# 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.
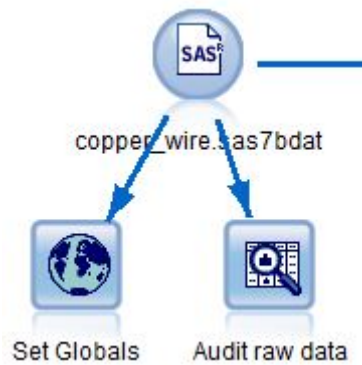
## 1.1 Data import



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.
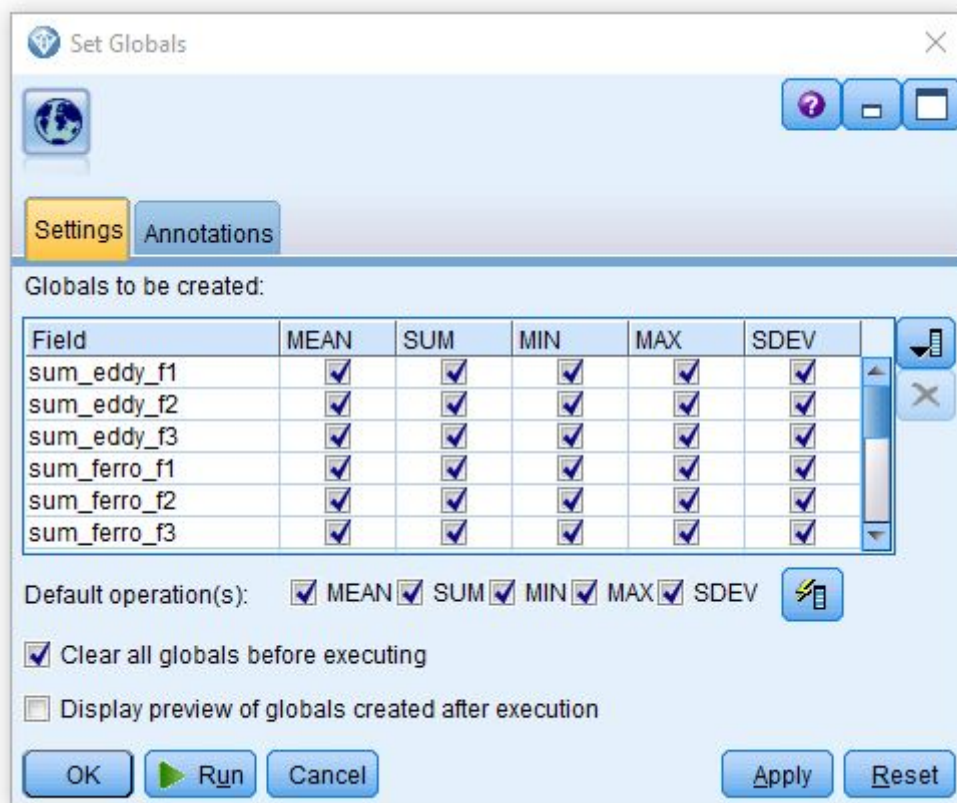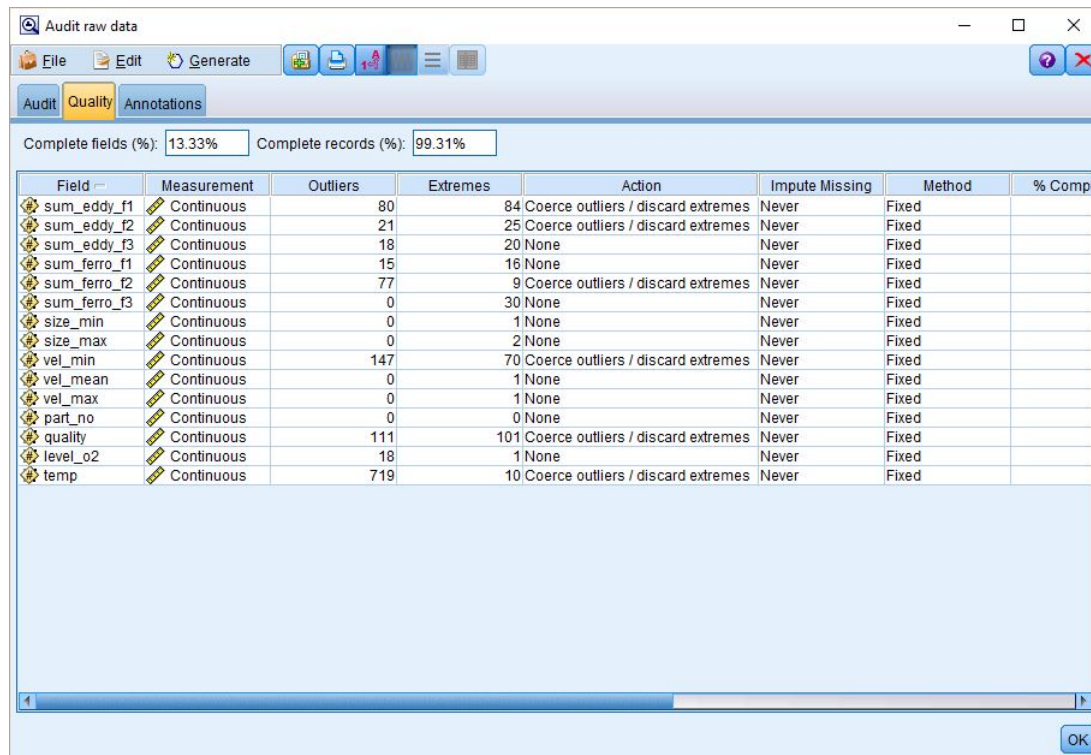


Abbildung 2: Set audit results as global values in the stream
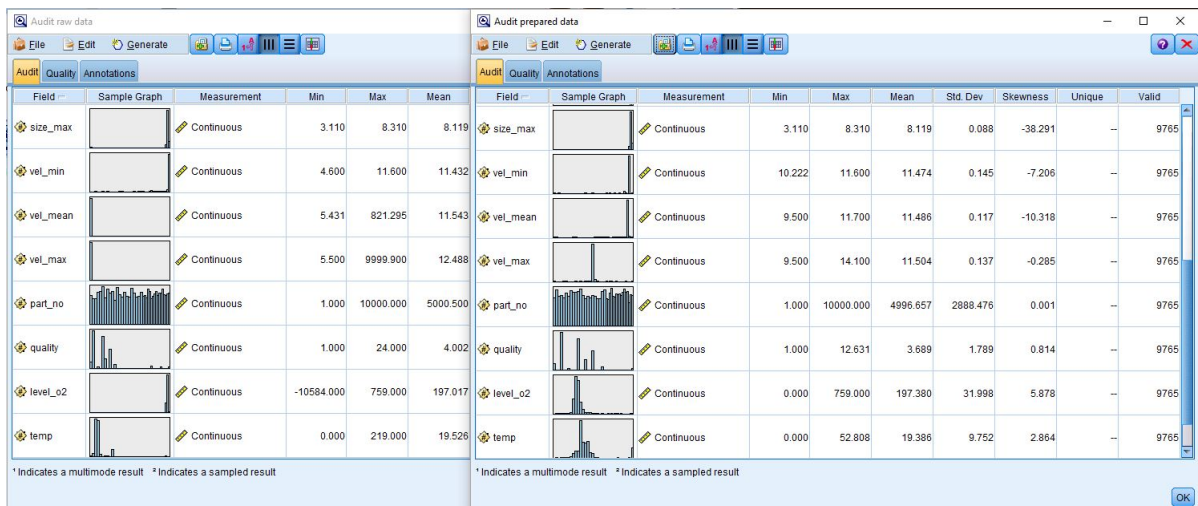
## 1.2 Data preparation

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 by splitting it into test and training data, additionally the *False qualtity_bool* will be boost to increase their representation in the data. At last the data gets prepared on the one hand with a *Feature Selection* and on the other hand with a *PAC Factor*.



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*

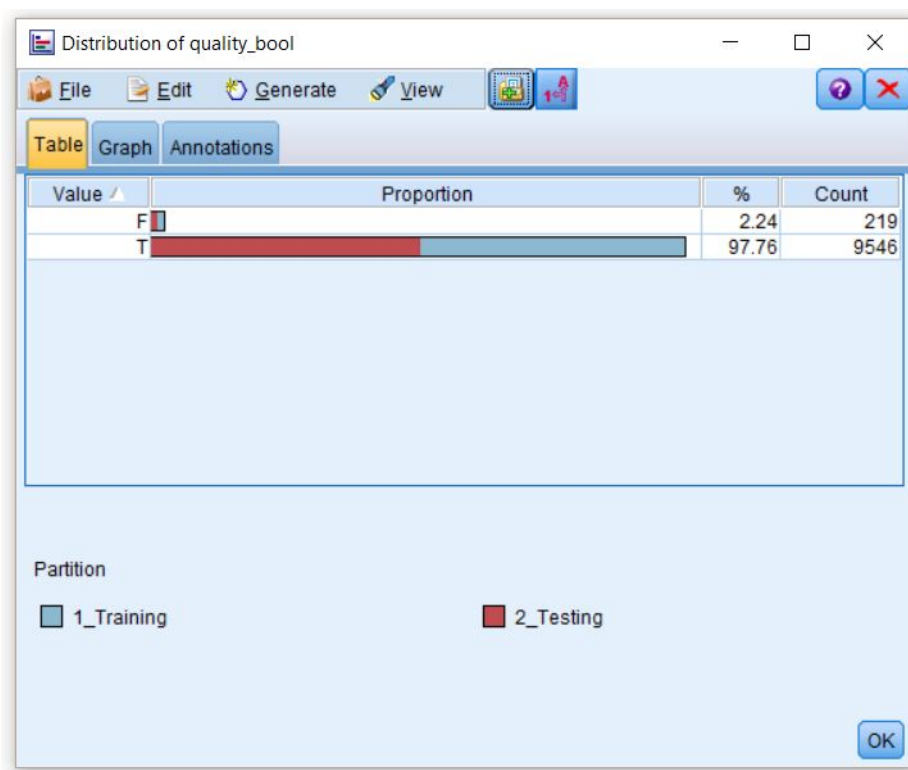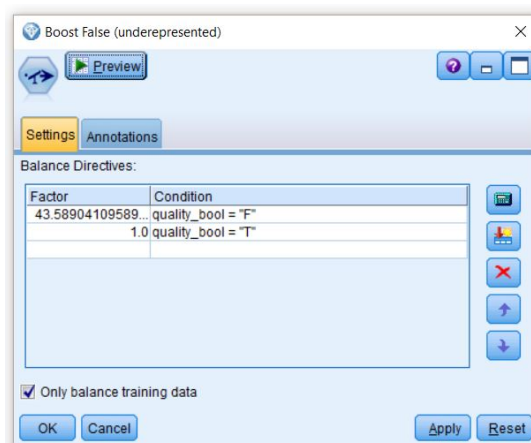The node *Balance* has been generated by the node *Distribution* and boost the representation of the *False* quality. After this nodes follows the node *Field selection* which removes fields which are not related to the *target*.

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

### 1.3.1 What are results of the use of *Feature Selection* and *PCA Factor* with defaults ?



Abbildung 9: C5.0 with *Feature Selection* compared to *PAC*

**NNet feature sel**

Results for output field quality_bool
Comparing $N-quality_bool with quality_bool

| 'Partition' | 1_Training | | 2_Testing | |
|---|---|---|---|---|
| Correct | 7,788 | 77.34% | 4,233 | 85.93% |
| Wrong | 2,282 | 22.66% | 693 | 14.07% |
| Total | 10,070 | | 4,926 | |

Coincidence Matrix for $N-quality_bool (rows show actuals)

| 'Partition' = 1_Training | F | T |
|---|---|---|
| F | 3,655 | 1,699 |
| T | 583 | 4,133 |
| 'Partition' = 2_Testing | F | T |
| F | 52 | 44 |
| T | 649 | 4,181 |

User Defined Score for $N-quality_bool

| 'Partition' | 1_Training | 2_Testing |
|---|---|---|
| Mean | -106.634 | 83.059 |
| Sum | -1073800.0 | 409150.0 |
| Minimum | -1000.0 | -1000.0 |
| Maximum | 100.0 | 100.0 |
| Standard Deviation | 403.139 | 104.302 |

**NNet pca**

Results for output field quality_bool
Comparing $N-quality_bool with quality_bool

| 'Partition' | 1_Training | | 2_Testing | |
|---|---|---|---|---|
| Correct | 8,264 | 81.98% | 4,319 | 87.68% |
| Wrong | 1,817 | 18.02% | 607 | 12.32% |
| Total | 10,081 | | 4,926 | |

Coincidence Matrix for $N-quality_bool (rows show actuals)

| 'Partition' = 1_Training | F | T |
|---|---|---|
| F | 4,012 | 1,353 |
| T | 464 | 4,252 |
| 'Partition' = 2_Testing | F | T |
| F | 66 | 30 |
| T | 577 | 4,253 |

User Defined Score for $N-quality_bool

| 'Partition' | 1_Training | 2_Testing |
|---|---|---|
| Mean | -69.834 | 86.774 |
| Sum | -704000.0 | 427450.0 |
| Minimum | -1000.0 | -1000.0 |
| Maximum | 100.0 | 100.0 |
| Standard Deviation | 366.984 | 86.729 |

Abbildung 10: Neuronal Net with *Feature Selection* compared to *PAC*

**LogReg feature sel**

Results for output field quality_bool
Comparing $L-quality_bool with quality_bool

| 'Partition' | 1_Training | | 2_Testing | |
|---|---|---|---|---|
| Correct | 7,312 | 72.61% | 3,769 | 76.51% |
| Wrong | 2,758 | 27.39% | 1,157 | 23.49% |
| Total | 10,070 | | 4,926 | |

Coincidence Matrix for $L-quality_bool (rows show actuals)

| 'Partition' = 1_Training | F | T |
|---|---|---|
| F | 3,611 | 1,743 |
| T | 1,015 | 3,701 |
| 'Partition' = 2_Testing | F | T |
| F | 60 | 36 |
| T | 1,121 | 3,709 |

User Defined Score for $L-quality_bool

| 'Partition' | 1_Training | 2_Testing |
|---|---|---|
| Mean | -113.366 | 79.974 |
| Sum | -1141600.0 | 393950.0 |
| Minimum | -1000.0 | -1000.0 |
| Maximum | 100.0 | 100.0 |
| Standard Deviation | 406.296 | 95.095 |

**LogReg pca**

Results for output field quality_bool
Comparing $L-quality_bool with quality_bool

| 'Partition' | 1_Training | | 2_Testing | |
|---|---|---|---|---|
| Correct | 7,582 | 75.21% | 4,199 | 85.24% |
| Wrong | 2,499 | 24.79% | 727 | 14.76% |
| Total | 10,081 | | 4,926 | |

Coincidence Matrix for $L-quality_bool (rows show actuals)

| 'Partition' = 1_Training | F | T |
|---|---|---|
| F | 3,531 | 1,834 |
| T | 665 | 4,051 |
| 'Partition' = 2_Testing | F | T |
| F | 64 | 32 |
| T | 695 | 4,135 |

User Defined Score for $L-quality_bool

| 'Partition' | 1_Training | 2_Testing |
|---|---|---|
| Mean | -120.93 | 85.15 |
| Sum | -1219100.0 | 419450.0 |
| Minimum | -1000.0 | -1000.0 |
| Maximum | 100.0 | 100.0 |
| Standard Deviation | 415.185 | 89.592 |

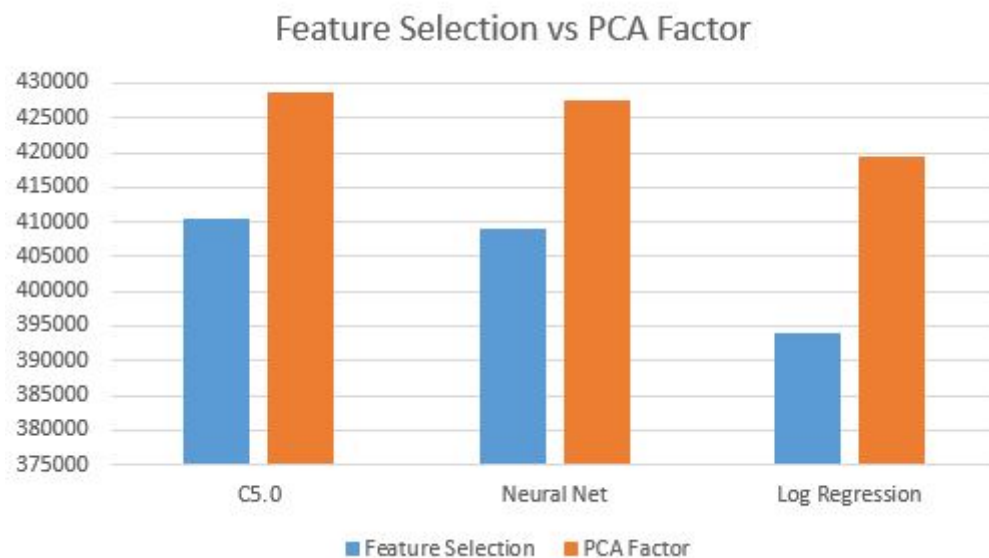Abbildung 11: Logistic Regression with *Feature Selection* compared to *PAC*

Abbildung 12: The chart shows the results of the comparison between *Feature Selection* and *PAC*

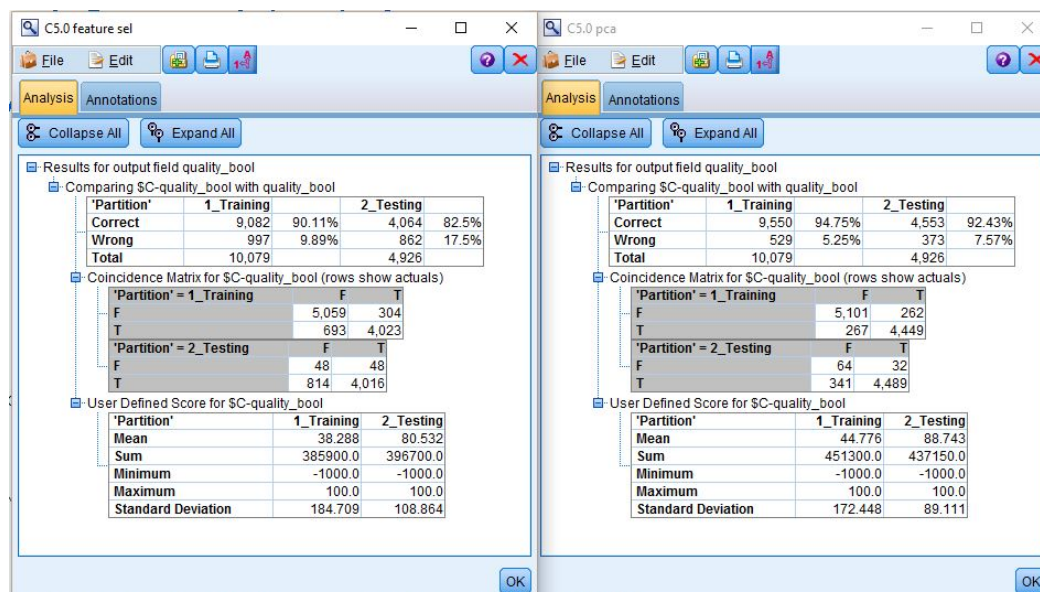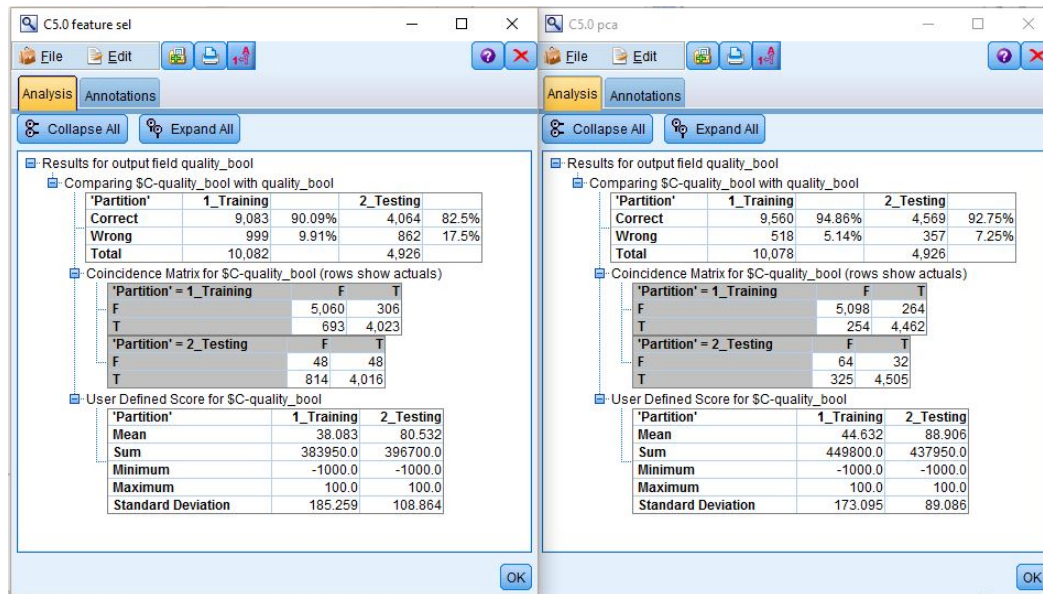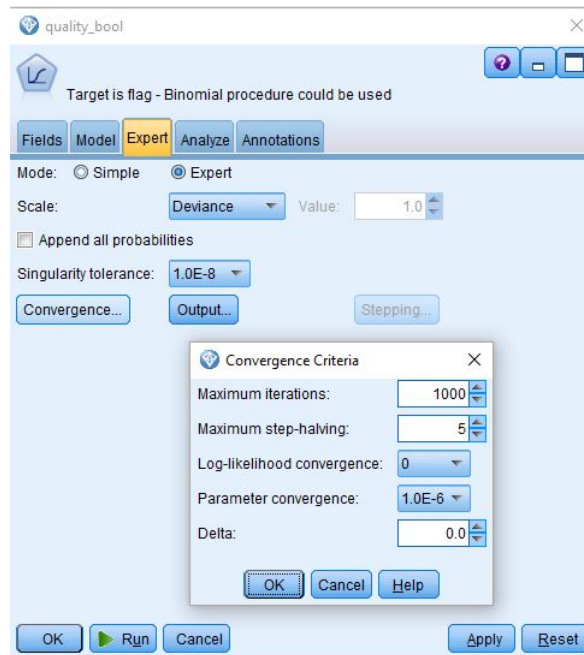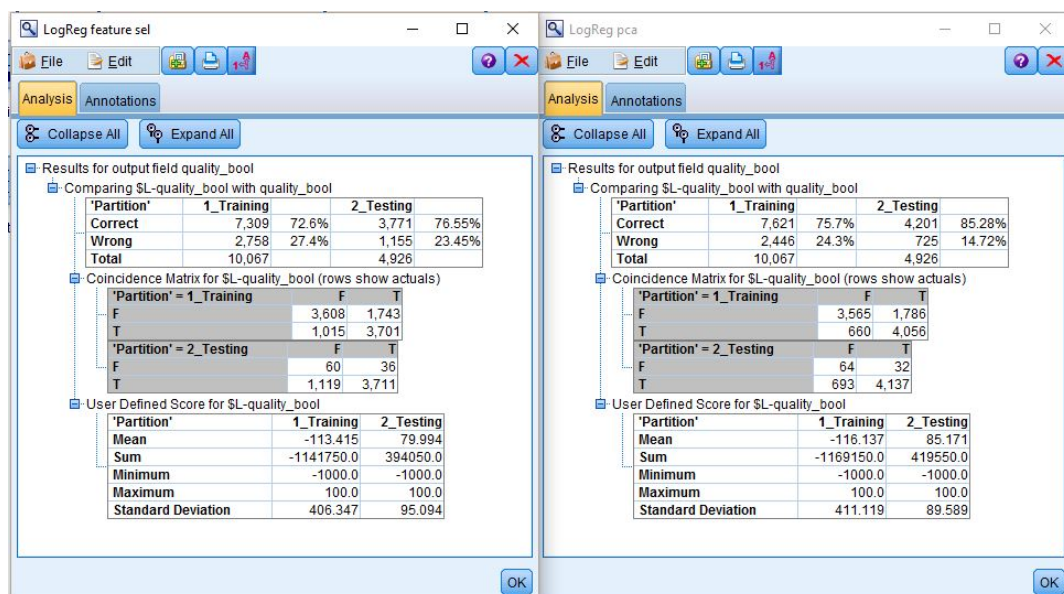### 1.3.2 What are results when C5.0 is modified ?



Abbildung 13: C5.0 (costs for TF=10) with *Feature Selection* compared to *PAC*

Abbildung 14: C5.0 (costs for TF=10, prun serv=10) with *Feature Selection* compared to *PAC*

### 1.3.3 What are results when the Neural Net is modified ?



Abbildung 15: Part one of Neural Net settings

Abbildung 16: Part two of Neural Net settings



Abbildung 17: Neural Net with *Feature Selection* compared to *PAC*

Abbildung 18: Logistic Regression Settings



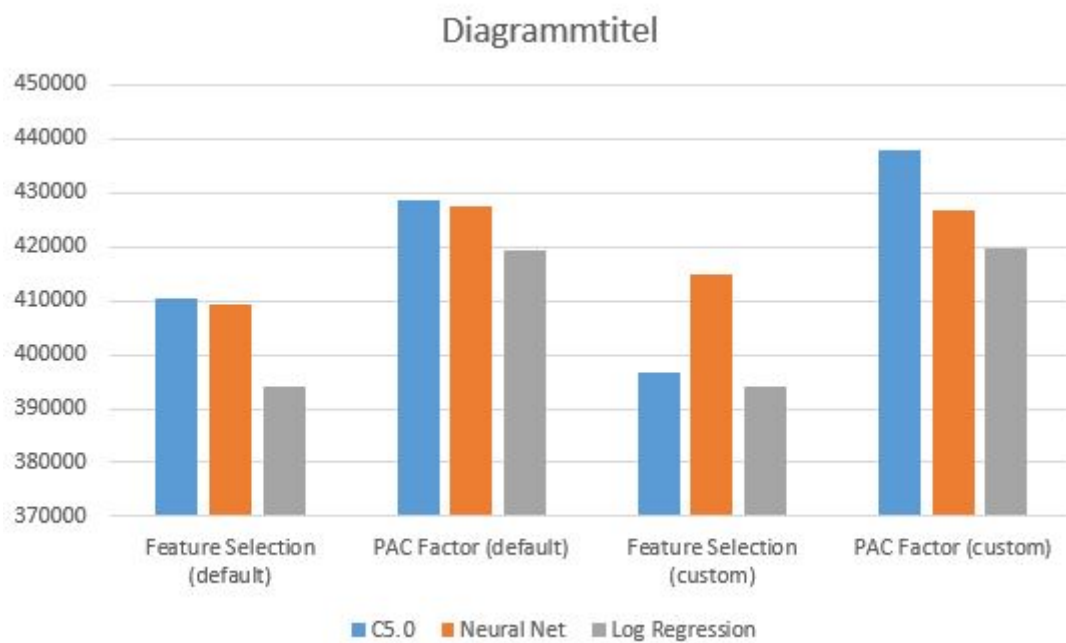Abbildung 19: Logistic regression with *Feature Selection* compared to *PAC*

Abbildung 20: All results of the experiments in a table