# Joint Data Analysis

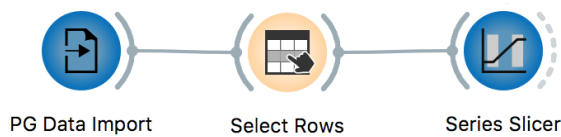Working notes for the hands-on course
at Procter & Gamble

This tutorial will combine data from the previous days — measured
instrument data and the consumer study data — and try to build
predictive models on it.

Handouts prepared by:
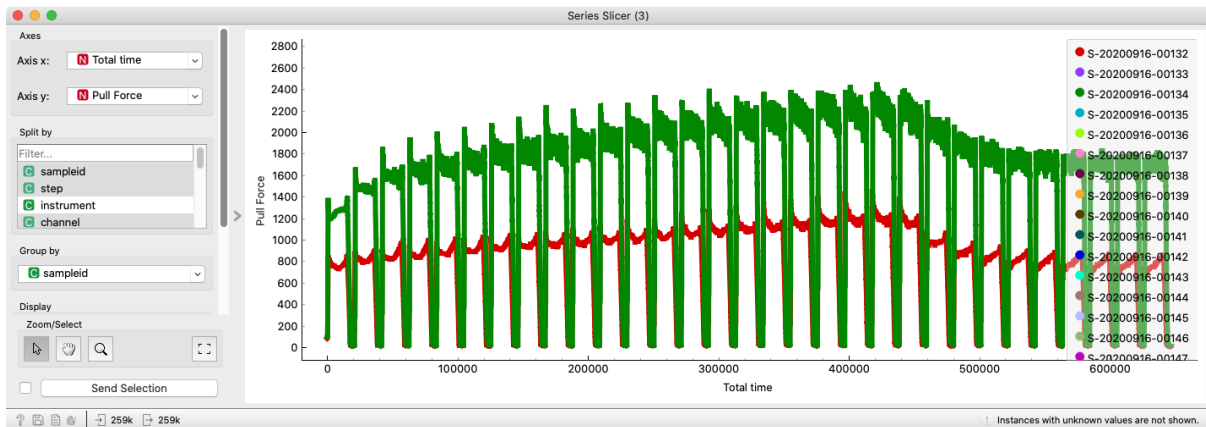Bruna Zupan, Lan Žagar and Primož Godec

# Lesson 1: DRF

This lesson will explore DRF data and its possible relationship with selected questions from the Consumer study data set.

In *PG Data Import*, we load all DRF data, Metadata, and Consumer study data. Note that the DRF experiment has a much bigger data set than the rest of the instruments (6.6 million measurements).
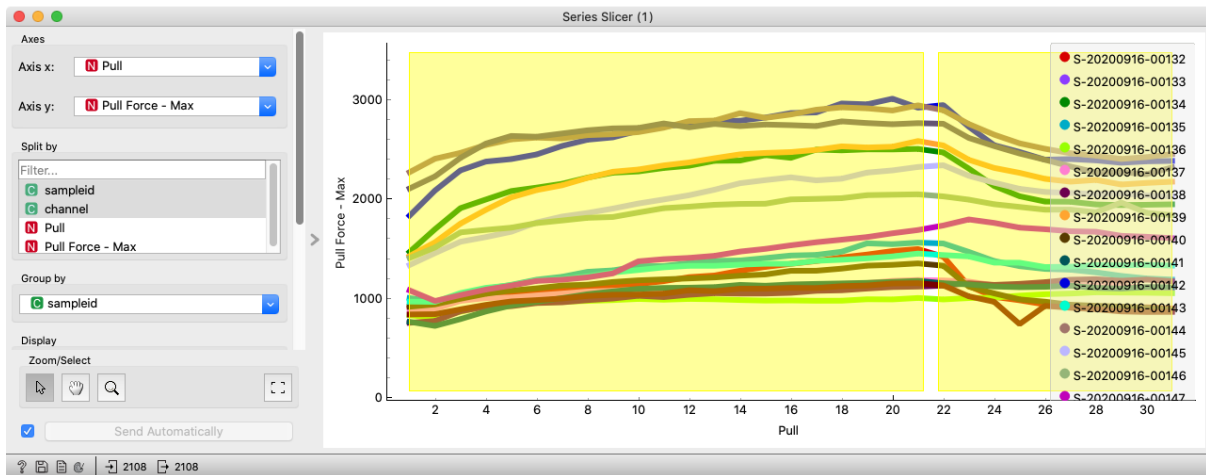
To first visualize a small sample, connect DRF data with the *Select Rows* widget and select two sample ids and step 8. Then, use the *Series Slicer* widget to visualize them.

For Axis x, select `Total Time` to plot all pulls sequentially or `Time` to plot all pulls simultaneously (overlapping). For Axis y, select `Pull Force`. In Split by select `sampleid`, `step`, and `channel`.



Going back to all data, use *Select Rows* to select just step 8 and all samples. Then use the *Aggregate* widget to aggregate each pull with the max `Pull Force` to obtain a summarized series of 31 pull forces.

Connect *Aggregate* widget to the *Series Slicer* widget to visualize aggregated data as a time series. Then construct two slices: shampoo phase (pull 1-21) and rinse phase (pull 22-31).
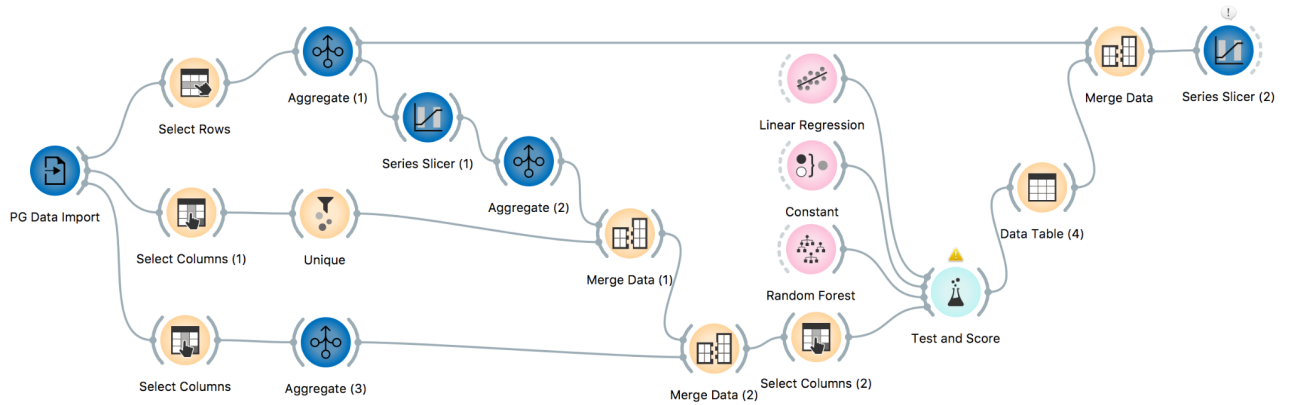
Use another *Aggregate* widget to model the `Pull Force - Max` trends in each slice with a quadratic function and summarize the overall forces with the area under the curve.

Use Metadata from *PG Data Import* to add batch information to the DRF data. First, use *Select Columns* to select `sampleid`, `batchid`, `batchdesc`, and `productdesc` from Metadata then connect it to the *Unique* widget to delete repeated instances. Then use the *Merge* widget to append extra columns to the main DRF (on Data channel) matching rows by `sampleid`.

Connect the *PG Data Import* widget to the *Select Columns* widget with the Consumer Study data. Select the question "Hair feels smooth while rinsing shampoo". Then use *Aggregate* to compute the average rating for each batch. Merge the ratings to DRF data.

Now we have prepared the dataset to explore the relationship between DFR measurements and the consumer's answers.

Note that Orange workflow for exploring the DRF and Consumer study data is available among the training documentation. Below are highlights and hints on how to start the investigation and what to observe.
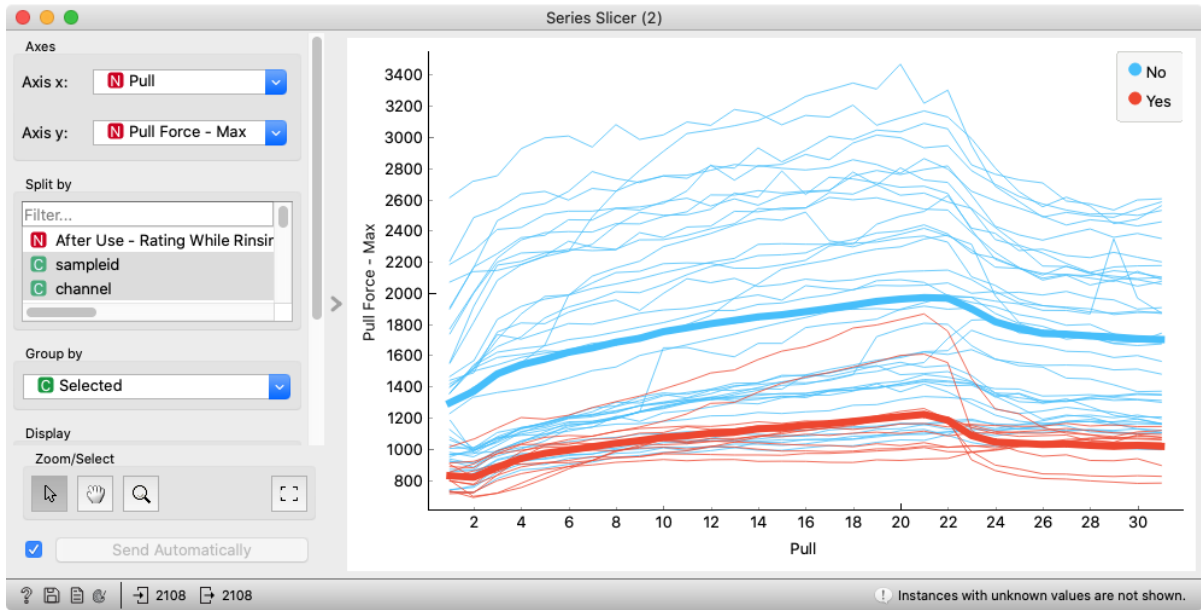
Train some models that can predict the consumer ratings from the profiled DRF data and assess their accuracy.

Instead of the standard cross-validation, try using Cross-validation by feature (in this example use `batchid`). We should use this type of cross-validation to exclude overfitting due to having the same sample tested multiple times.
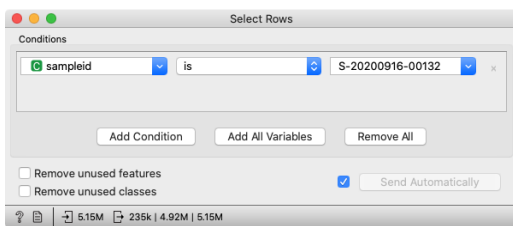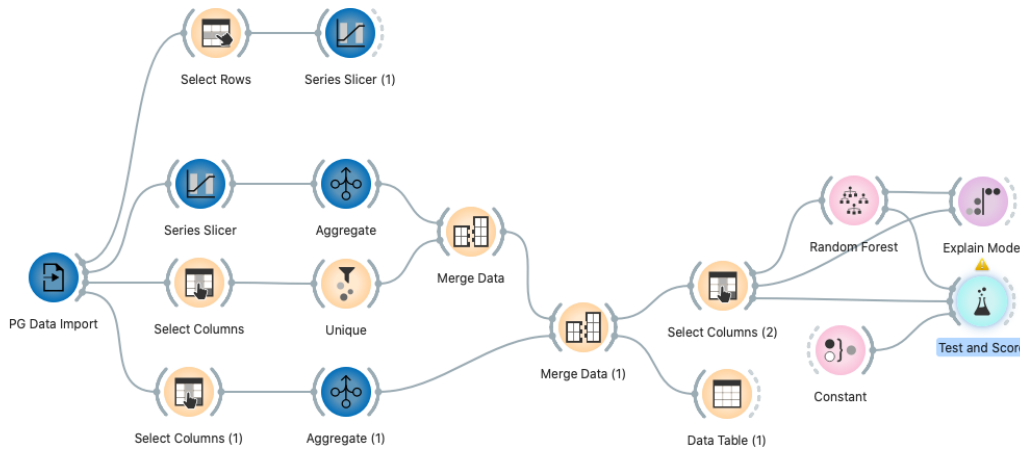
Check the predicted vs. true ratings and plot the profiles for some selected examples.

# Lesson 2: Instron

In this lesson, we will load the Instron data, inspect them with the *Series Slicer*, merge them with the consumer study data, and predict the attributes from the consumer study. With the prediction model, we will observe which attributes from the Instron data contribute the most to predicting two combing-related consumer ratings.





In *PG Data Import* we load all Instron data, Metadata, and Consumer study data. We then select one sample with *Select Rows* and inspect it in the *Series slicer (1)* widget. In the slicer, we can see that there are significant differences between steps 3, 9, and steps 4, 10.

We use another *Series slicer* widget on all data, where we select different slices in the data for the aggregation.



In the *Aggregate* widget that follows we select `sampleid` and `Tress ID` as Rows, step as Columns, and `Load` as Value to aggregate. We use mean and max aggregations. We use *Merge Data* to add `batchid` and `productdesc` columns to the data.

In the bottom branch we use Select Columns to select "After Use - Rating While Hair Damp Outside of Shower - Easy to Comb (Detangle) Hair When Wet" and "After Use - Rating While Hair Dry and Styled - E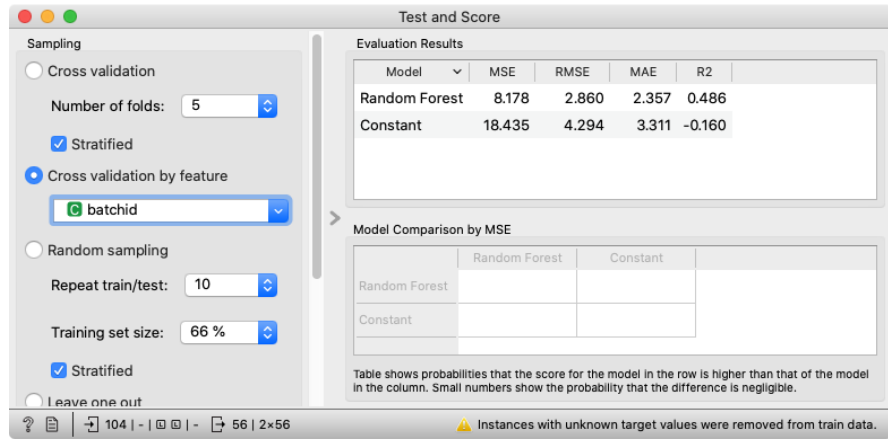asy to Run Fingers/comb Thru My Hair" as attributes and add `BatchID` attribute to meta part of the table. We then average both attributes on the batch level (`BatchID` is set as Rows, both attributes as Values, and Mean as the aggregation). We merge Consumer study aggregations to Intron data with *Merge Data (1)* widget via `BatchID`.

In the last part of the workflow, we use Select Columns to switch between both selected consumer study attributes as a target variable. In *Test and Score* we see that Random Forest with a Mean absolute error (MAE) of 2.357 predict the target variable "After Use - Rating While Hair Damp Outside of Shower - Easy to Comb (Detangle)".

Finally, we will open the *Explain Model* widget which reveals that the mean load from step 3 of the slice Half1 of Instron aggregation is the most important attribute for predicting the target attribute.



We could also try to predict the other selected attributes from the consumer study — "After Use - Rating While Hair Dry and Styled - Easy to Run Fingers/comb Thru My Hair". What is MAE in this case and does the explanation change and how?

# Lesson 3: IFF

We have already performed some explorative analysis of IFF data and seen that the area under the curve and the absolute area under the curve profile the data well. Now we will use that profiling to predict consumer ratings.

Start as before and select a single step with *Select Rows* to simplify the task. After that, use *Aggregate* to profile the data for each `sampleid` and `dataset` combination. Merge in the batch information from Metadata and batch-aggregated responses from the Consumer study. You can try to predict the responses to "Easy to run fingers/comb thru my hair" and "Leaving my hair feeling soft" or find another question you think to be related to IFF measurements.

Because we have just 2 independent features, you might want to try some non-linear models in addition to Linear regression (e.g. *Random Forest*, *kNN*). Use *Test and Score* to evaluate different models and choose Cross validation by feature `batchid`.

*Test and Score* can output predictions. Visualize them in *Scatter Plot* and select the predicted and actual values for axes x and y. You can review the predictions in the *Data Table* as well. As a bonus try to select some rows of interest and merge the Data output (containing the Selected variable) to the raw data and visualize the selection.

| | led – Easy 1 ^ | kNN | batchdesc |
|---|---|---|---|
| 31 | 65.1934 | 61.4453 | BC100C DS25K55D100 1.1% Florabelle (EXP-20... |
| 30 | 65.1934 | 61.0606 | BC100C DS25K55D100 1.1% Florabelle (EXP-20... |
| 29 | 65.1934 | 61.4005 | BC100C DS25K55D100 1.1% Florabelle (EXP-20... |
| 20 | 66.0294 | 64.9293 | BC100C DB25J55D025 1.1% Florabelle (EXP-20... |
| 19 | 66.0294 | 65.034 | BC100C DB25J55D025 1.1% Florabelle (EXP-20... |
| 18 | 66.0294 | 64.9432 | BC100C DB25J55D025 1.1% Florabelle (EXP-20... |
| 17 | 66.0294 | 64.758 | BC100C DB25J55D025 1.1% Florabelle (EXP-20... |
| 24 | 66.568 | 62.5025 | BC100C DB25J25 1.1% Florabelle (EXP-20-ES0... |
| 23 | 66.568 | 64.056 | BC100C DB25J25 1.1% Florabelle (EXP-20-ES0... |
| 22 | 66.568 | 63.0843 | BC100C DB25J25 1.1% Florabelle (EXP-20-ES0... |
| 21 | 66.568 | 65.93 | BC100C DB25J25 1.1% Florabelle (EXP-20-ES0... |
| 4 | 68.1287 | 69.206 | BC189T 1.0% Florabelle BC189T A17NKP301015... |
| 3 | 68.1287 | 69.4383 | BC189T 1.0% Florabelle BC189T A17NKP301015... |
| 2 | 68.1287 | 69.5327 | BC189T 1.0% Florabelle BC189T A17NKP301015... |
| 1 | 68.1287 | 69.5374 | BC189T 1.0% Florabelle BC189T A17NKP301015... |
| 8 | 71.0366 | 67.4915 | BC189T 1% Florabelle BC189H A20N25G15S10 ... |
| 7 | 71.0366 | 67.4318 | BC189T 1% Florabelle BC189H A20N25G15S10 ... |
| 6 | 71.0366 | 67.3136 | BC189T 1% Florabelle BC189H A20N25G15S10 ... |
| 5 | 71.0366 | 65.6891 | BC189T 1% Florabelle BC189H A20N25G15S10 ... |