# 1 Data import and preparation

This part of the document deals with the data preparation of the provided cooper wire data before the data analysis.
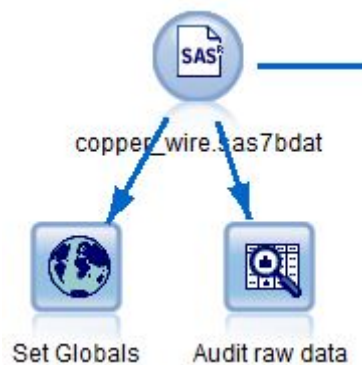
## 1.1 Data import



Abbildung 1: Data import in Stream

The data is imported via the node *SAS file*. The node *Set Globals* is used for setting the audited data results of the raw imported data as global values, which get used later on for the data preparation. The node *Data Audit* is used for analyzing the raw data.
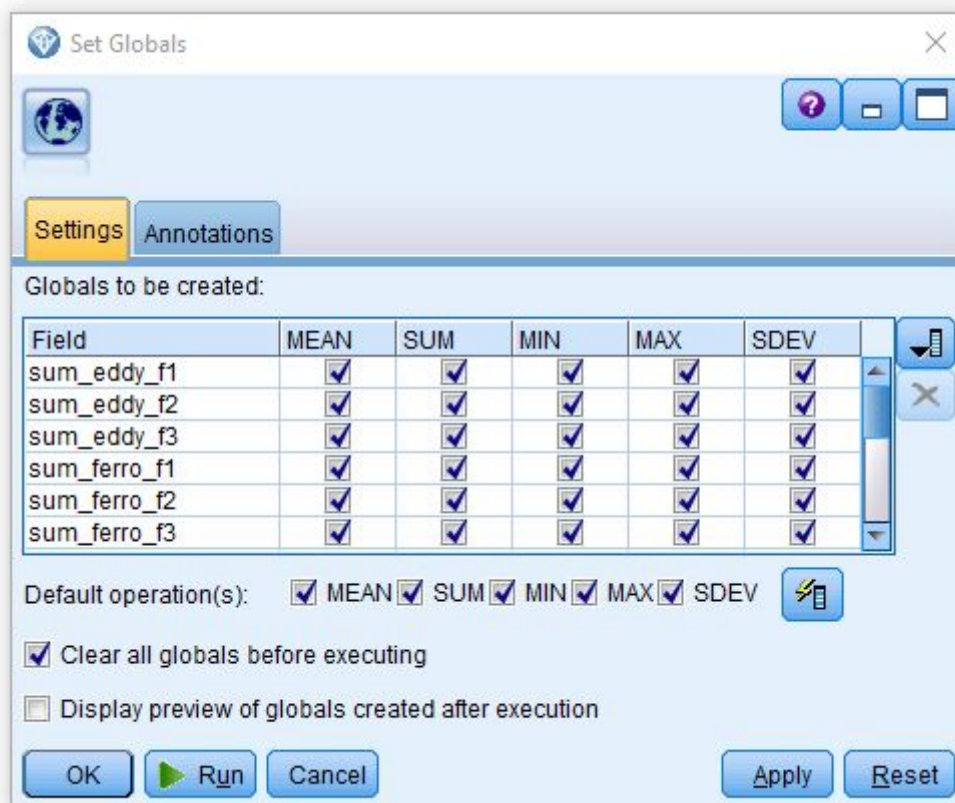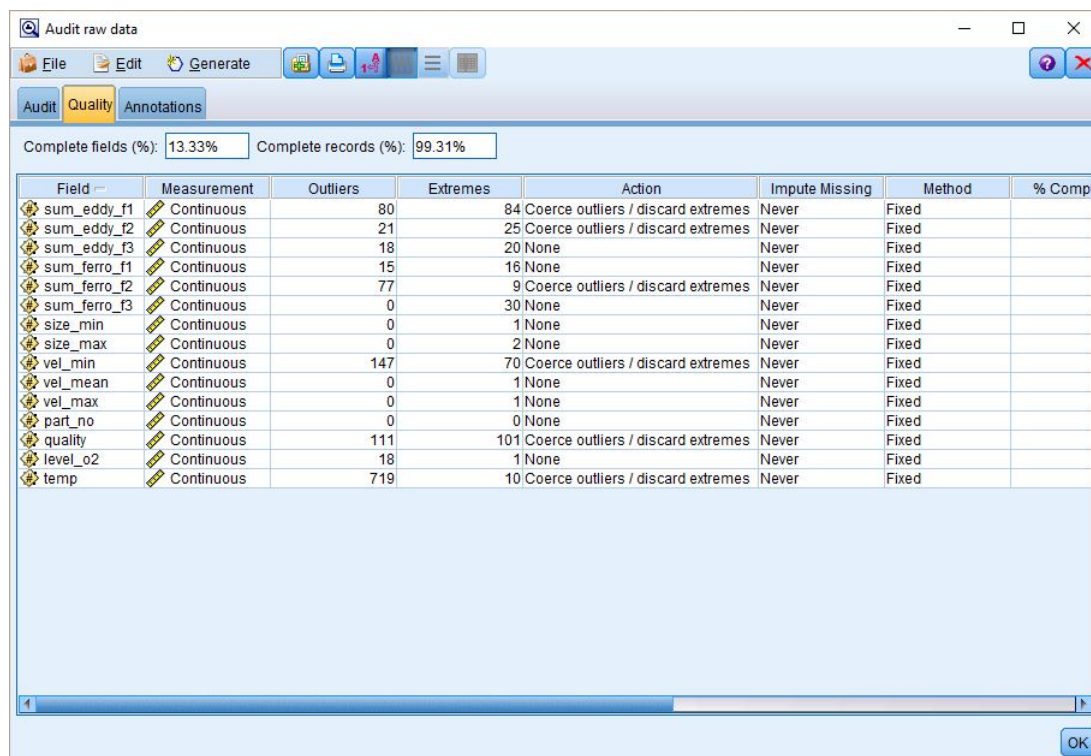


Abbildung 2: Set audit results as global values in the stream
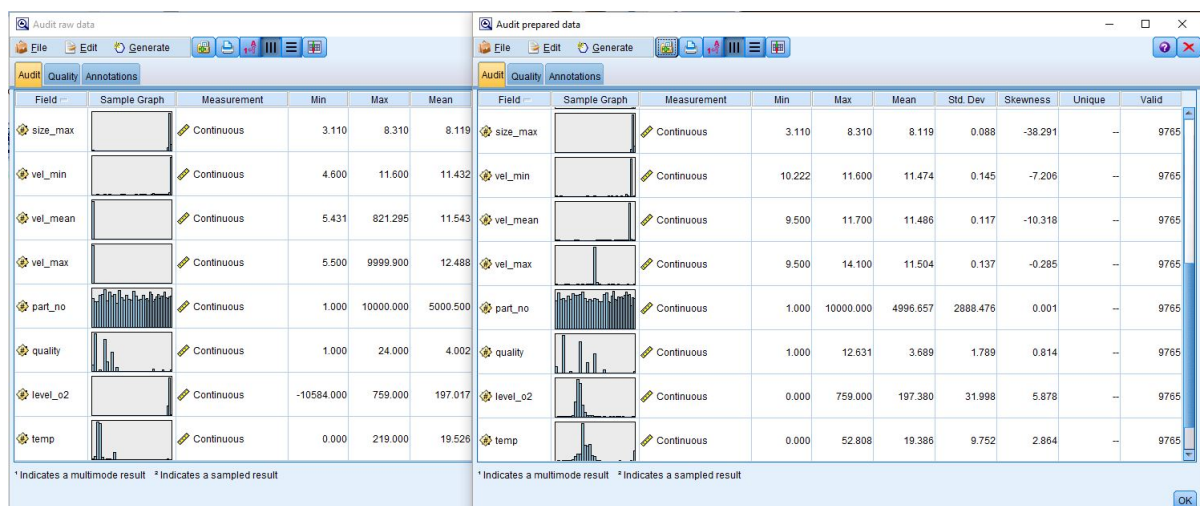
---

## 1.2 Data preparation

The outliers and extremes where determined during the audit of the raw data.



Abbildung 3: Audit of the raw data
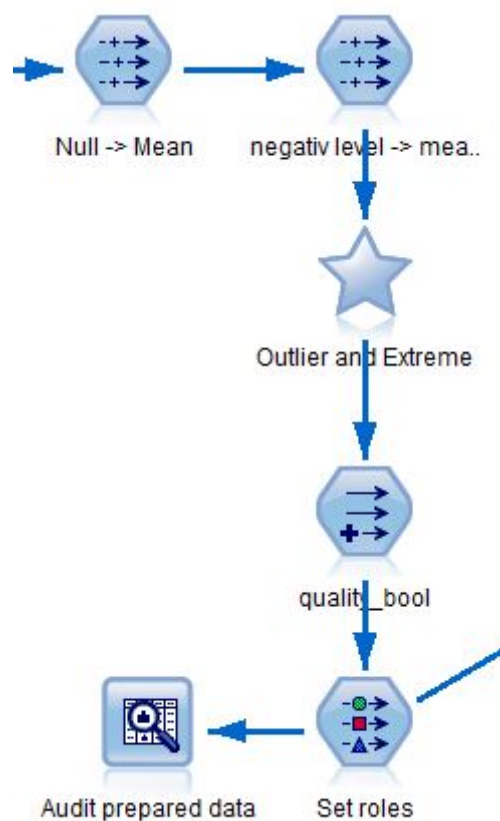


Abbildung 4: Audit of the raw data

Abbildung 5: Flow of data preparation tasks

This flow prepares the data for the later analysis. The following tasks are performed:

- Null values will be replaced with the mean value set by the *Set Global* node

- The negative value of the field *temp* will be replaced with the global mean of this field

- The outliers and extremes will be handled as you can see in image 3

- A new filed will be created *quality_bool* which represents the quality state good or false

- The fields which are not considered to be relevant will be set as ignored and the field *quality_bool* will be set as the target field for the further analysis
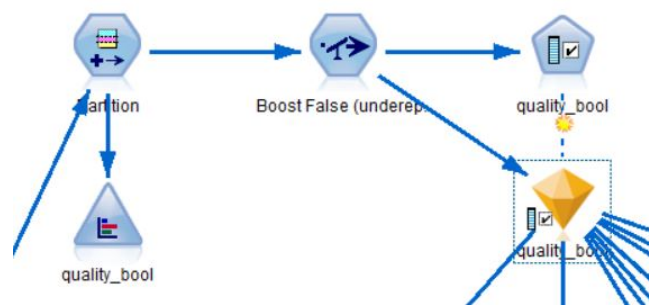
## 1.3 Predictive Model



Abbildung 6: Flow of further data preparation

This part of the stream prepares the data in the following way:

- The node *Partition* splits the data in a training and testing data.

- The node *Distribution* shows us that the field *quality_bool* is very bad distributed. *(False=2.24, True=97.76)*
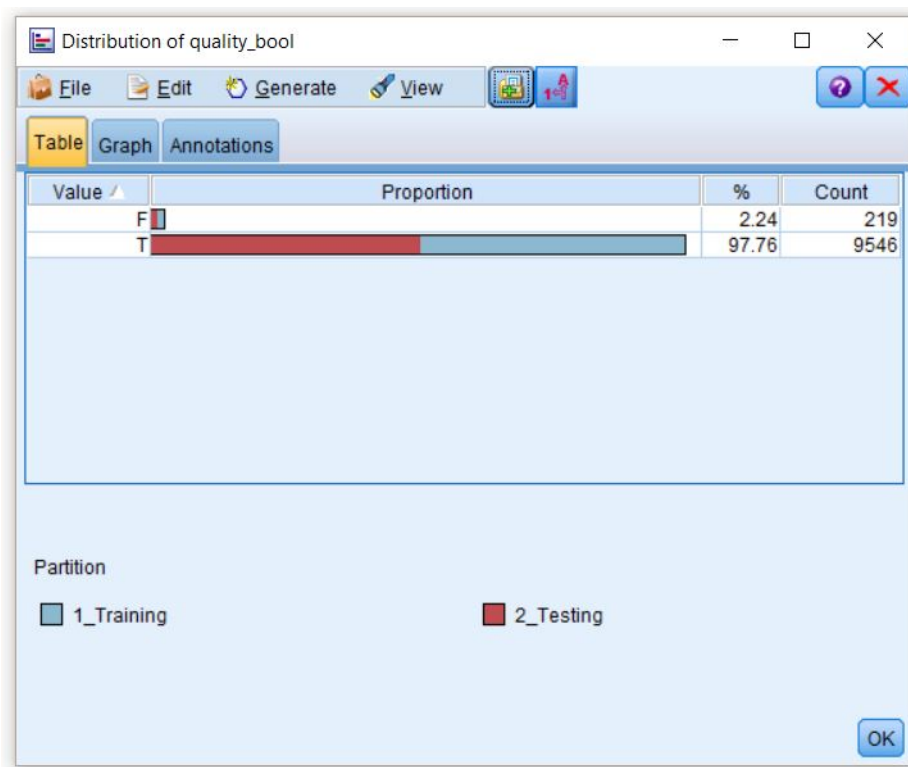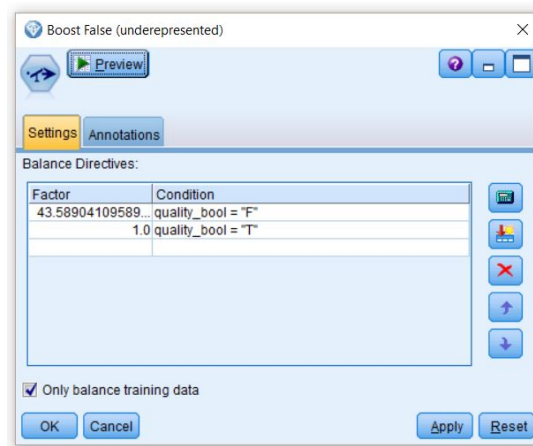
- 



Abbildung 7: Badly distributed *quality_bool*

As we can see that the *False* quality is underrepresented compared to the *True* quality.

Abbildung 8: Boost of quality *False*

The node *Balance* has been generated by the node *Distribution* and boost the representation of the *False* quality. After this nodes follows the node *Field selection* which removes fields which are not related to the *target*.

The node *Field selection* has reduced the count of fields from *15* down to *10*, therefore has removed *5* fields.

### 1.3.1 Not partitioned data to feature selection and boost



Abbildung 9: C5.0 with no partitioned and partitioned data

Abbildung 10: Neuronal Net with no partitioned and partitioned data



Abbildung 11: Logistic Regression with no partitioned and partitioned data

## 1.3.2 Feature selection to PCA



Abbildung 12: C5.0 feature selection to PCA



Abbildung 13: Neuronal Net feature selection to PCA

Abbildung 14: Logistic Regression feature selection to PCA

### 1.3.3 Fine tuning of PCA part