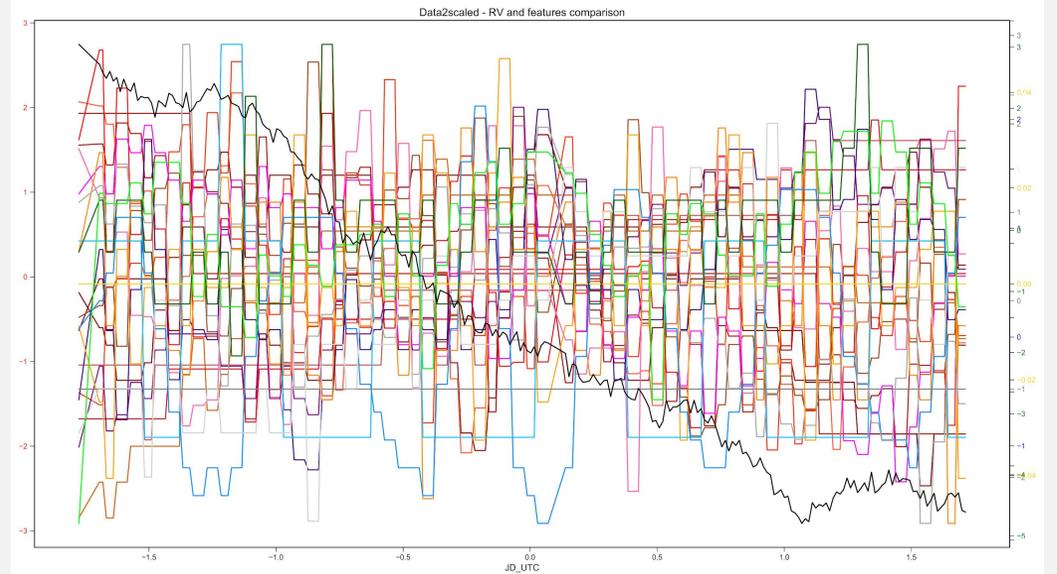
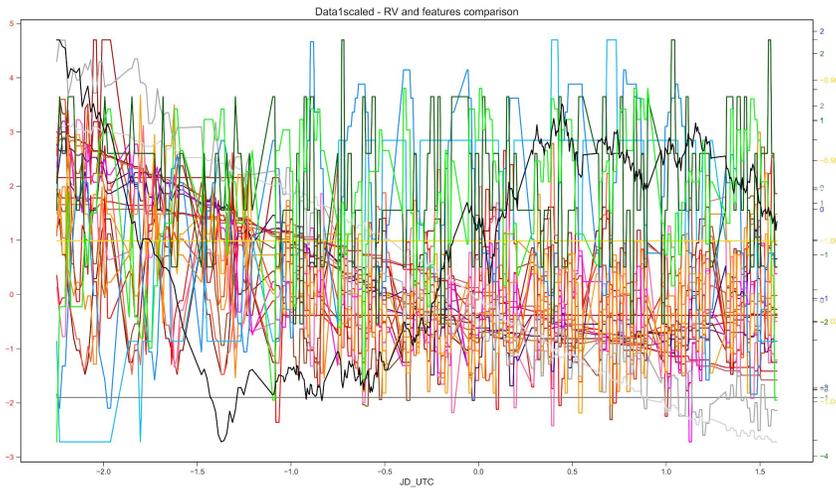
Visualization of the scaled telemetry and RV - Top 5 SHAP features



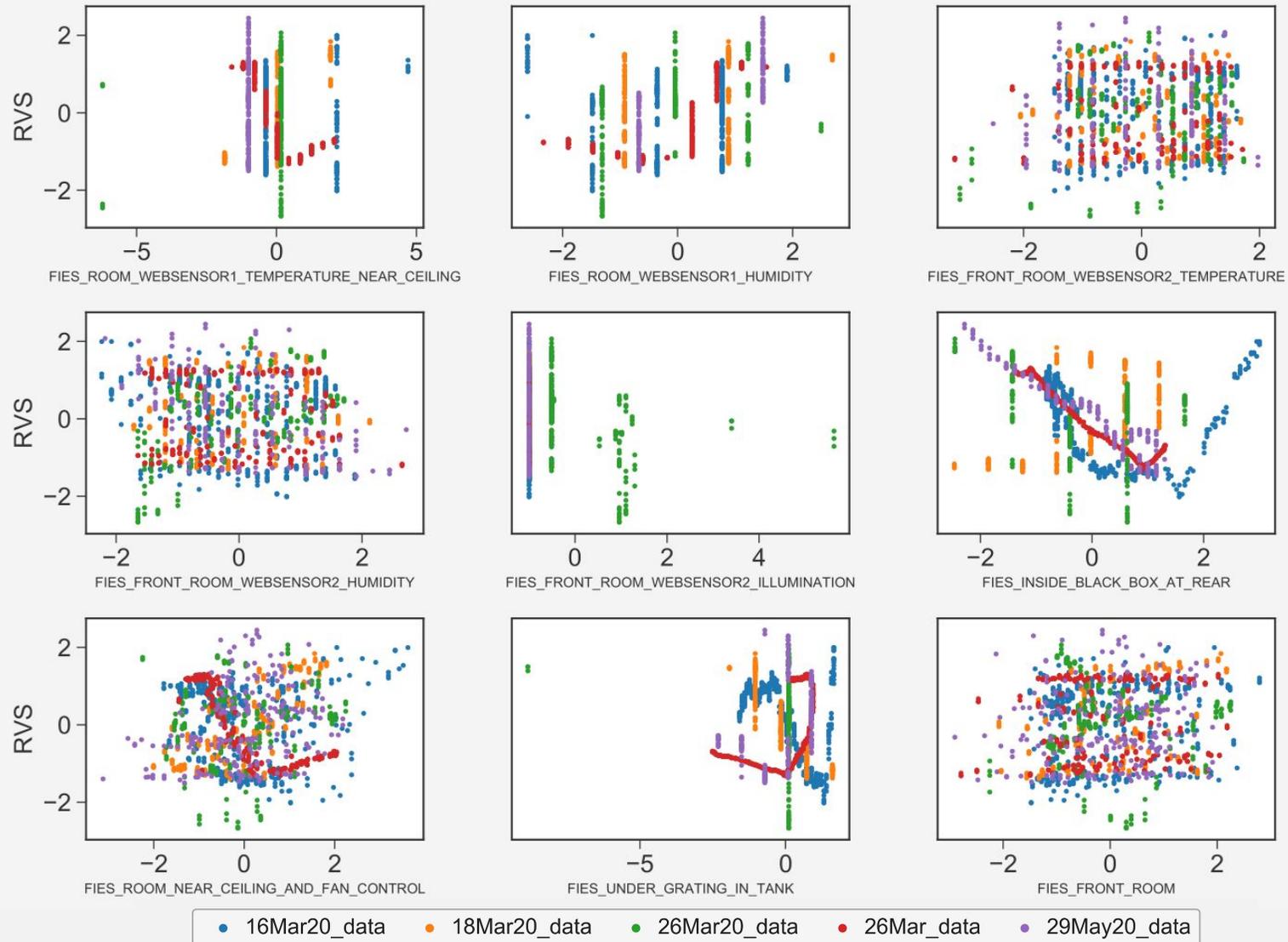
Visualization of the scaled telemetry and RV - Top 5 SHAP features



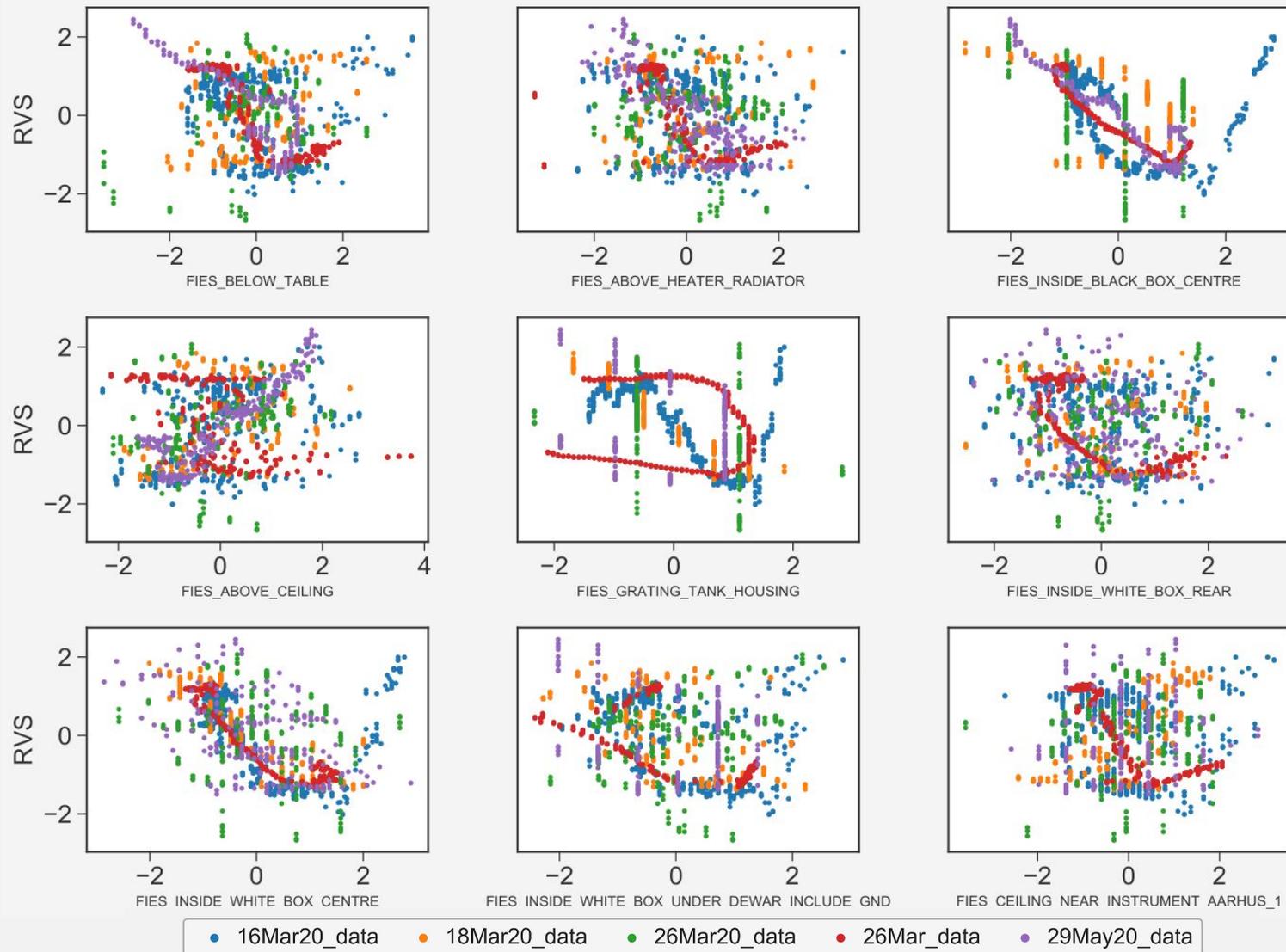
Visualization of the scaled telemetry with scaled RV data



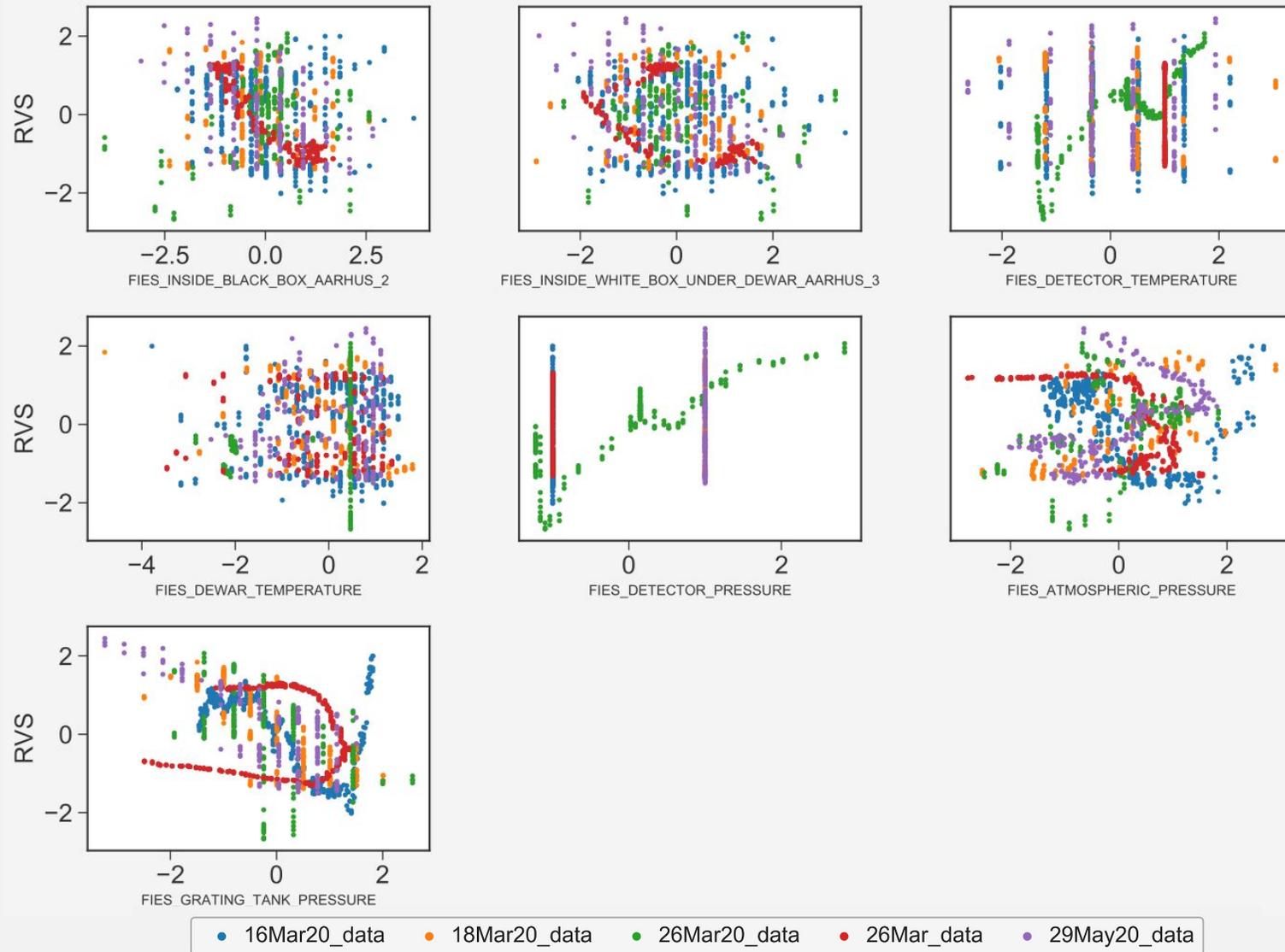
Correlations between RVS and parameters (Scaled) - Part 1/3



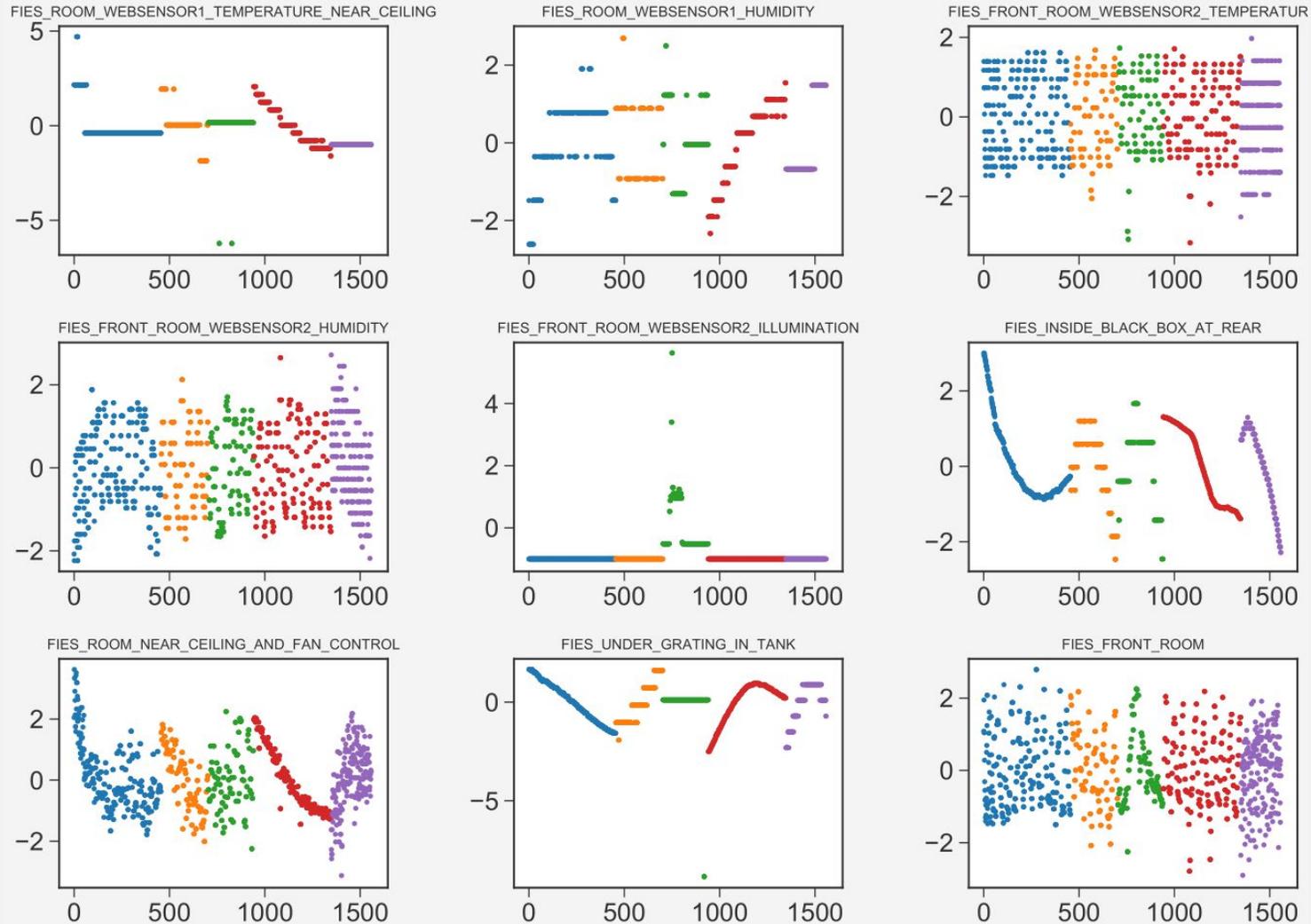
Correlations between RVS and parameters (Scaled) - Part 2/3



Correlations between RVS and parameters (Scaled) - Part 3/3

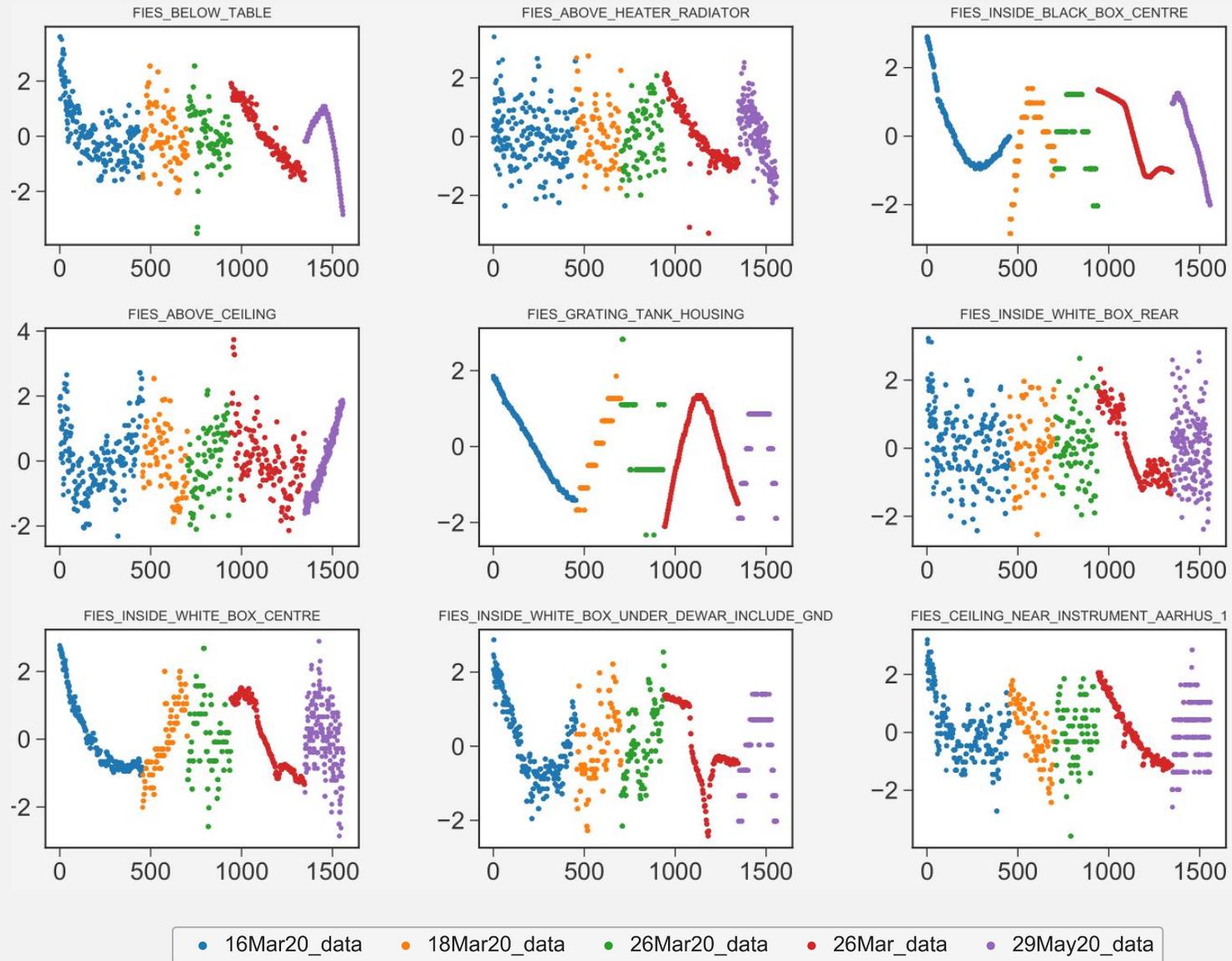


RVS vs. Time (Scaled) - Part 1/3

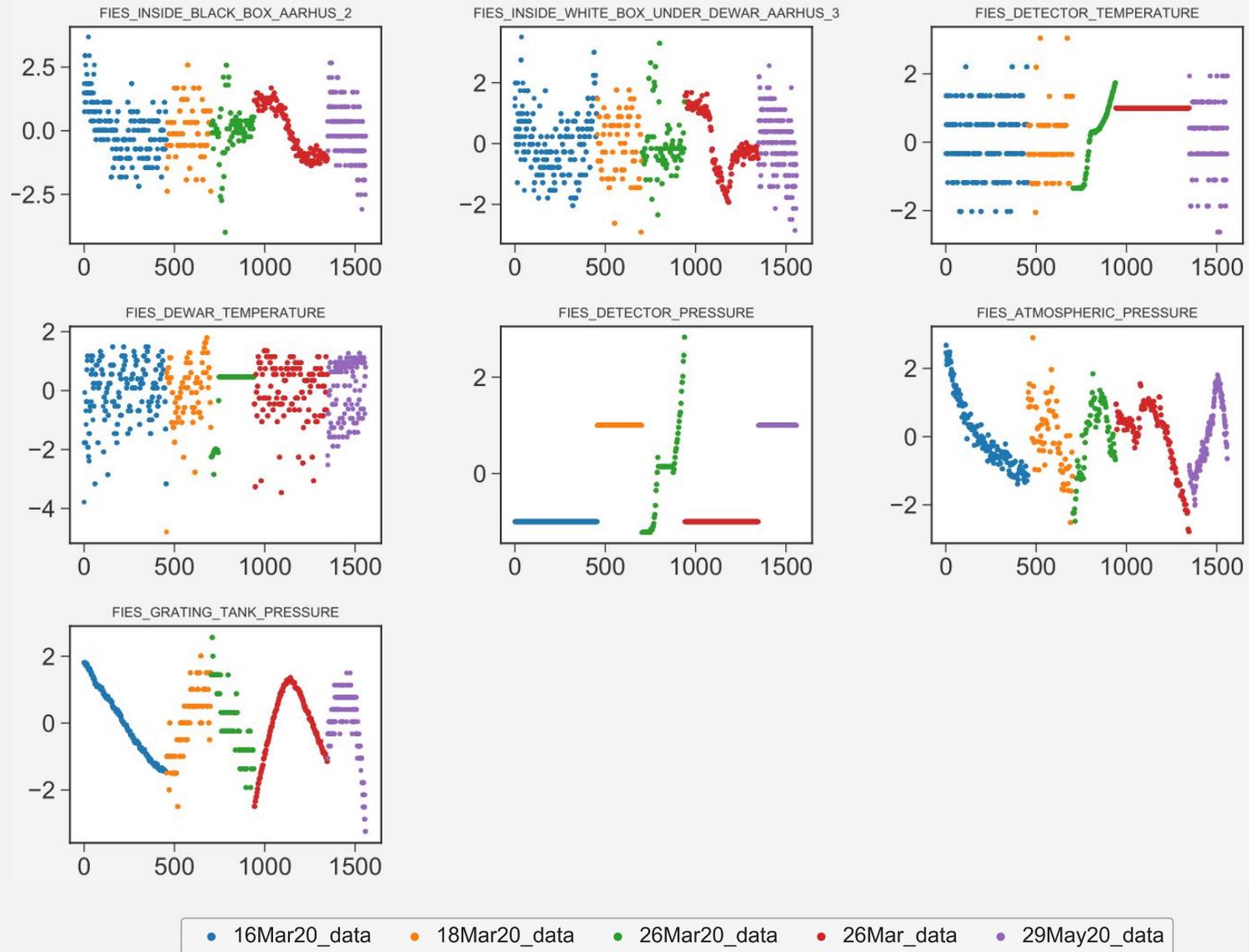


● 16Mar20_data ● 18Mar20_data ● 26Mar20_data ● 26Mar_data ● 29May20_data

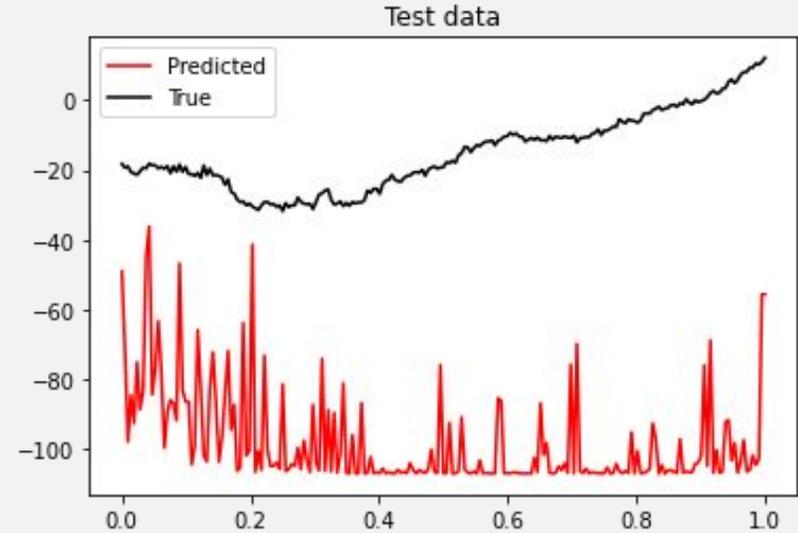
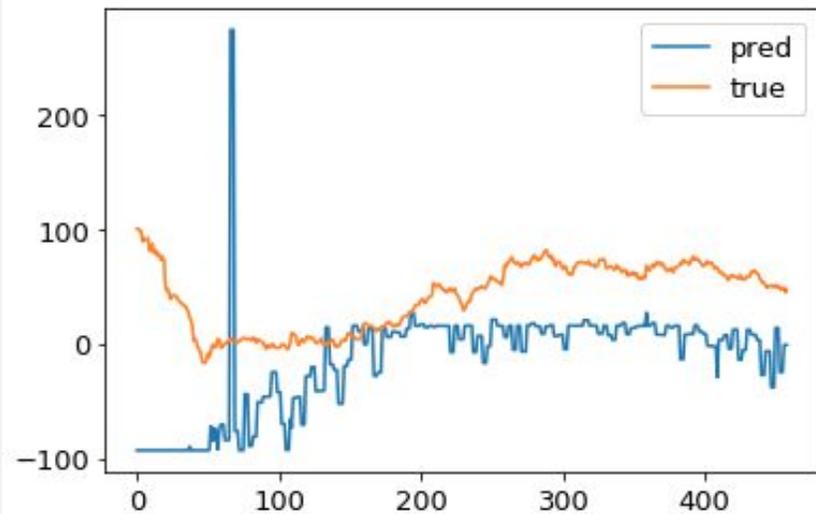
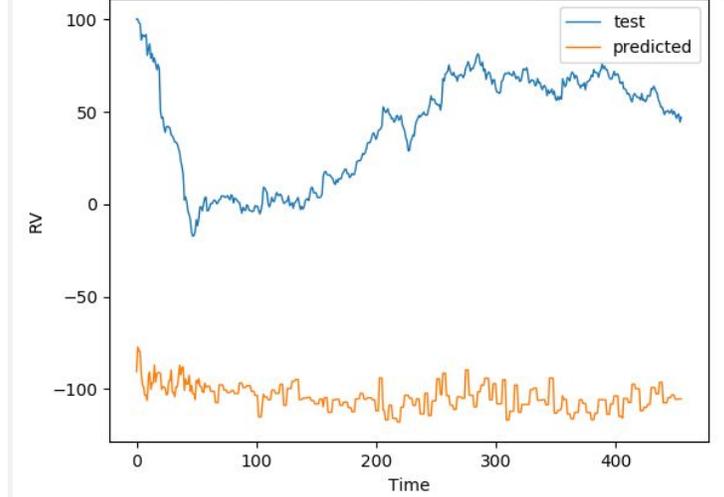
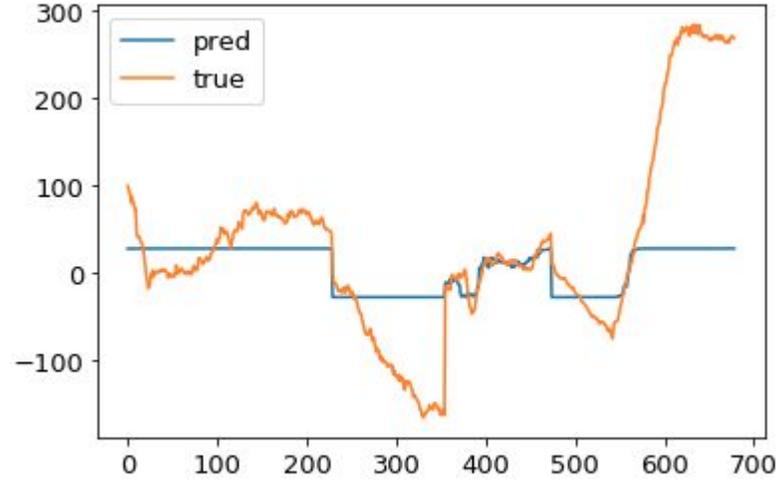
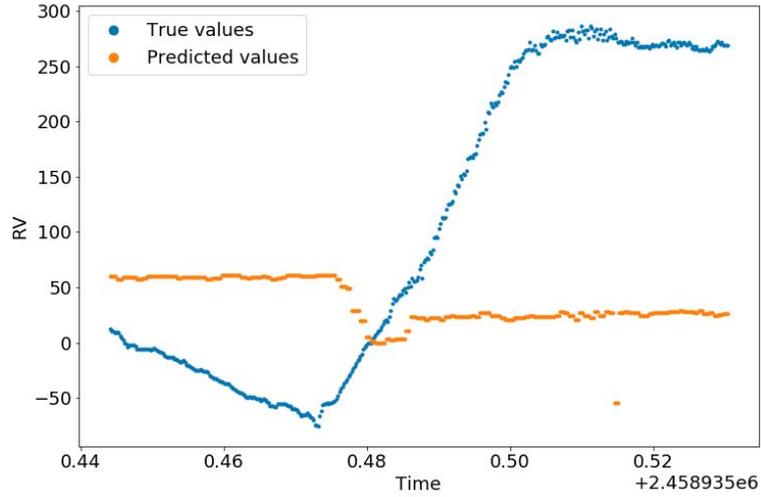
RVS vs. Time (Scaled) - Part 2/3



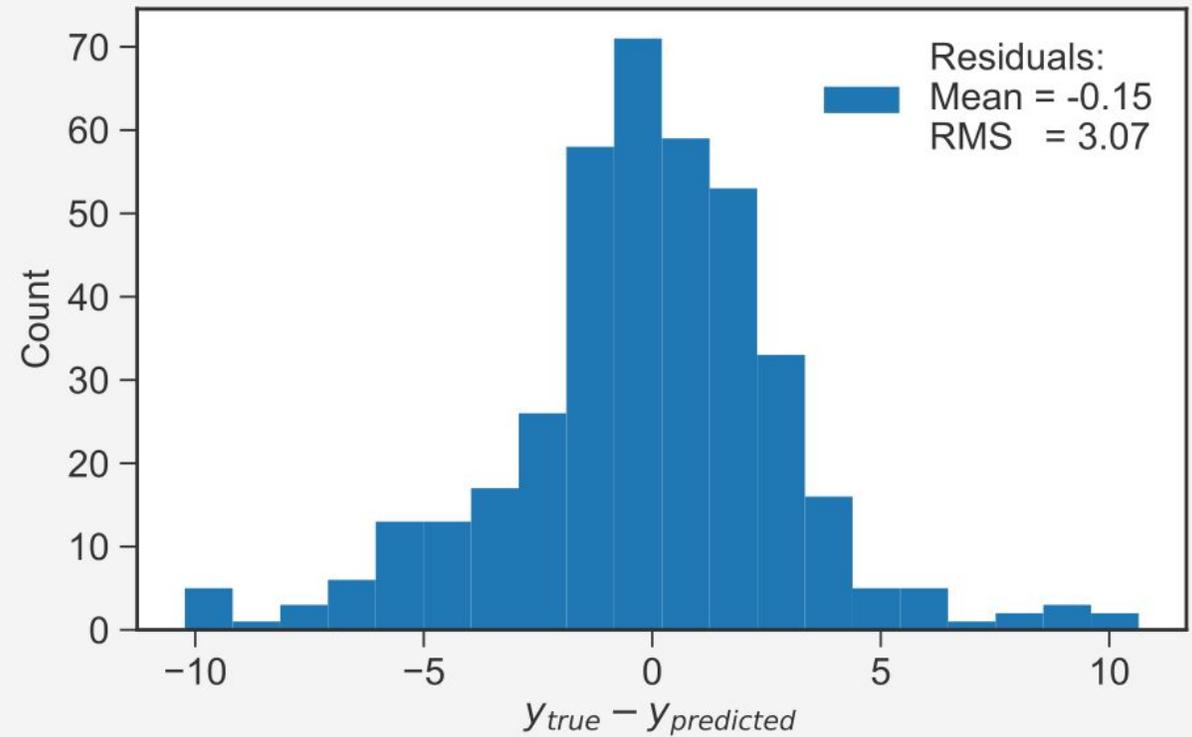
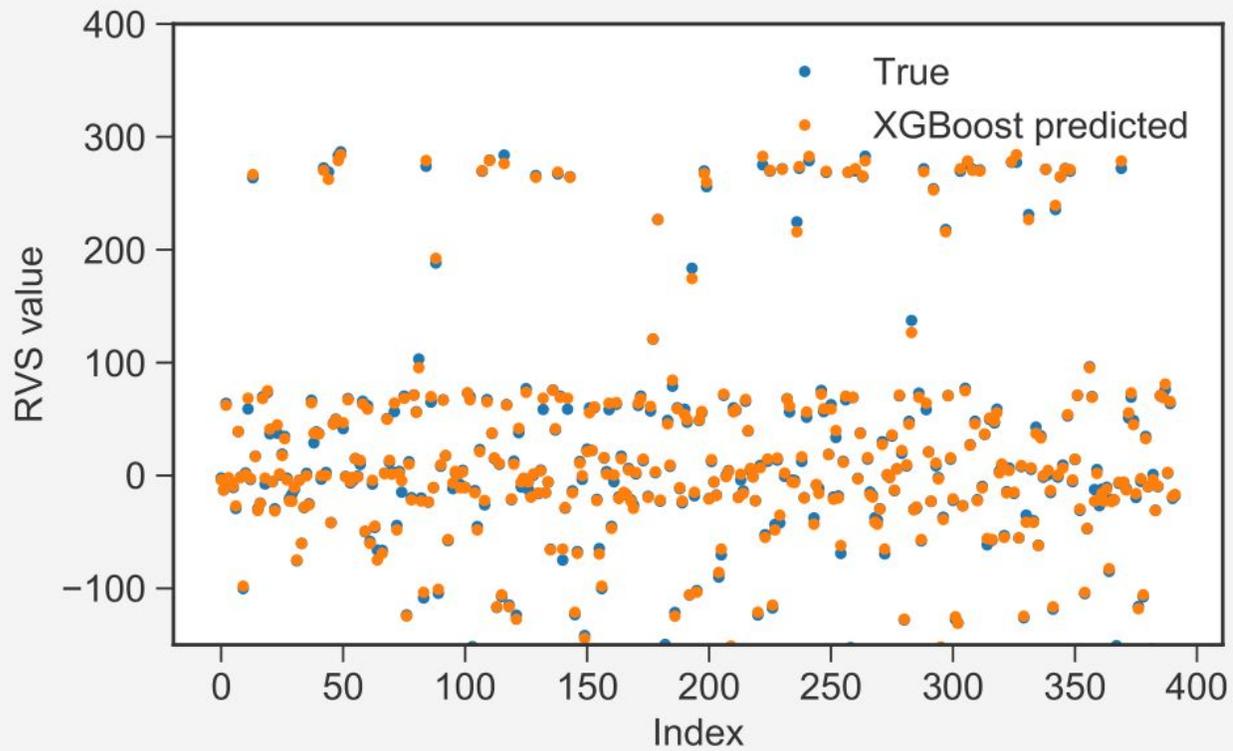
RVS vs. Time (Scaled) - Part 3/3



Different early failed attempts at machine learning

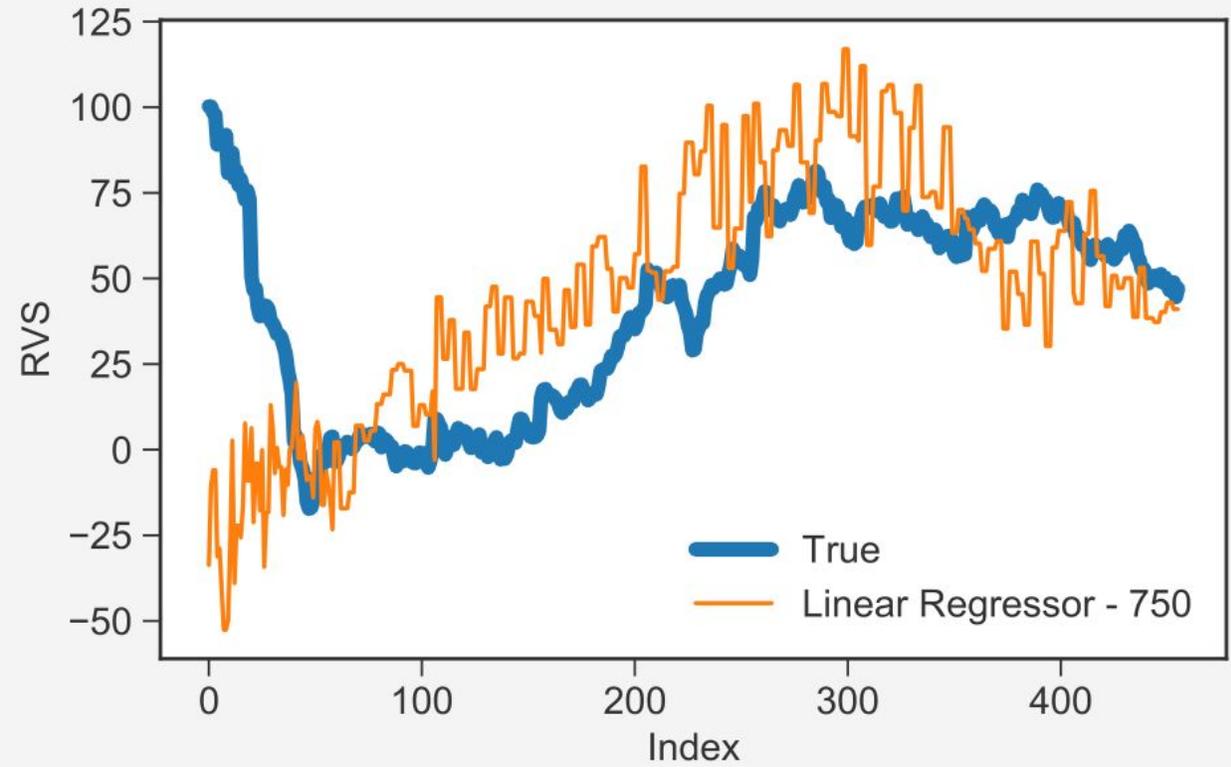
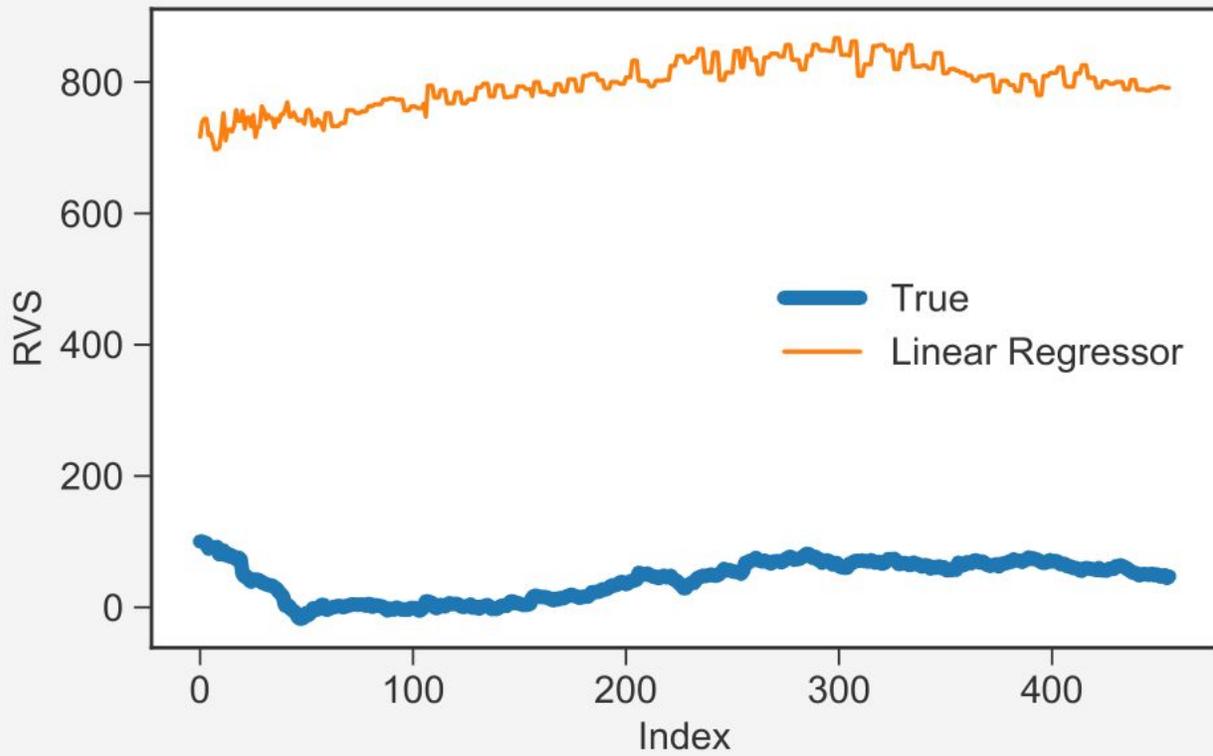


Predictions from randomly sampling train and test data from all sets using XGBoost (non-scaled)



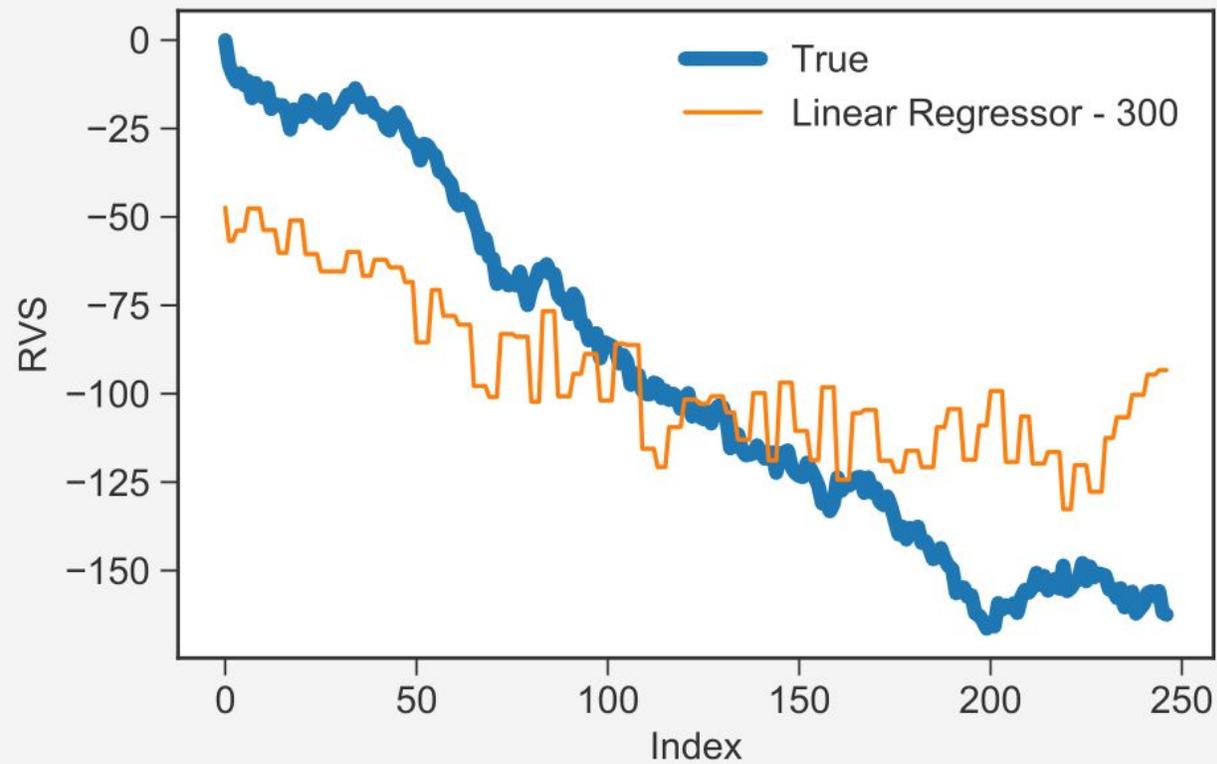
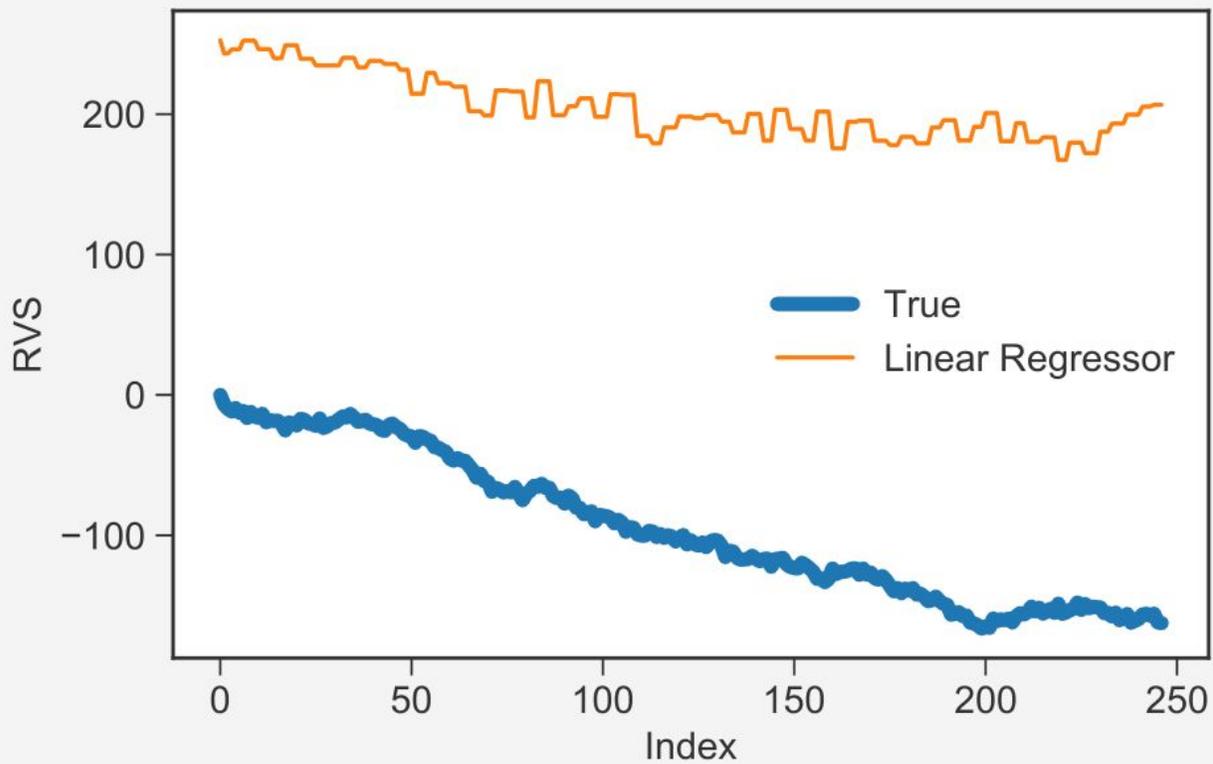
Predictions before scaling

LinReg on test 1/5

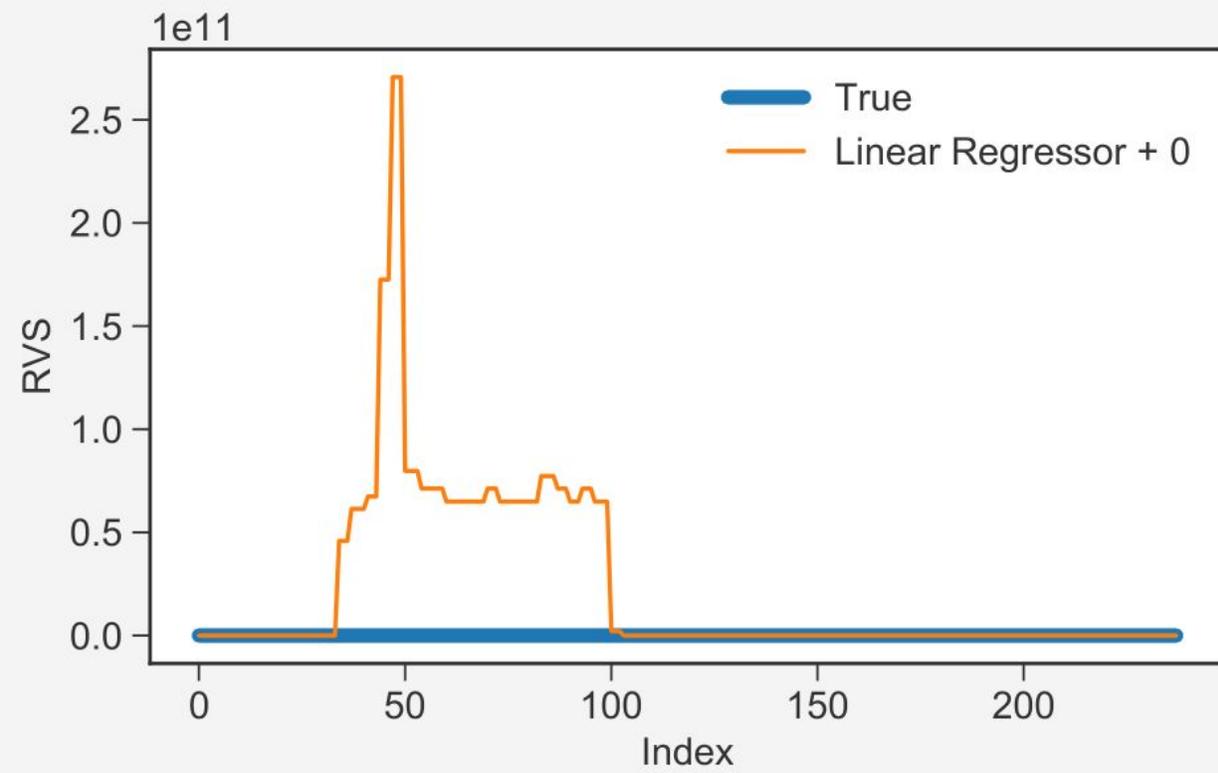
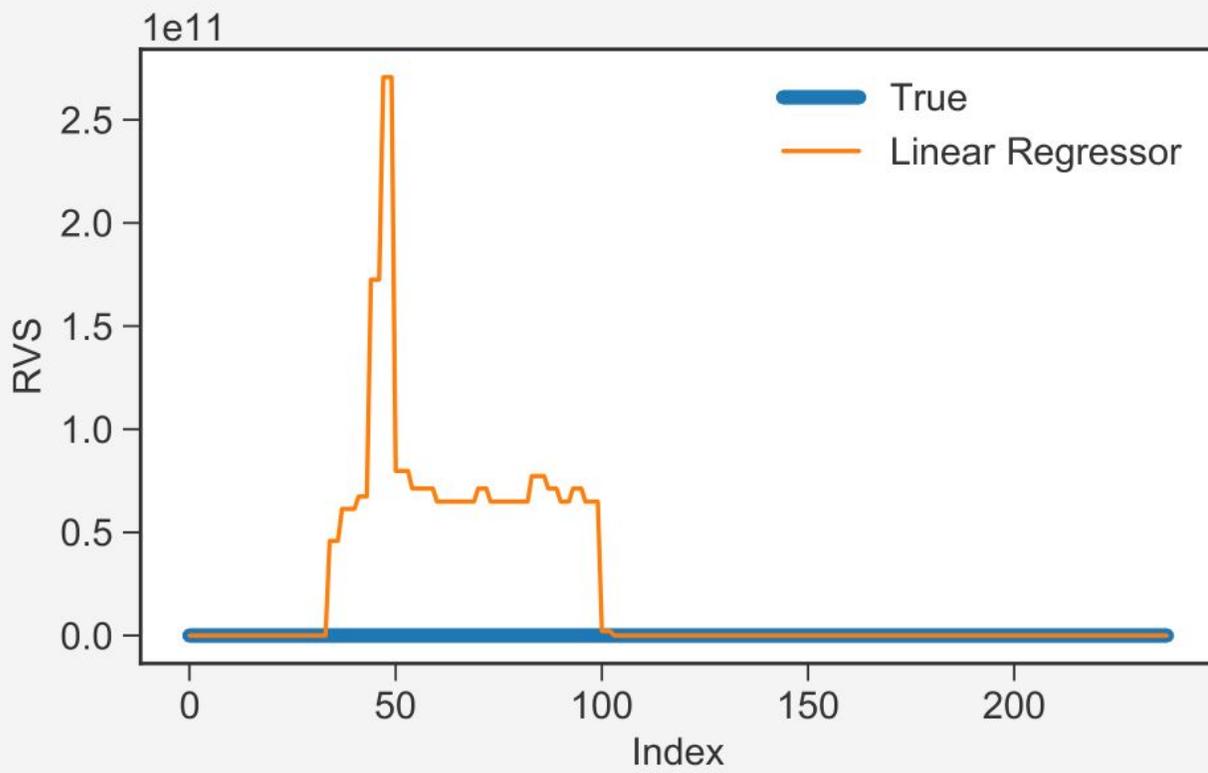


Predictions before scaling

LinReg on test 2/5

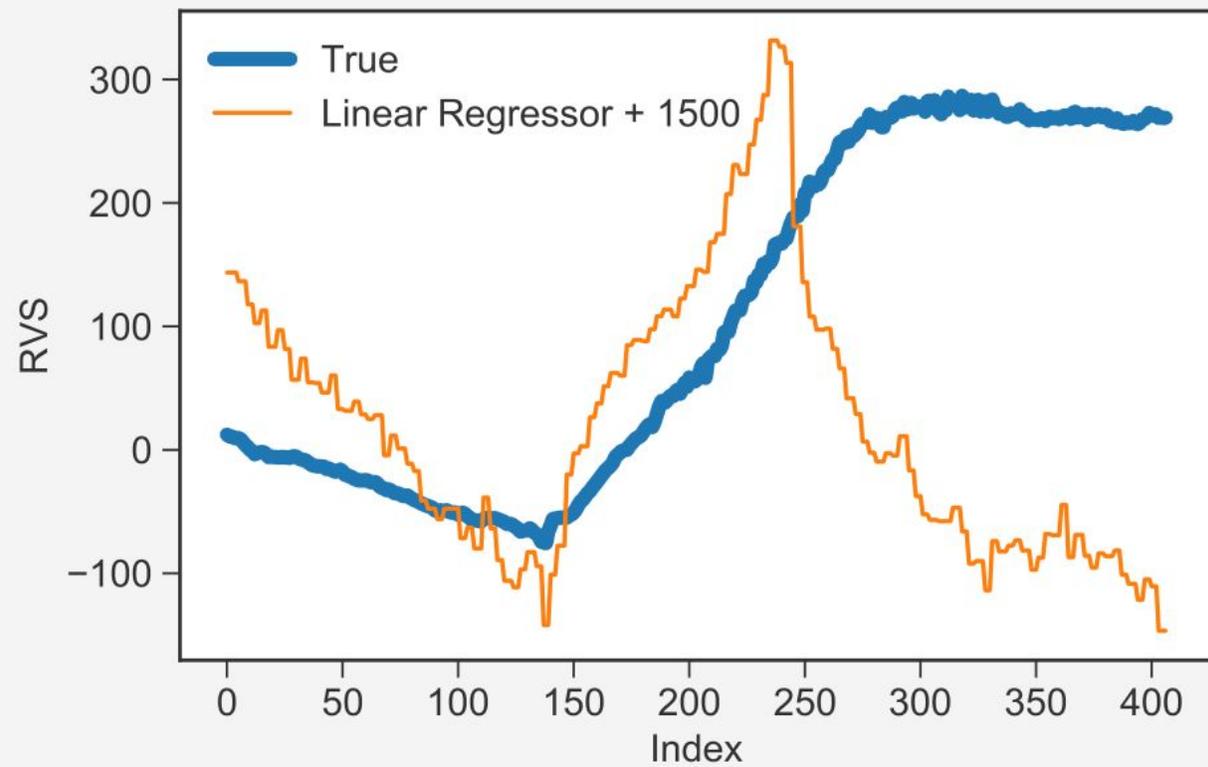
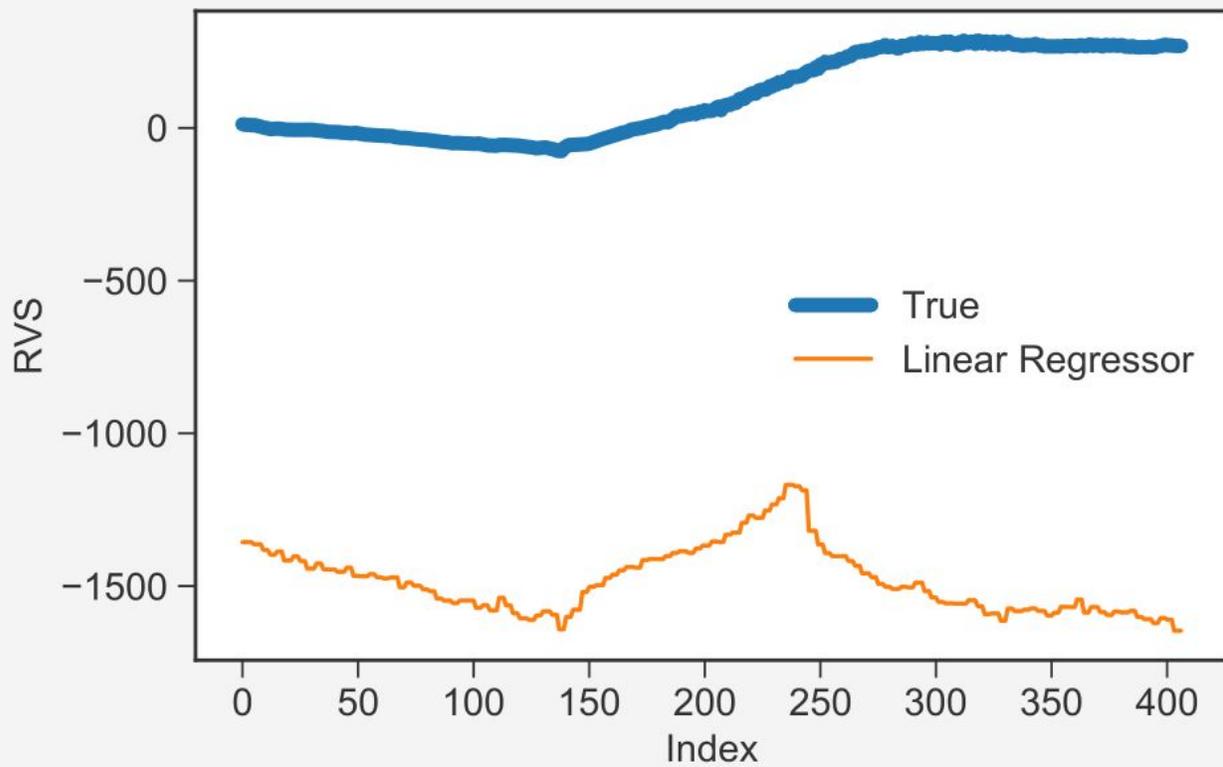


Predictions before scaling LinReg on test 3/5



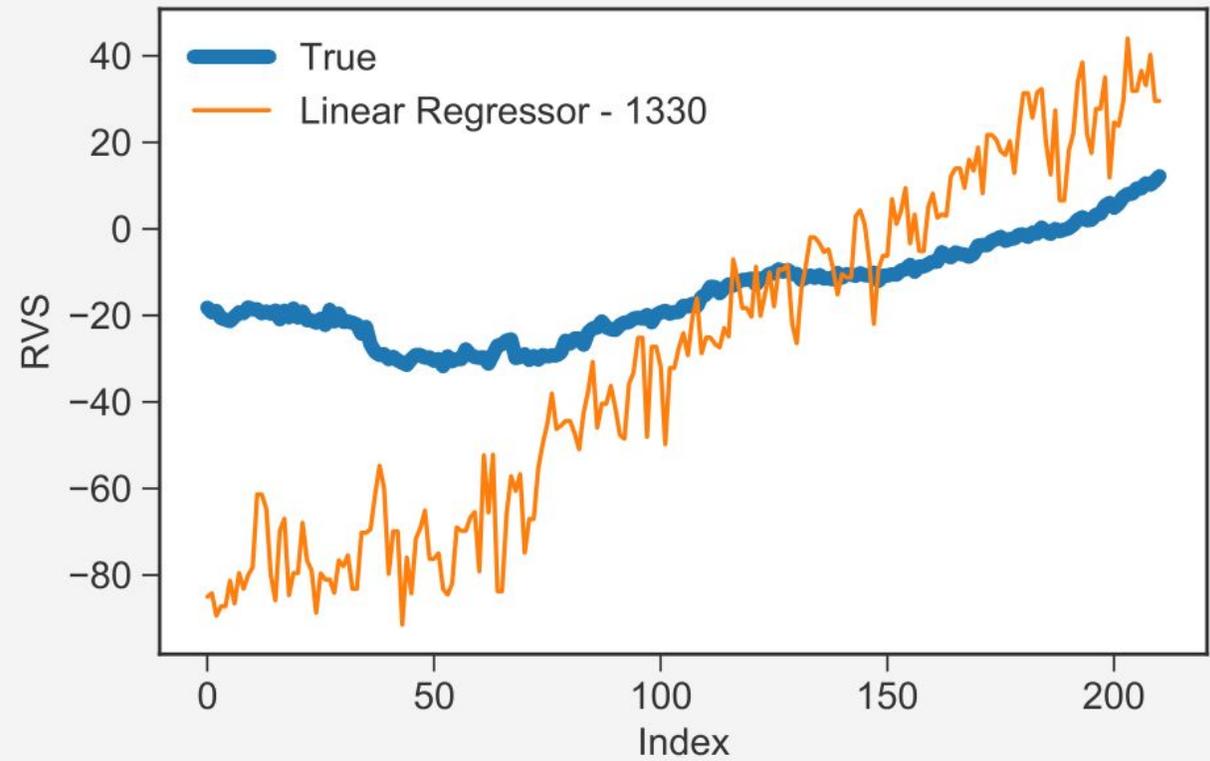
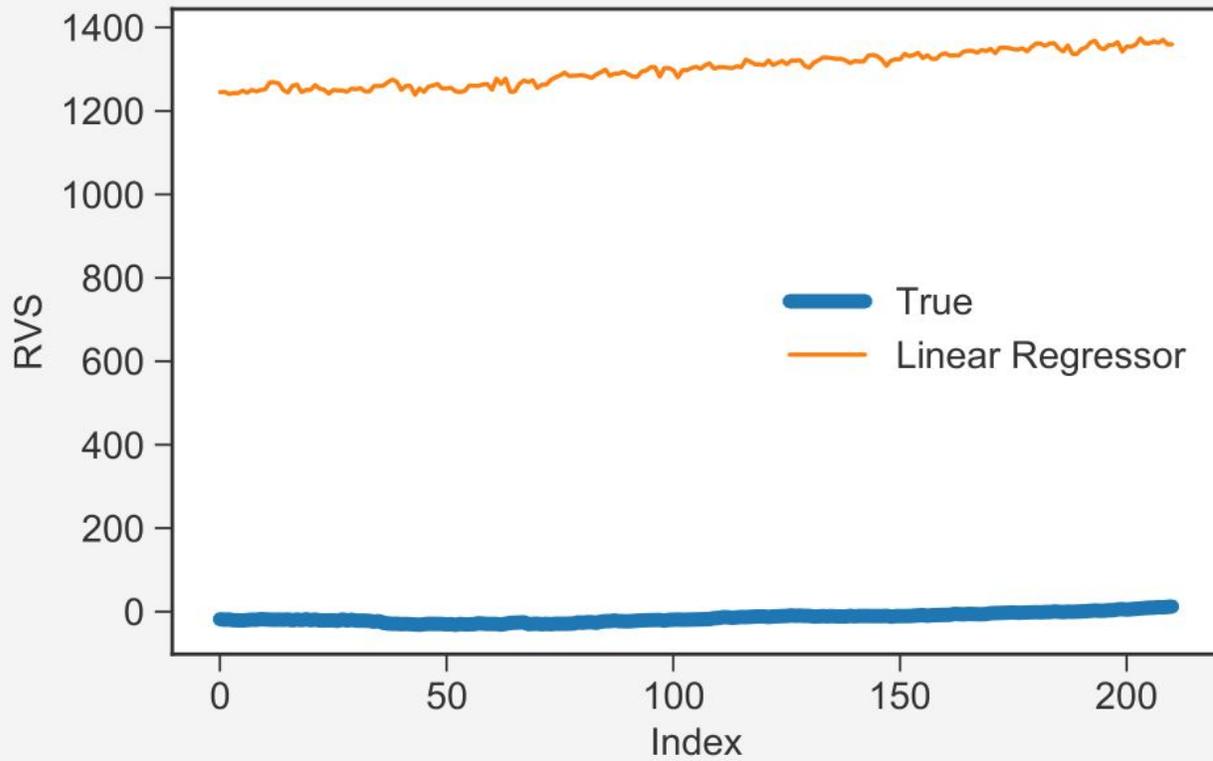
Predictions before scaling

LinReg on test 4/5



Predictions before scaling

LinReg on test 5/5



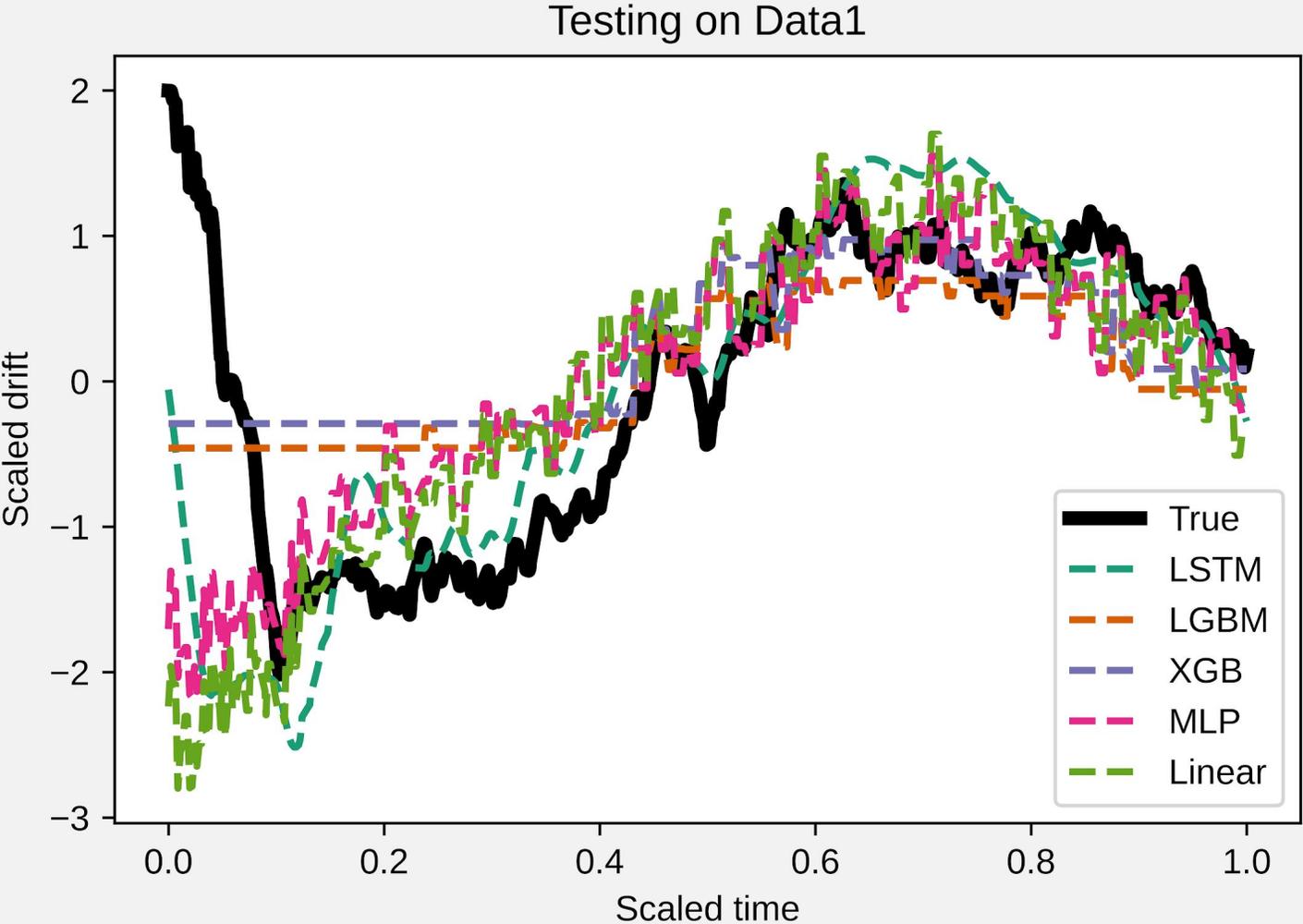
Optimization and results

The following slides present the results that were not presented during the presentation.

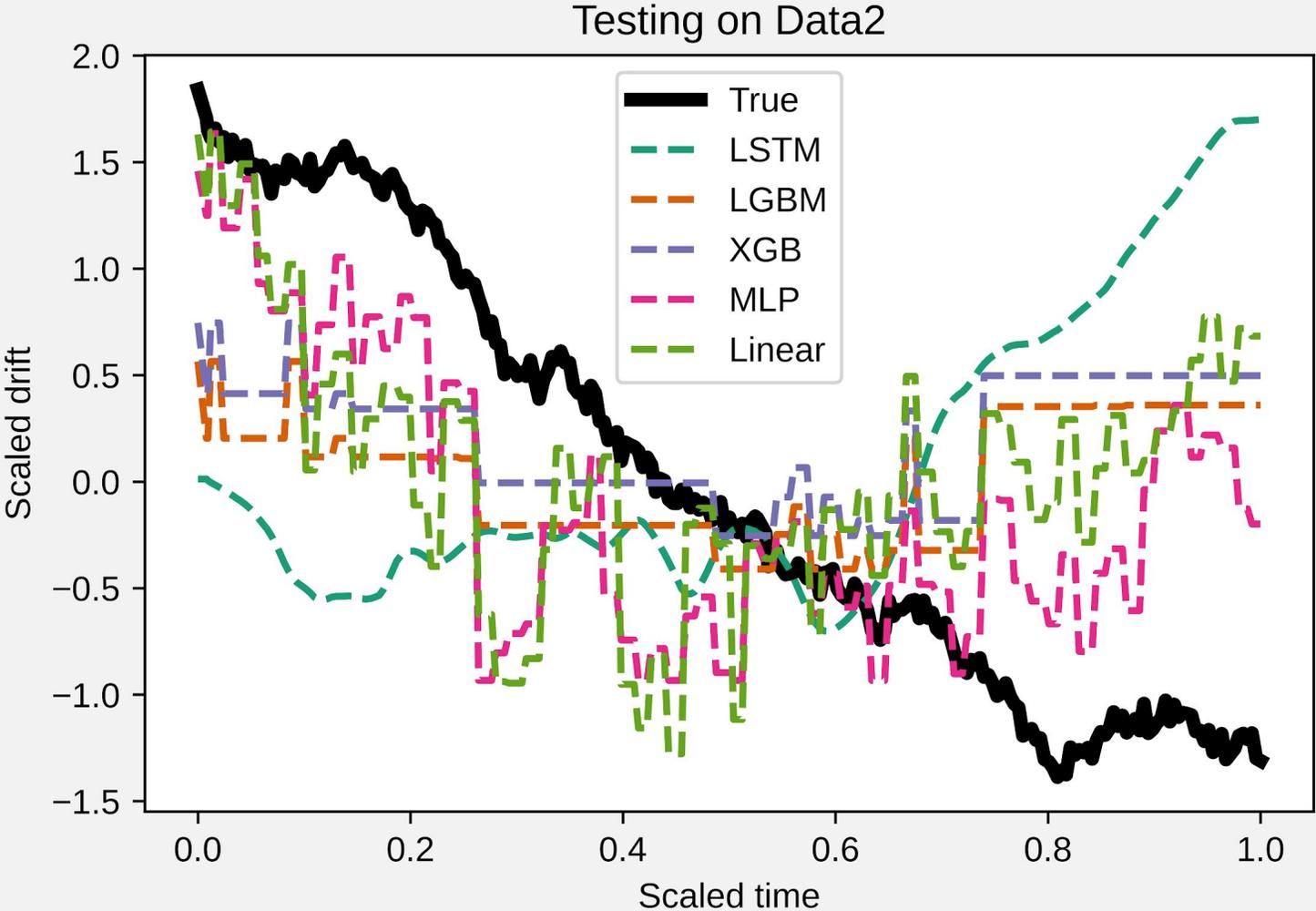
Each algorithm has been trained on all other dataset than the displayed one: Hence for the testing on dataset 1, we have been training on dataset 2, 3, 4, and 5.

To optimize the performance, we have exposed each model to hyperparameter optimization in the form of Bayesian Optimization with the cross validated estimate of the error as target.

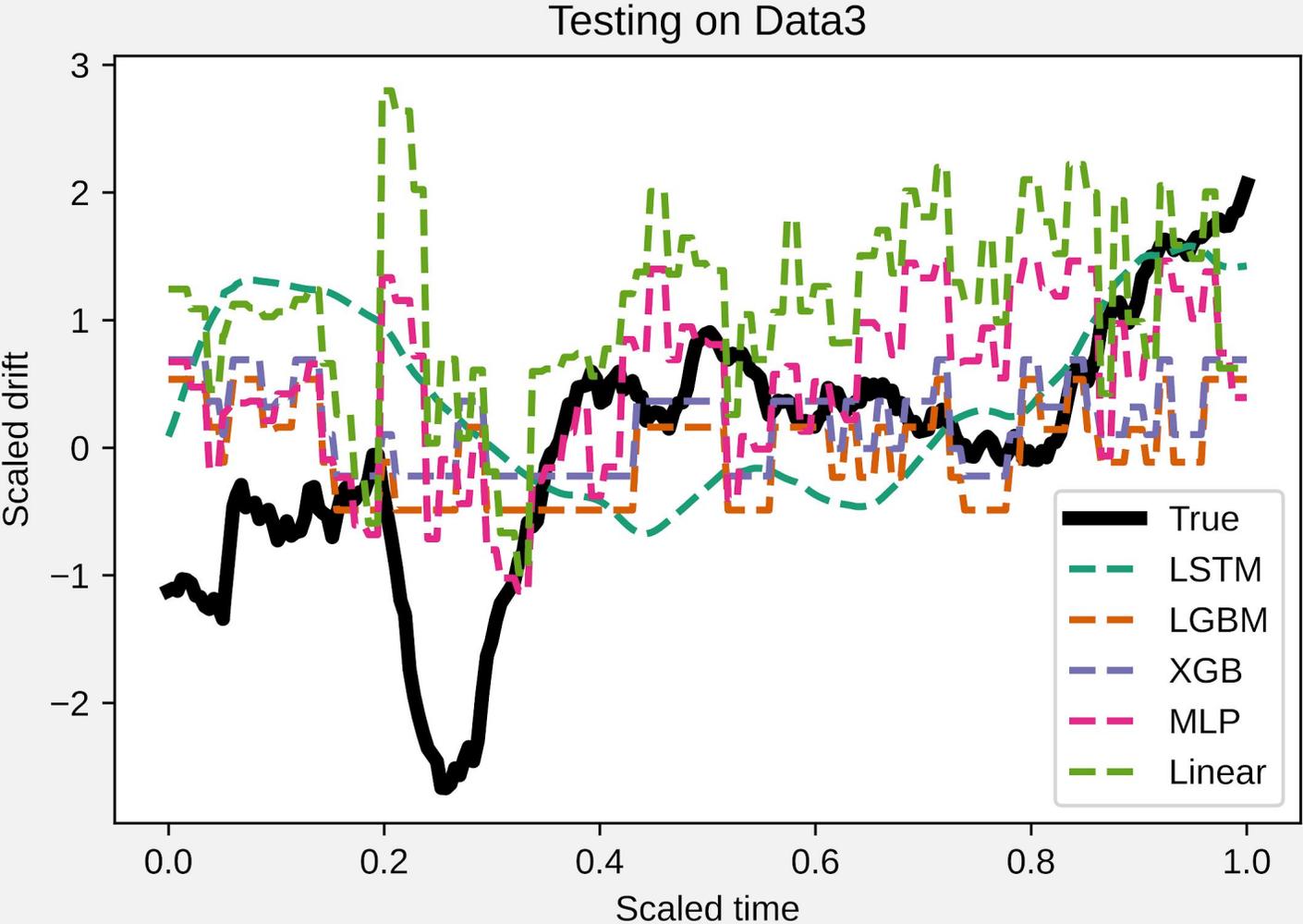
Results from testing on dataset I



Results from testing on dataset 2



Results from testing on dataset 3



Results from testing on dataset 4

