



 **BASF**

We create chemistry

Discovering hidden relationships in production data

Elie Maricau - BASF Antwerp

Context: the BASF Antwerpen production site



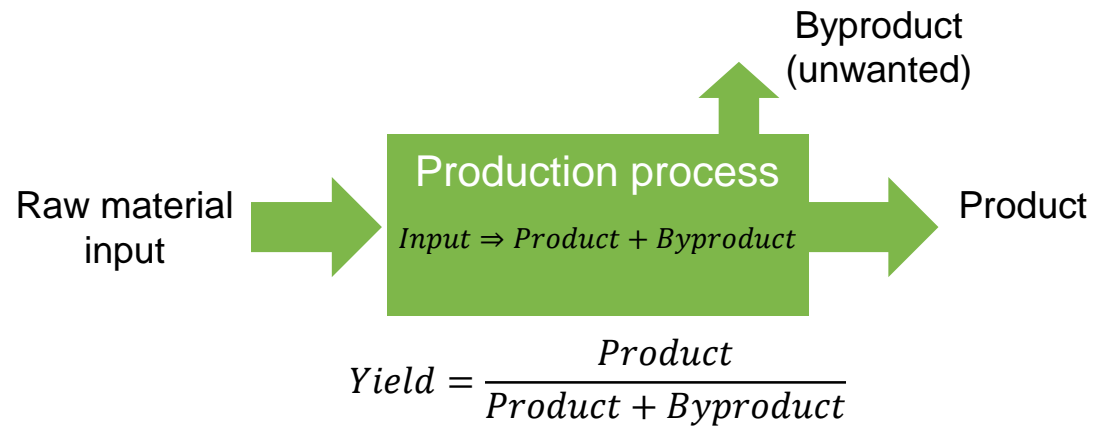
The BASF Antwerpen site is the **second largest production site of the BASF Group**. It's 55 production plants mainly consist of large scale continuous processes and produce commodity chemicals.

It's size provides sufficient critical mass for sustaining site-central expertise teams related to manufacturing and it's directly supporting functions. The highly integrated site (product streams, utilities, logistics) in combination with the presence of third parties provides a unique set of challenges

Context: the problem

Project context

- Continuous production process
- High raw material cost
- Production of unwanted byproducts reduces process efficiency
- Process efficiency varies over time



Illustrative picture of continuous production plant in BASF Antwerp

What data do we have?

- PIMS system: sensor data, stored in a big database with time series (vectors)
- LIMS system: lab data, stored in a separate database with sample times and lab values
- Operator logbook: manual text entries logging specific actions with a timestamp

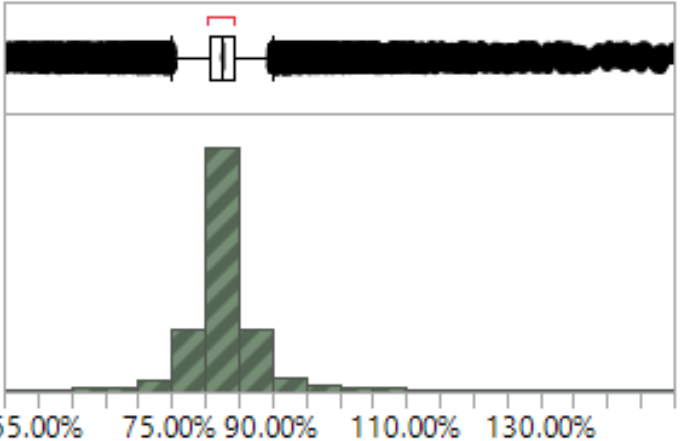
Context: the problem in numbers

Project goal

- Improve average production yield – improvement potential unknown at start of the project

Distributions

Yield

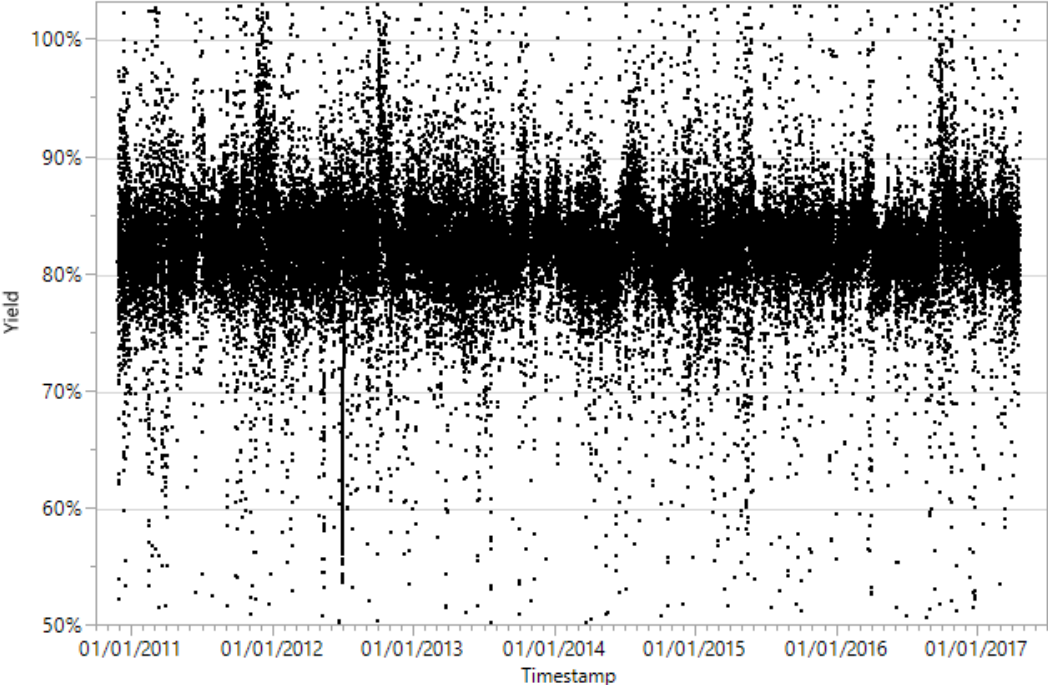


Summary Statistics

Mean	0.8272471
Std Dev	0.0731346
Std Err Mean	0.000312
Upper 95% Mean	0.8278586
Lower 95% Mean	0.8266355
N	54943

Graph Builder

Yield vs. Timestamp



Let's solve the problem – attempt 1

- Lots of data available
 - ▶ 5 years data (hour values)
 - ▶ 250+ sensors
- Use statistical algorithms and data analytics techniques to identify key variables for the yield

$$\text{Yield} = f(X1, X2, X3, \dots)$$

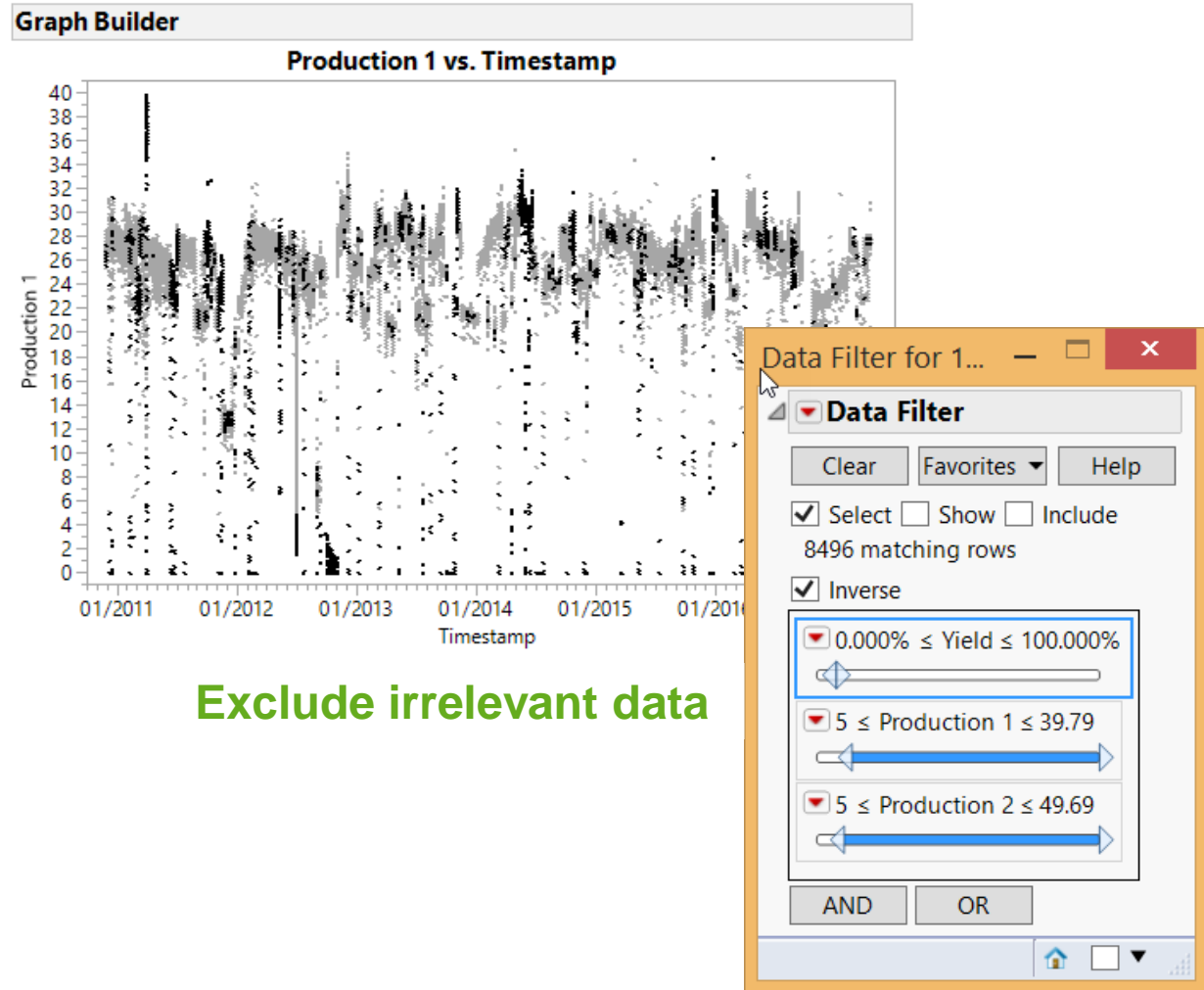
Select X1, X2, X3,... from a set of 250+ variables



Discovering key process variables

Attempt 1 – data crunching

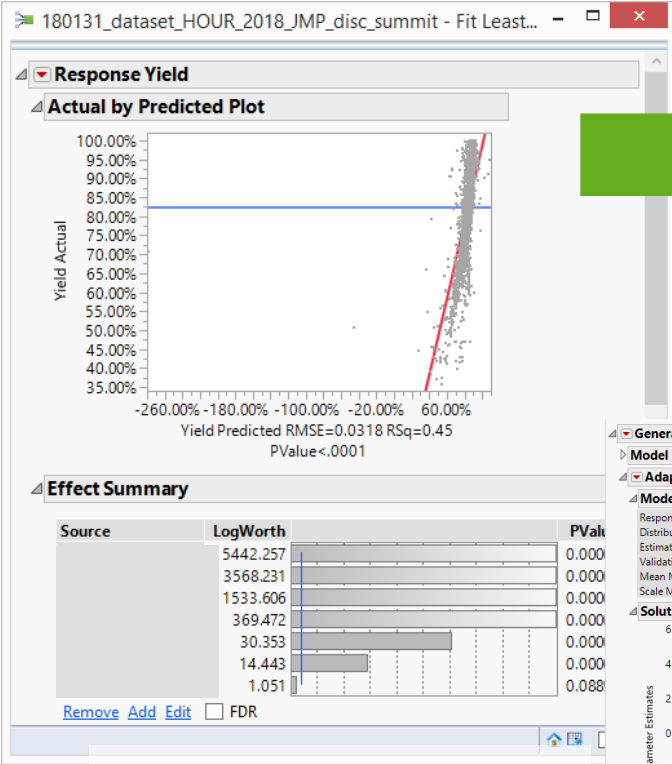
- Collect 5 years data for all measurements related to that part of the production process (online sensor + offline lab data)
→ over 250 variables
- Exclude irrelevant data: no or very low production output, yield < 0% or yield > 100%
- Find root causes for yield variation: try various statistical models
 1. Stepwise OLS
 2. Partial Least Squares
 3. Generalized regression Enet (JMP pro)
 4. Bootstrap forest (JMP pro)



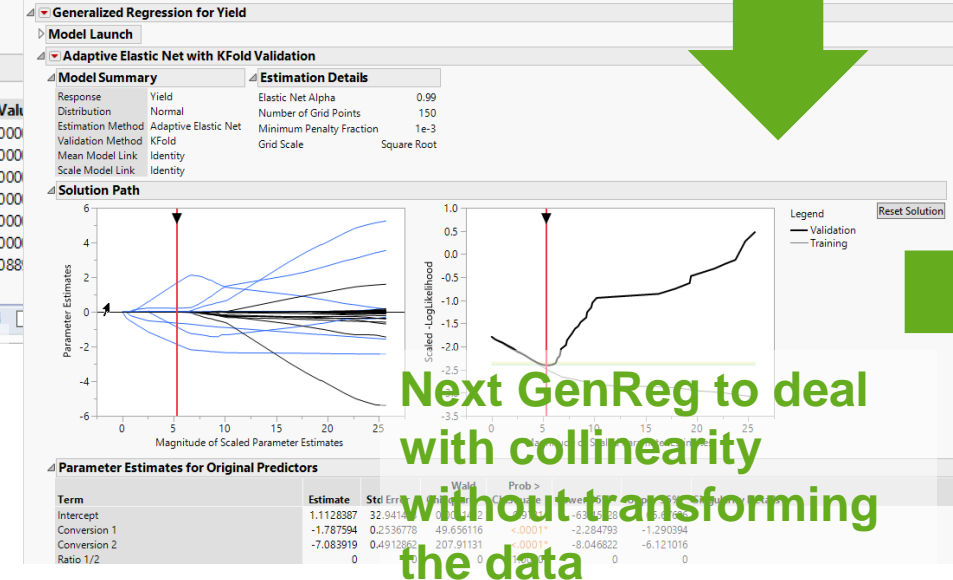
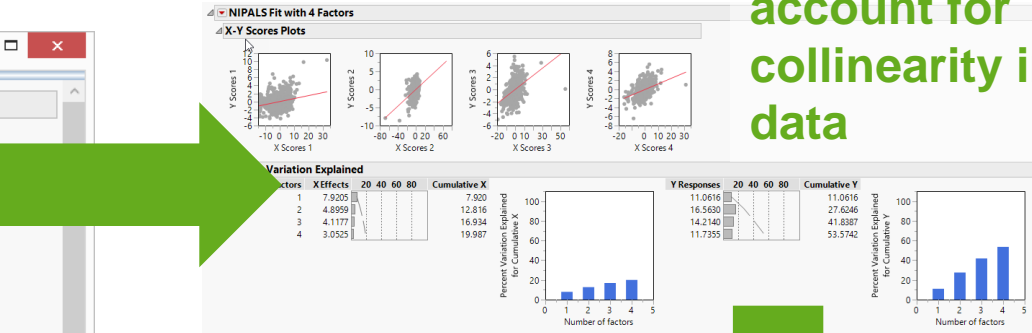
Discovering key process variables

Attempt 1 – data crunching

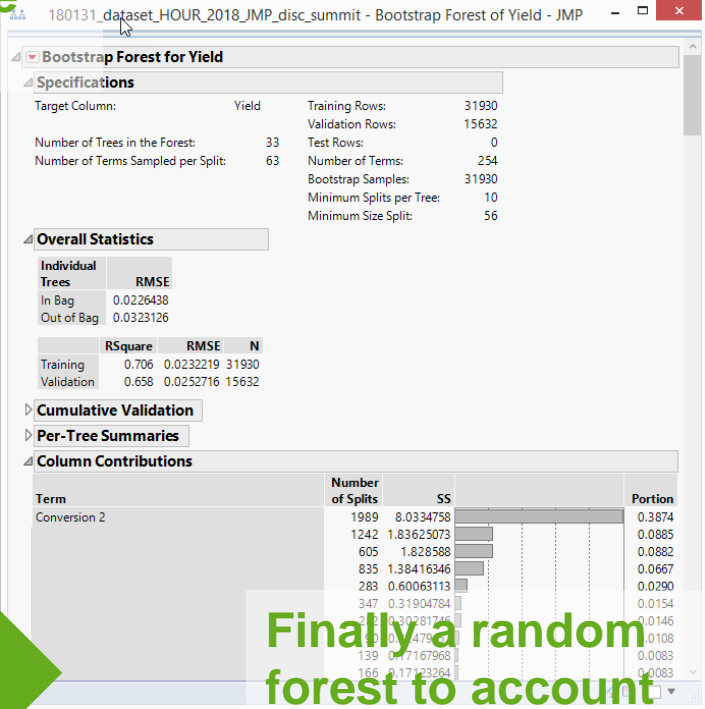
Then PLS to account for collinearity in the data



Start with OLS



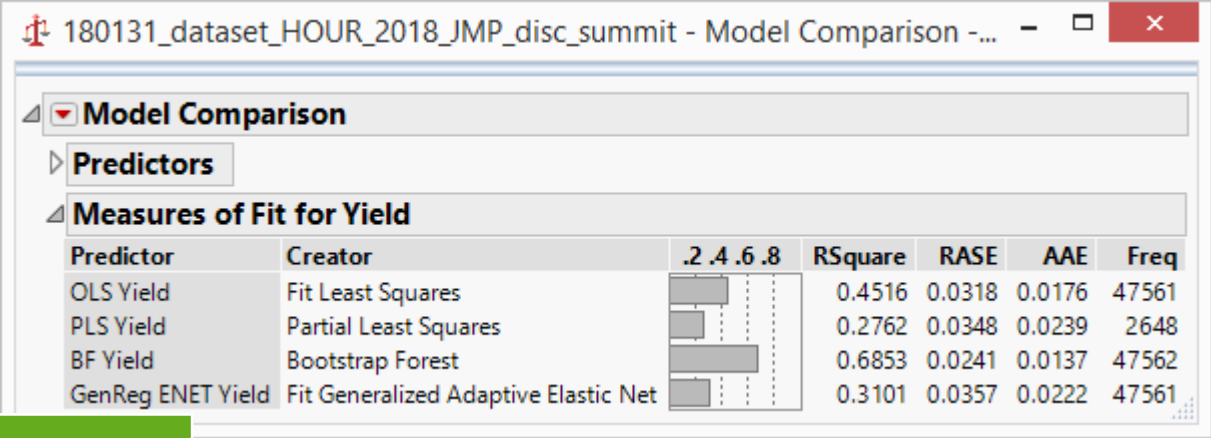
Next GenReg to deal with collinearity without transforming the data



Finally a random forest to account for interactions and non-linearities

Discovering key process variables

Attempt 1 – data crunching



OLS	PLS	GenReg	RF
Conversion 2	Conversion 2	Conversion 2	Conversion 2
Production 2	Total production	Conversion 1	Total production
Flow_1	Conversion 1	Production 2	Ratio 1/2
Level_1	Level_1	Feed 2	Conversion 1
Pressure_1	Temp_1	Quality_1	Level_2

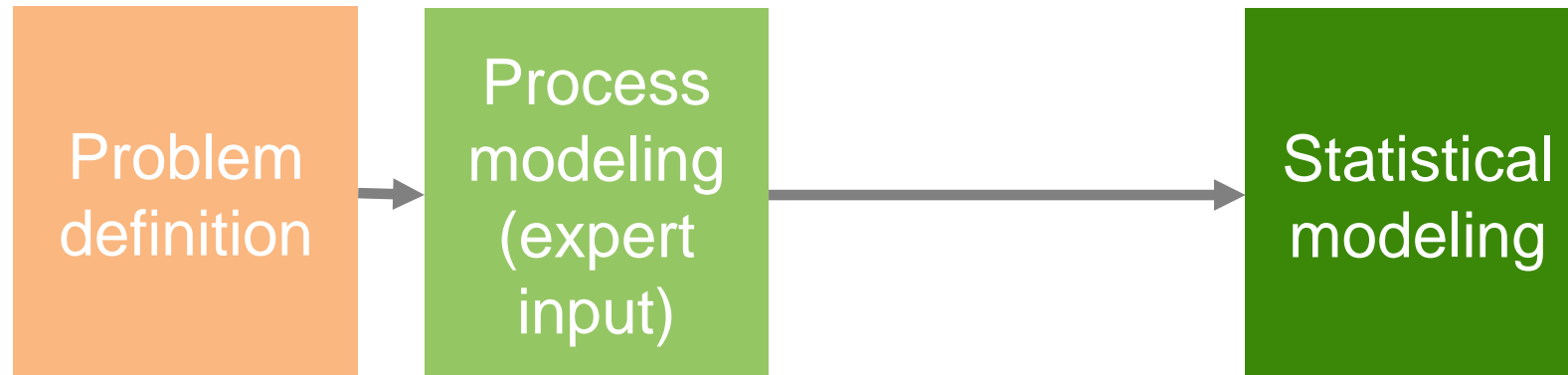
- None of the models performs great (especially the RASE/RMSE is too large compared to the target) – note in the actual study training/validation and test data has been used
- Key process variables in each model are different (except for “conversion 2”)
- Some of the parameters cannot be explained from a expert point of view (creates skepticism)
- Although there is some predictive power, none of the models is good enough to optimize the process (what are the ideal process settings?)

Let's solve the problem – attempt 2

- Ask the subject matter expert: what do they think X1, X2, etc. is?

$$\text{Yield} = f(X1, X2, X3, \dots)$$

Select X1, X2, X3,... from a set of preselected (SME input) variables



Discovering key process variables

Attempt 2 – ask the process expert

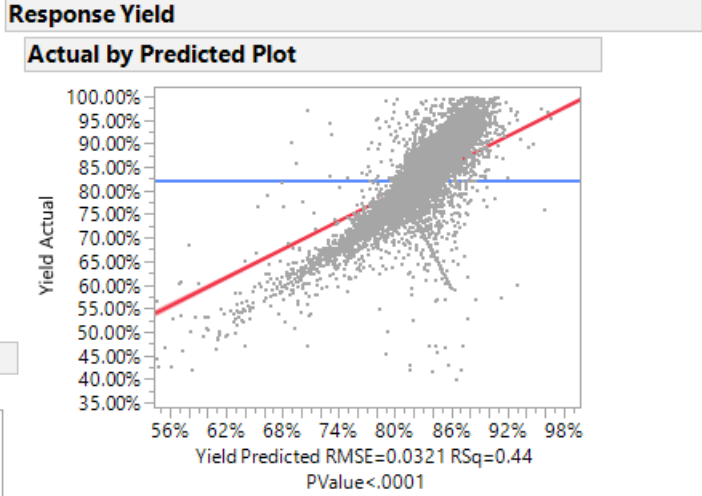
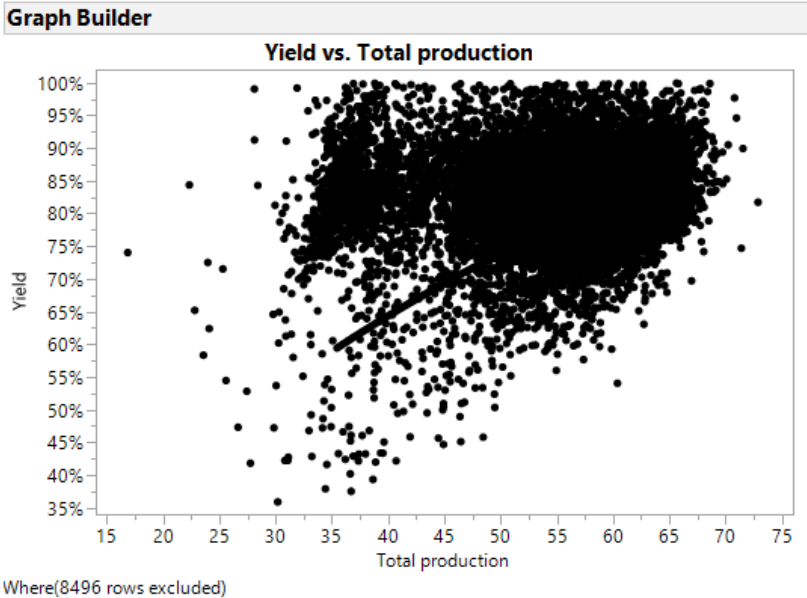
- Extensive interviews with plant management, plant operators and technology experts
 - ▶ Shortlist of suspected key process variables
 - ▶ Three categories
 - Measurement noise
 - Production (production planning, hard to change)
 - Other process settings

System effect	
Control range	Eq. to addition of noise when using IT controller
Process gain	Concentration not measured
Time constant	Aggregated response time to control changes by use of control
Dead time	
Production	Concentration control
Reaction section (R1)	
Temp. cooling water inlet	Reaction exothermic
Temperature control	Low temperature to reduce exothermic reaction
Pressure control	Controlled mainly by off-gas section & input to ammonia synthesis?
Flow rate	Reaction is slow production - distribution with and among sections
Concentration in air	
CO ₂ control	→ dependent with water in feed
Temperature control	→ 100 °C
Distillation section (D1)	
Temp. control	→ 100 °C in all temperature steps
Pressure control	→ 100 - 1000 hPa in feed
Control range	

Discovering key process variables

Attempt 2 – ask the process expert

- In an OLS model, a subset of expert defined parameters are relevant.
- Performance of the model is on par with the data crunching OLS model, which is still not good enough
- Parameters suspected to have the biggest impact (total production) is not as relevant as suspected



Effect Summary

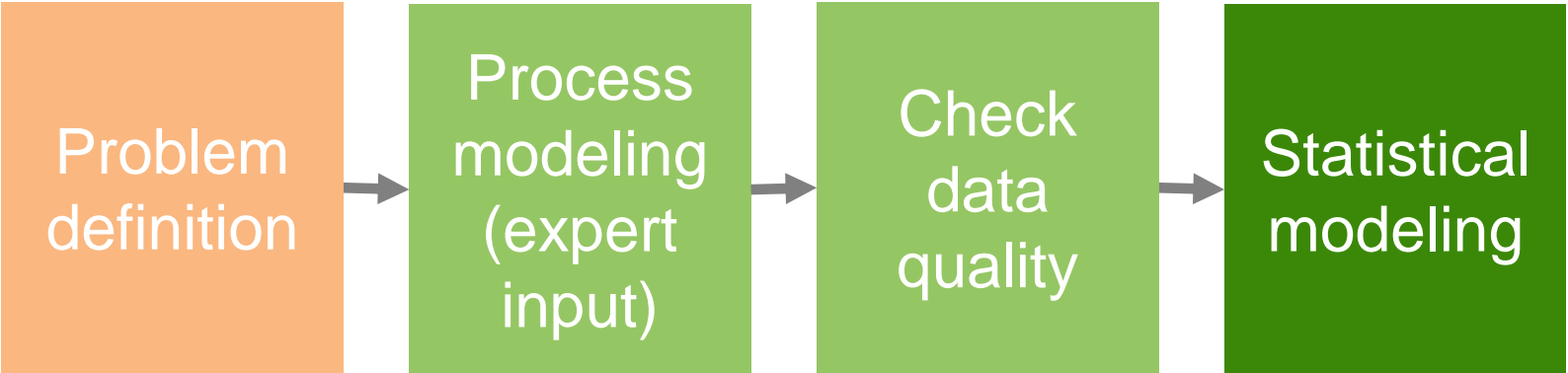
Source	LogWorth	PValue
Conversion 2	5253.518	0.00000
Production 2	2134.664	0.00000
Feed 2	1835.876	0.00000
Feed 1	27.497	0.00000
Production 1	8.943	0.00000
Conversion 1	2.827	0.00149
Ratio 1/2	2.771	0.00169

Let's solve the problem – attempt 3

- Improve the data quality.

$$\text{Yield} = f(X1, X2, X3, \dots)$$

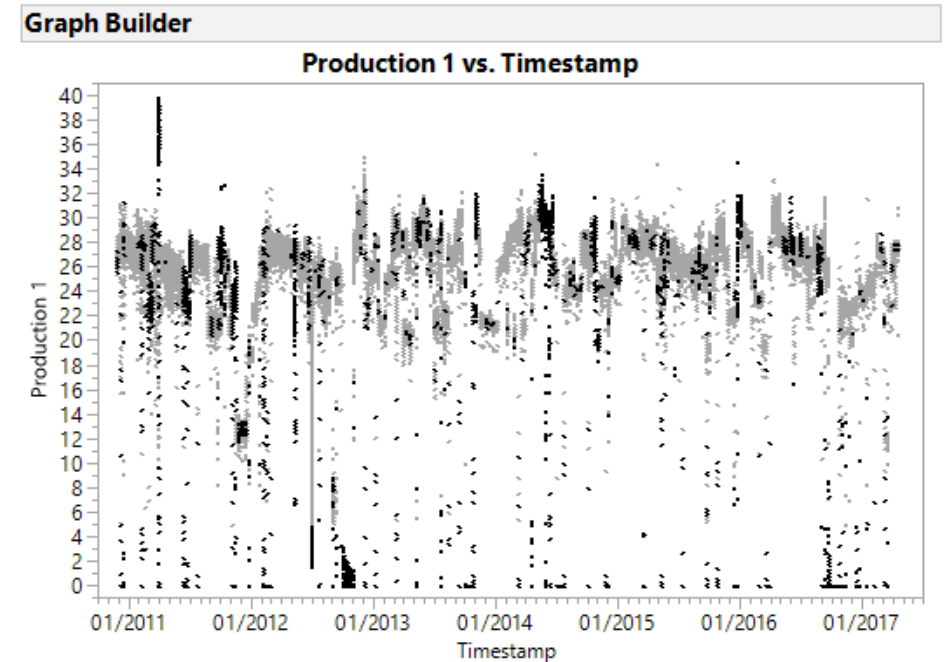
Select X1, X2, X3,... from a set of preselected (SME input) variables



Discovering key process variables

Attempt 3 - data quality evaluation

- Step 1: only look at periods with “normal” operating regimes (until know, our data cleaning was limited to this)



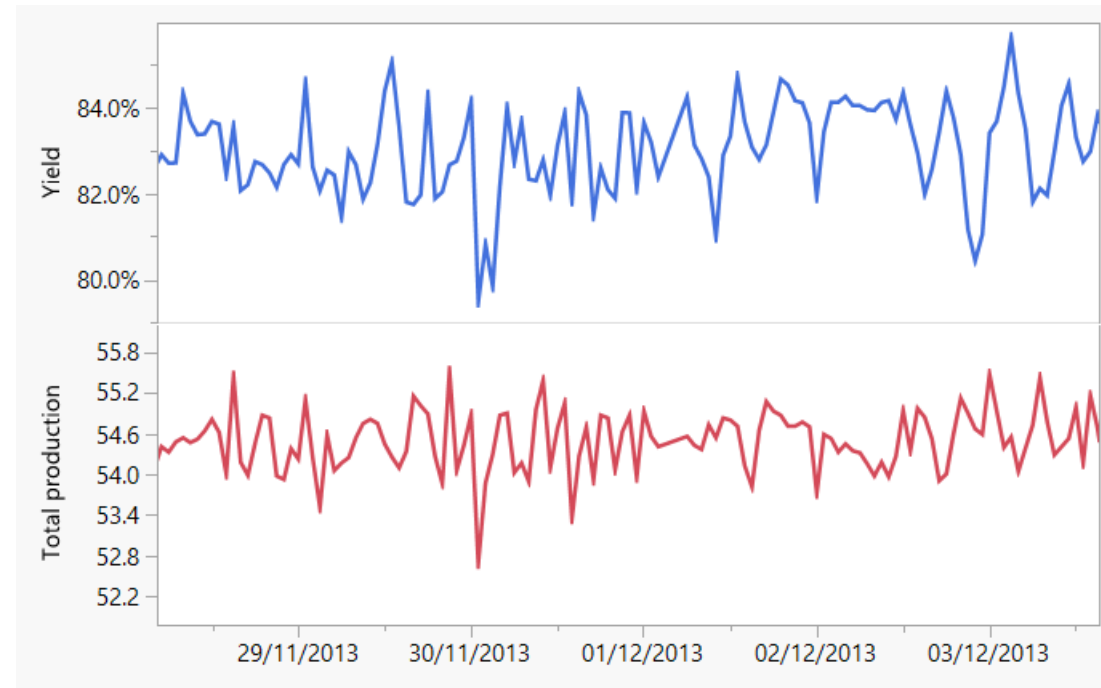
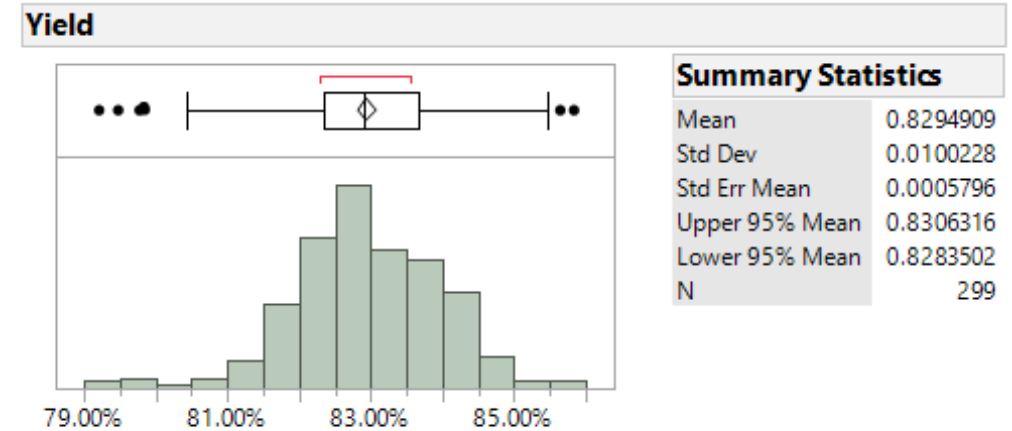
Discovering key process variables

Attempt 3 - data quality evaluation

■ Step 2: measurement noise

- ▶ Measurement system analysis: 1% variation (stddev – 15% of overall variation) in yield measurement due to variations in flow measurement

➔ Solution: look at daily (24H) averages (or medians) instead of hourly values (reducing the measurement error by approximately 5)



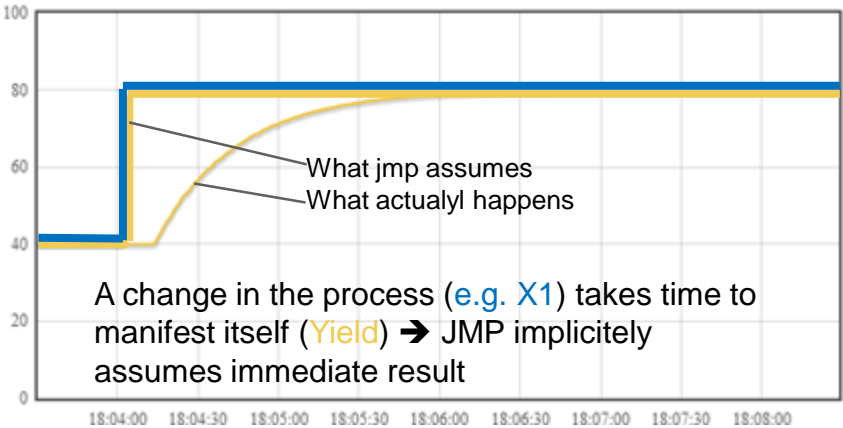
Discovering key process variables

Attempt 3 – data quality evaluation

- Step 3: dynamic effects

- Dynamic effects: after changing the process, it takes up to 48hours to get to a new steady state condition (and often another change is made within that time → seldom at steady state)

→ Solution: formula column to identify moments where the process is stable for at least 48H (look at overall 48H stddev of all major production flows)



Stability raw material 1 last 48H²

+ Stability production 1 last 48H²

+ Stability raw material 2 last 48H²

+ Stability production 2 last 48H²

+ Stability offspec last 48H²

+ Stability

Std Dev

(Raw_material_1 ,

Lag (Raw_material_1 , 1) ,

Lag (Raw_material_1 , 2) ,

Lag (Raw_material_1 , 3) ,

Lag (Raw_material_1 , 4) ,

Lag (Raw_material_1 , 5) ,

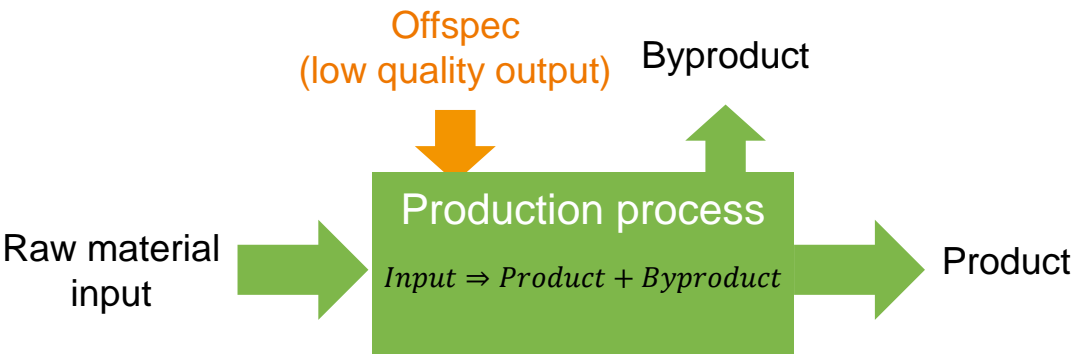
Lag (Raw_material_1 , 6) ,

Lag (Raw_material_1 , 7) ,

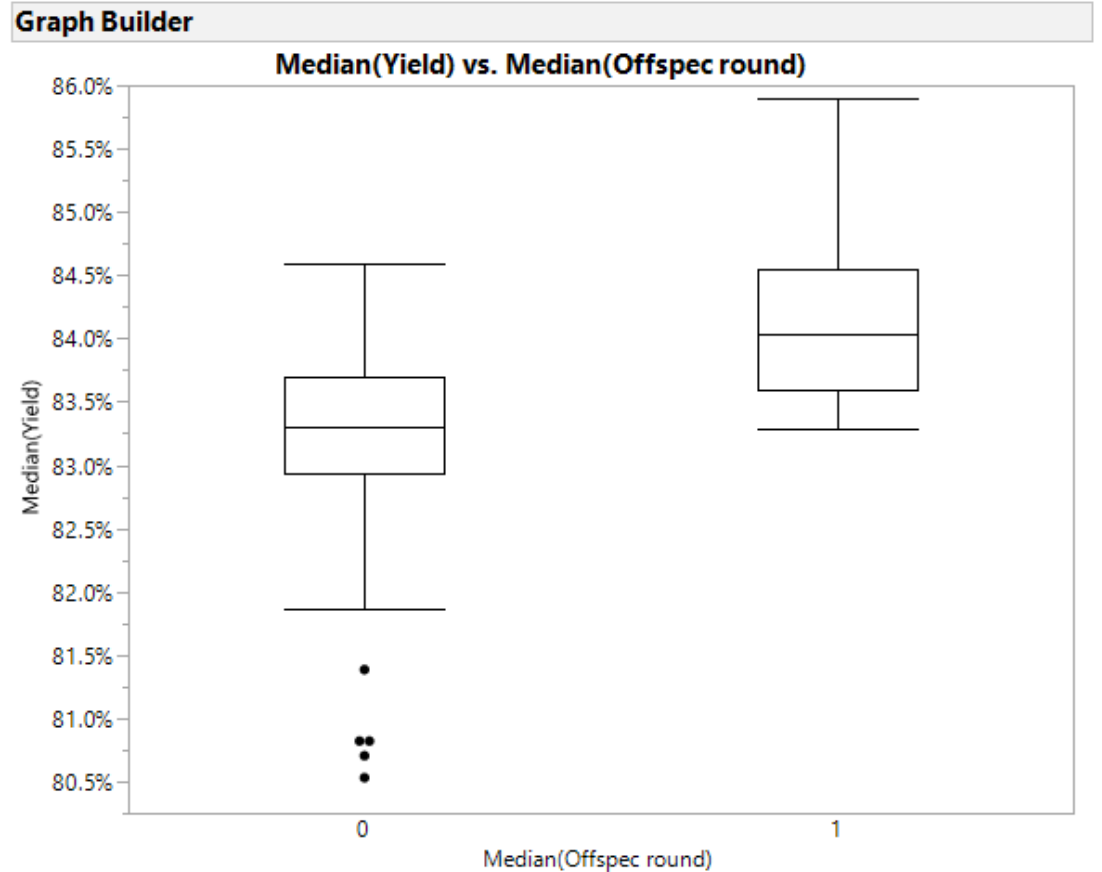
Discovering key process variables

Attempt 3 – data quality evaluation

- Step4: offspec intake
 - ▶ Intake of offspec product affects yield calculations (artificial increase of output wrt what is expected from amount of input product)
 - ▶ Note: offspec composition unknown → impact on yield cannot be quantified



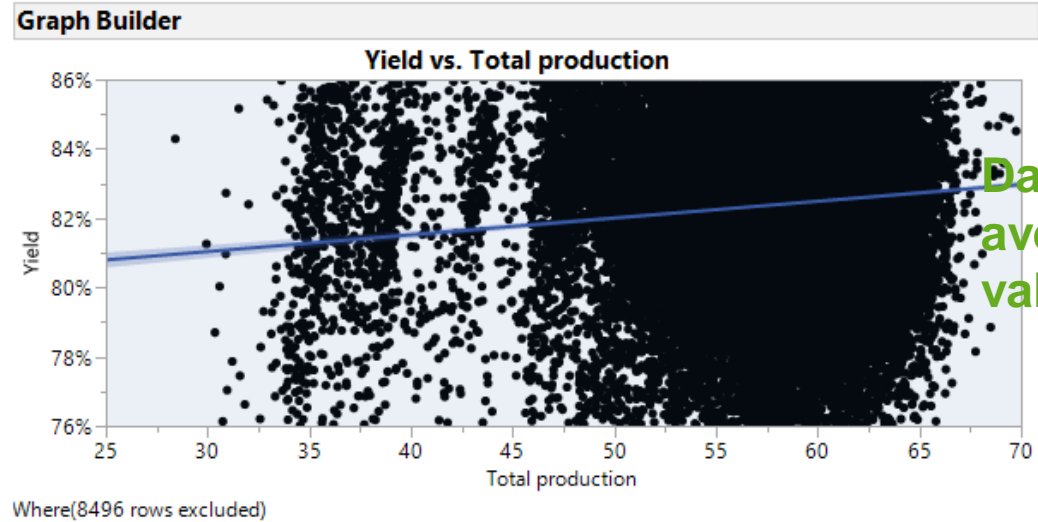
$$Yield = \frac{Product + Offspec}{Product + Byproduct}$$



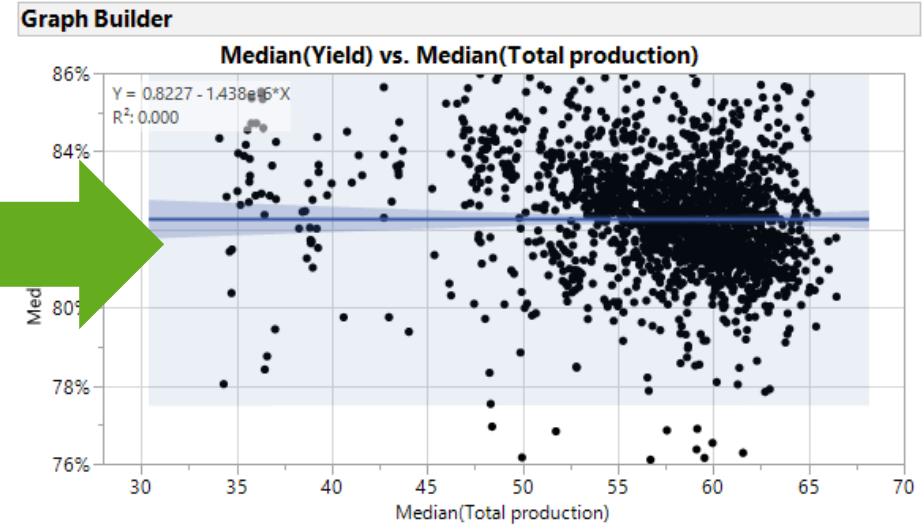
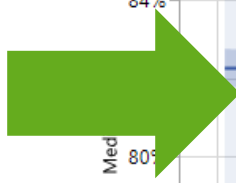
Where((Median(Offspec round) = 0, 1) and (Name("Median(Total production)") >= 45 & Name("Median(Total production)") <= 55)...))

Discovering key process variables

Attempt 3 - data quality evaluation (summary)



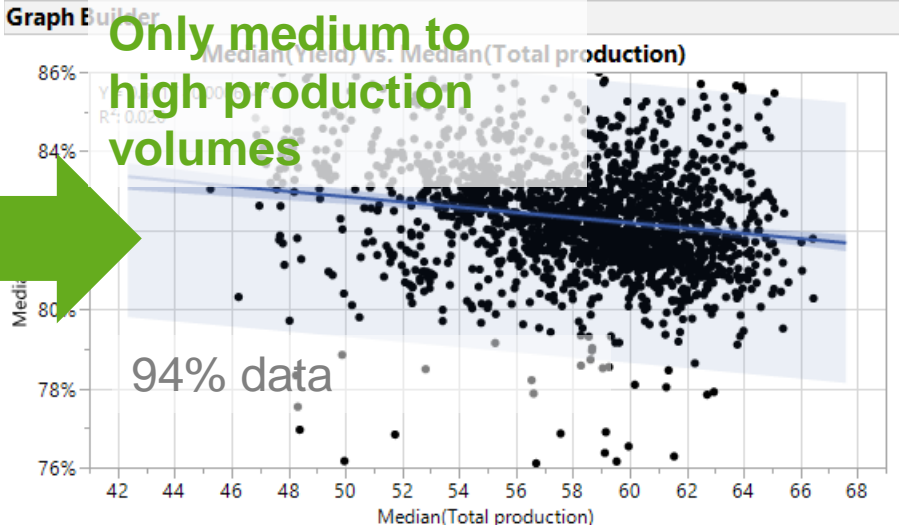
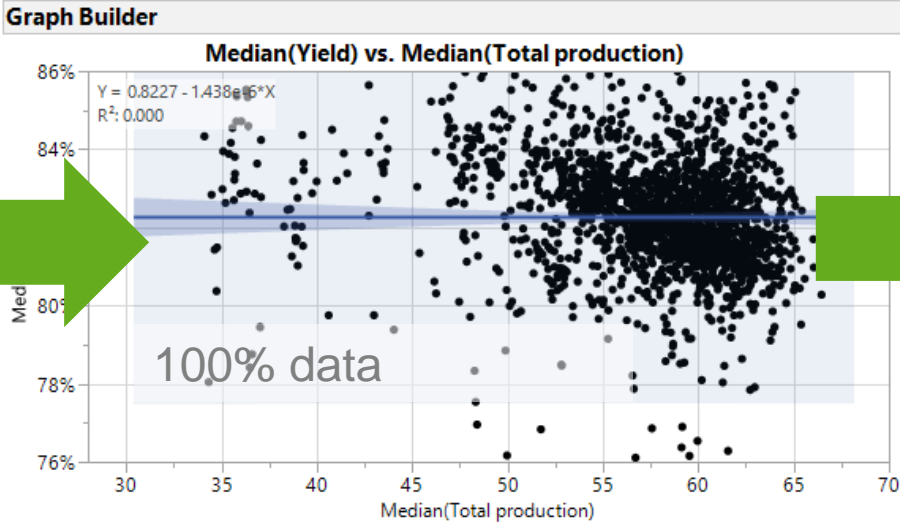
Data: 24H
average
values



Discovering key process variables

Attempt 3 - data quality evaluation (summary)

Data: 24H
average
values

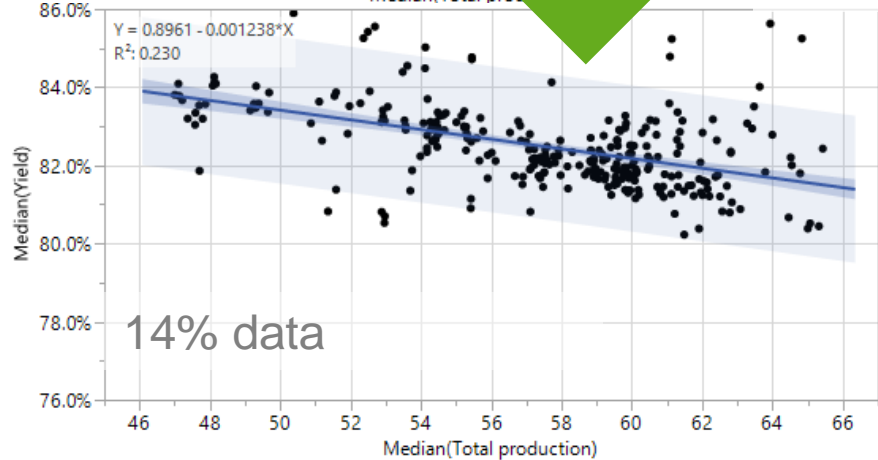
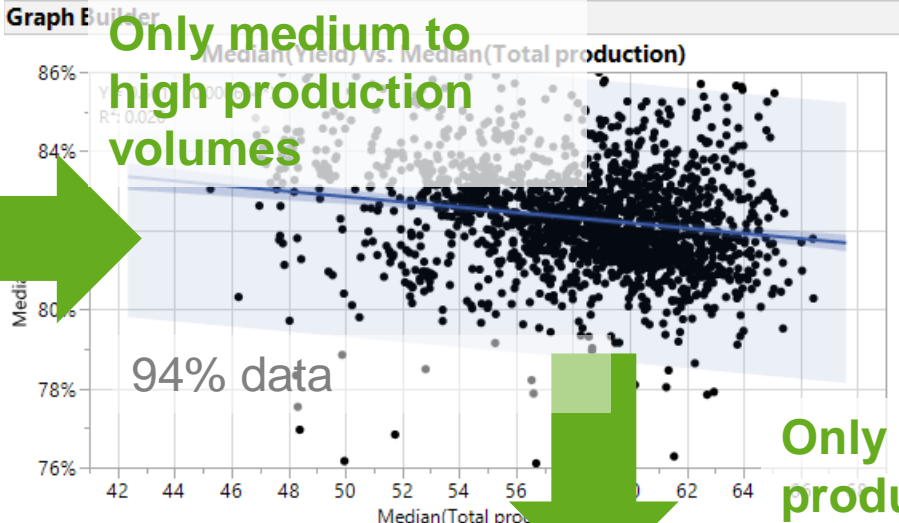
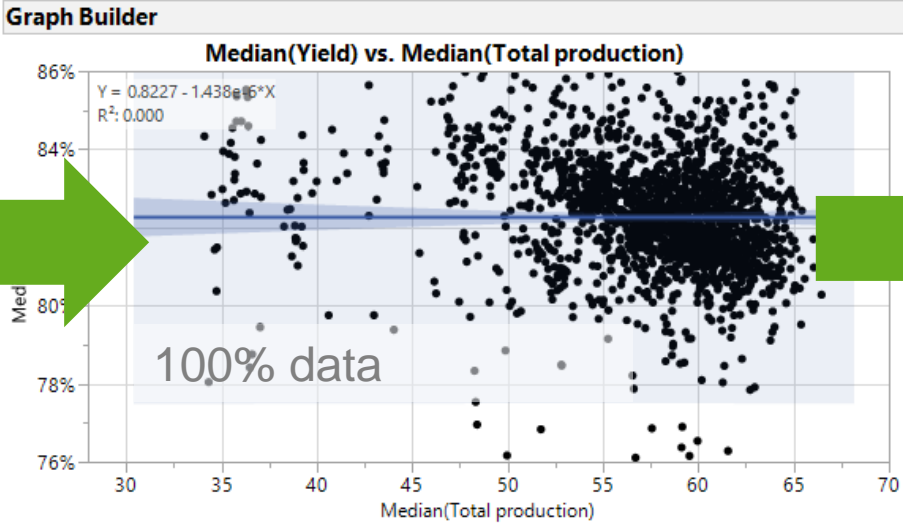


Only medium to
high production
volumes

Discovering key process variables

Attempt 3 - data quality evaluation (summary)

Data: 24H average values



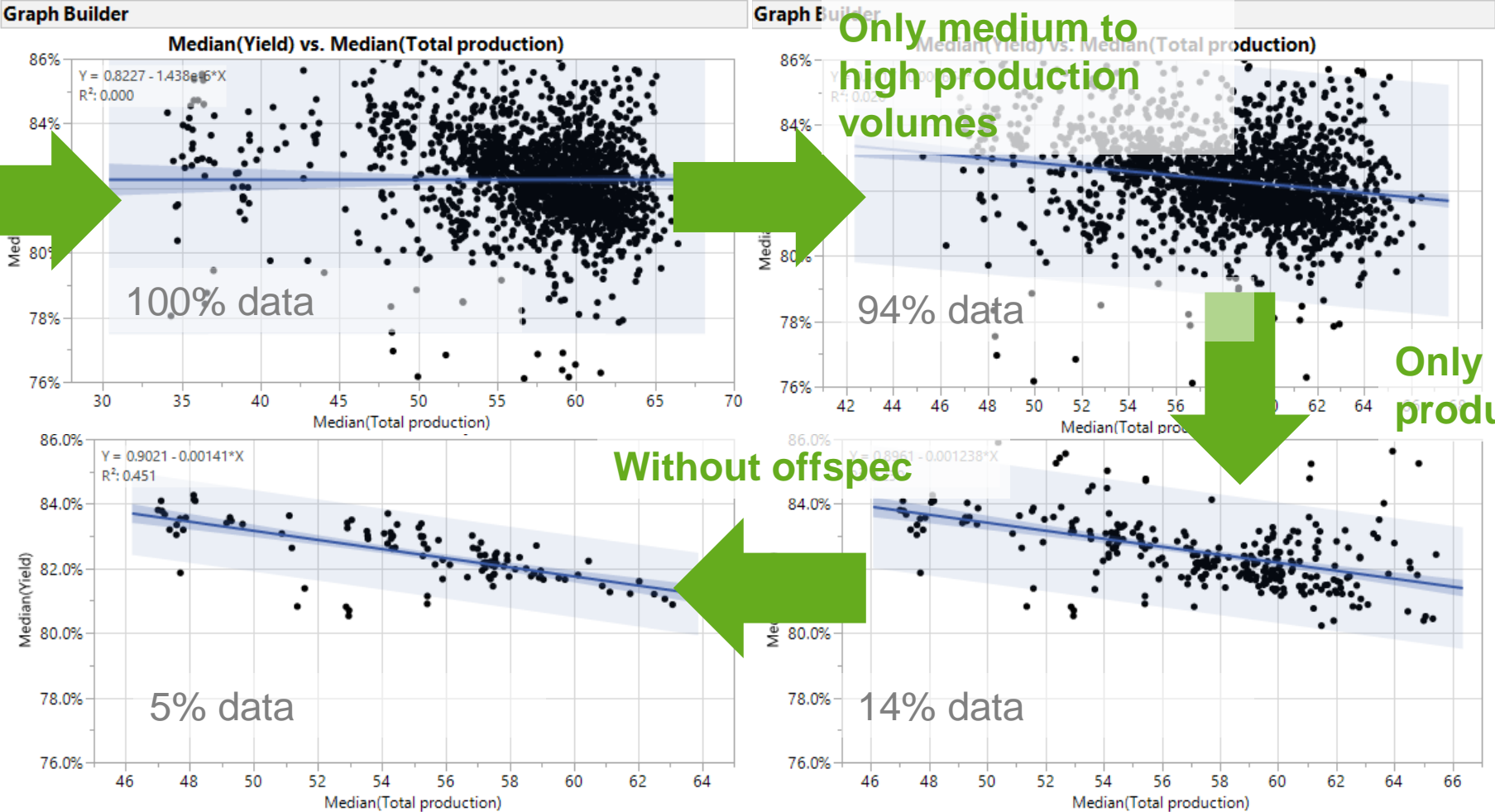
Only medium to high production volumes

Only "stable" production days

Discovering key process variables

Attempt 3 - data quality evaluation (summary)

Data: 24H average values



Only medium to high production volumes

Only "stable" production days

Without offspec

■ Only 5% of the data is high quality data!

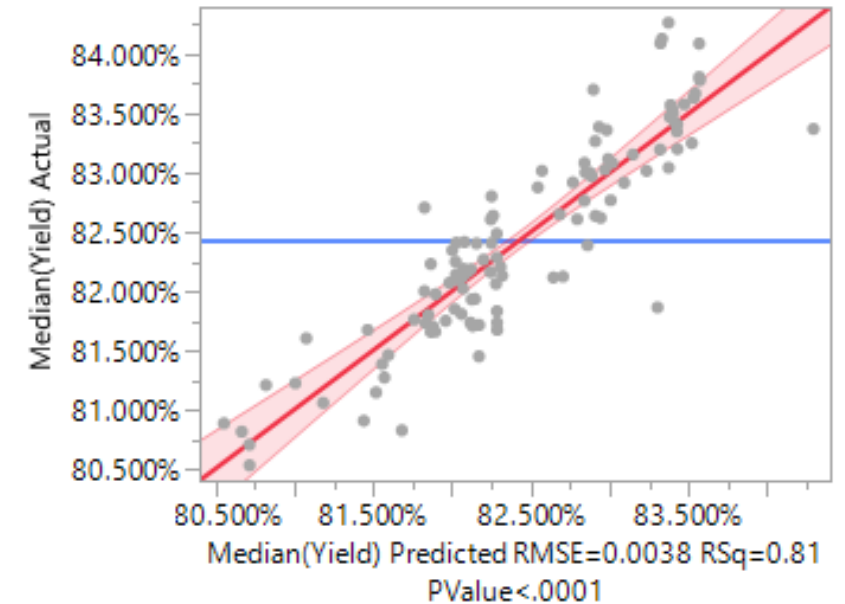
Discovering key process variables

Attempt 3 – data quality evaluation

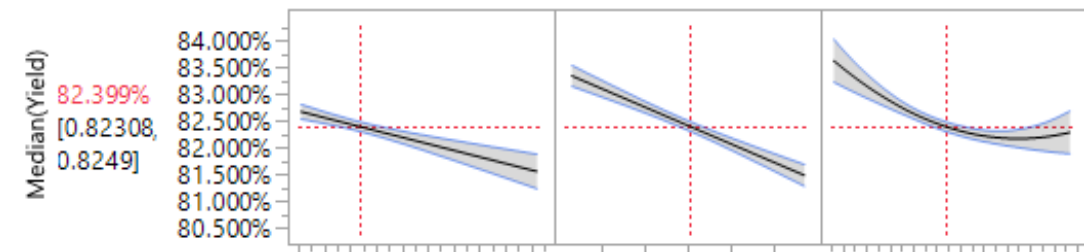
- OLS model on high quality data
- Key process variables can be explained by process expert
- Impact of key process variables can be accurately estimated (including an interaction and a non linear effect)
- Result
 - ▶ Optimal settings for yield (at a certain load)
 - ▶ Prediction of expected yield (an detection of deviations)

Response Median(Yield)

Actual by Predicted Plot

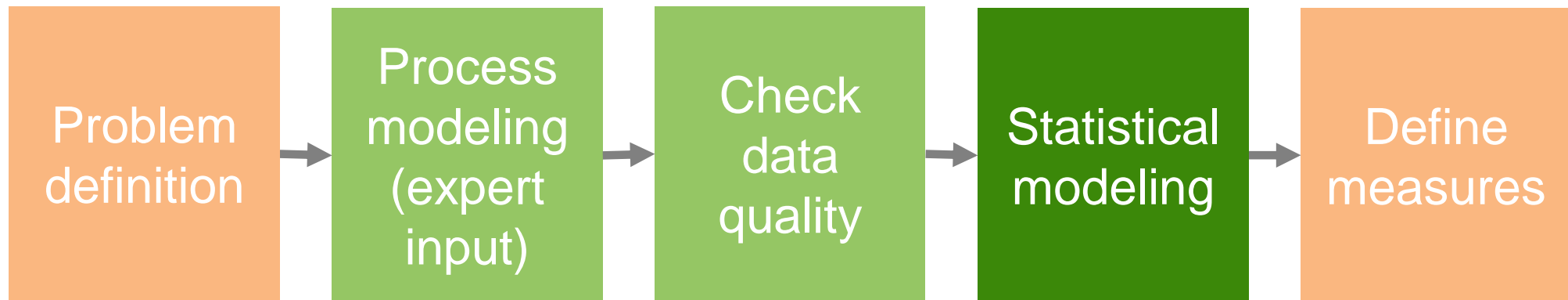


Prediction Profiler



Knowledge based action

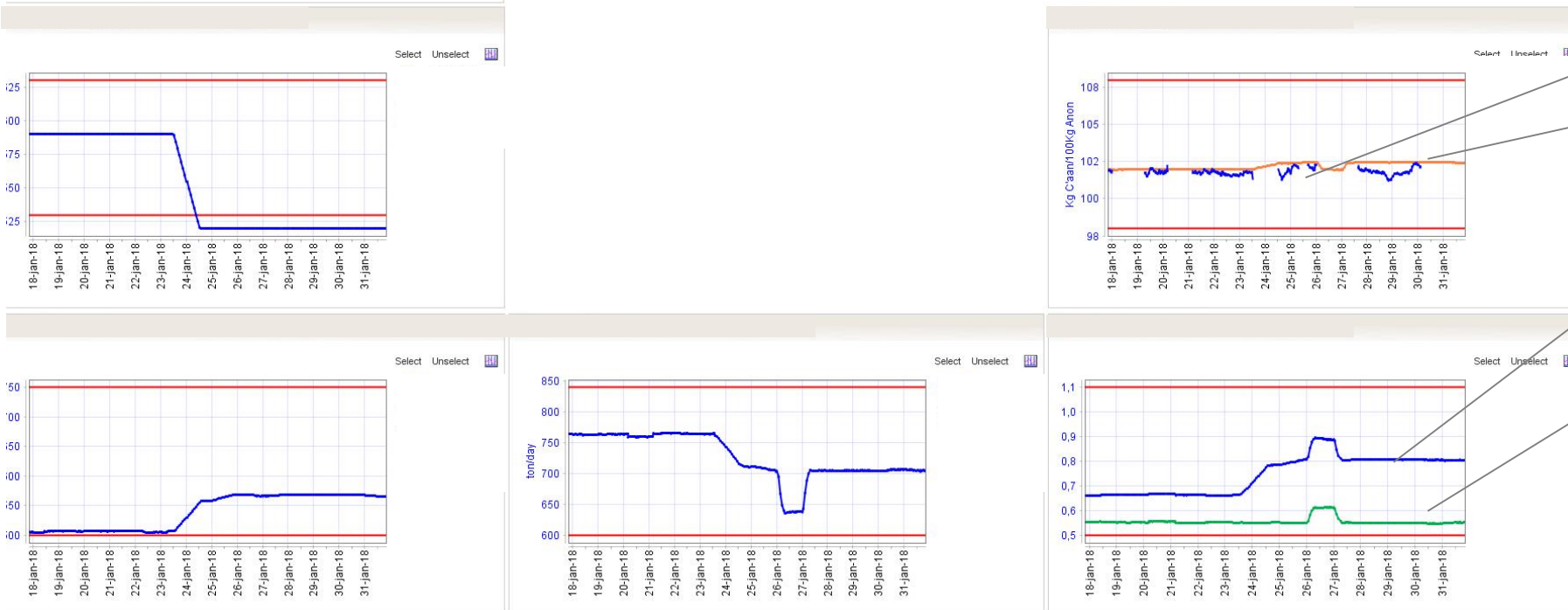
- Now we converted data into knowledge → let's act on that knowledge



Knowledge based action

- Implementation of a dashboard in the control room
- Targets show where the process should be (for maximal yield)
- The actual yield WITH indication of reliability of that value is displayed
- A corrected yield value based on other process parameters (model)

Operation Dashboard

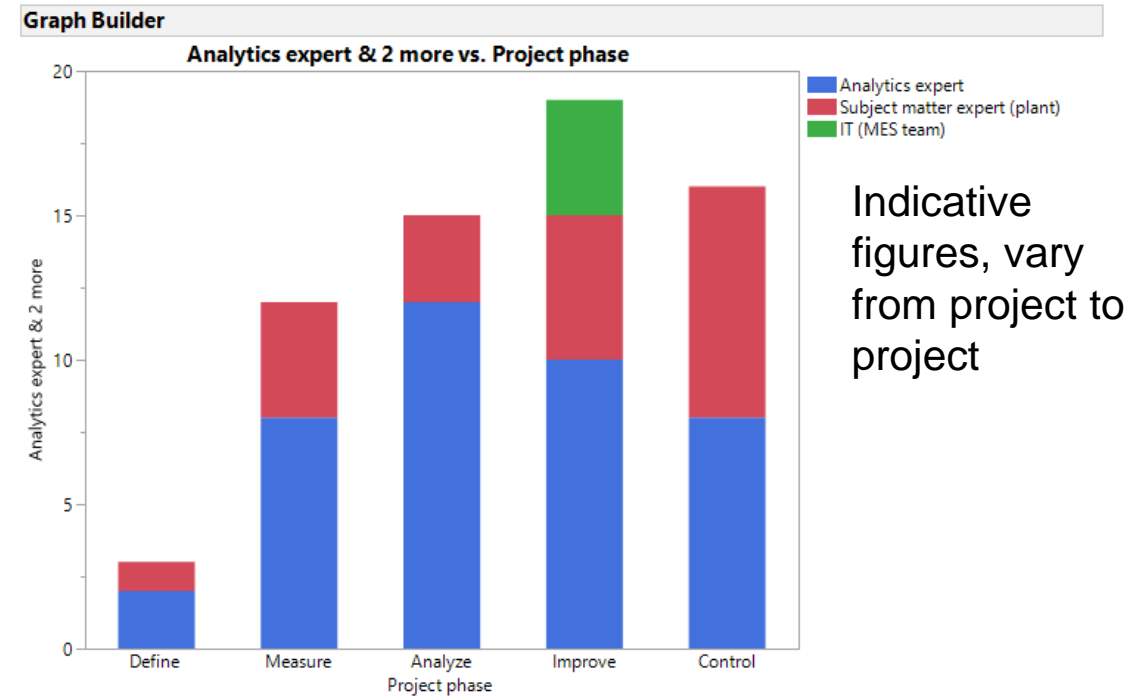


Measured yield
Corrected yield prediction

Actual process parameter
Target process parameter

Its not only about technology...

- To be succesful, the analytics part (doing the datamining and modeling) is only 25% of the work!
- Communication (in all directions) and change are major succes factors
- A project lead (in our case the data scientist) must oversee the project from end to end (clear problem definition → sustainable benefits)
- A LSS DMAIC project workflow (with a good amount of advanced analytics sauce in the measure and analyze phase) is a best practice



CONTROL

Key take aways

- Succes (creating value) =

- ▶ technology (JMP) +
- ▶ data science +
- ▶ expert input +
- ▶ thorough data cleaning +
- ▶ project based approach



- Handling time series data (which is typical for process industry) requires some specific approaches

- ▶ Data preprocessing: measurement noise, dynamic effects, ...
- ▶ Modeling: colinearity, autocorrelation between consecutive data points, ...



We create chemistry