

Industrial Data Science Workshop (L1)

19-21 September 2023

Dr. Mattia Vallerio

Manufacturing Excellence Site Manager at Solvay and
Advanced Process Control Specialist [[in](#)]

Dr. Carlos Perez

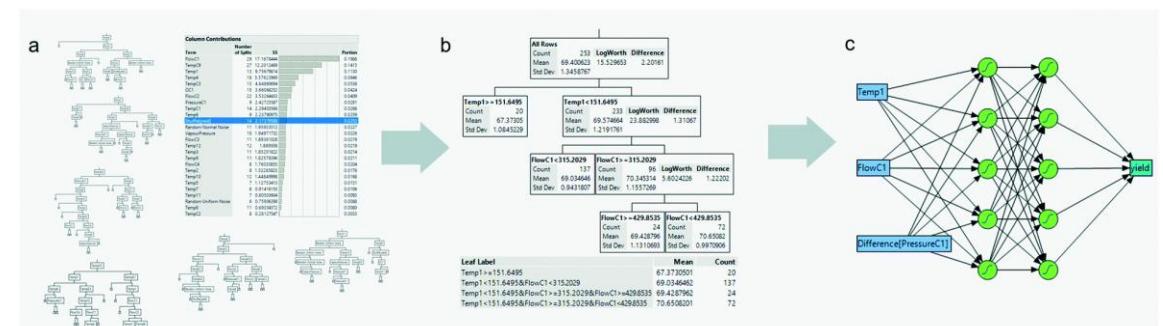
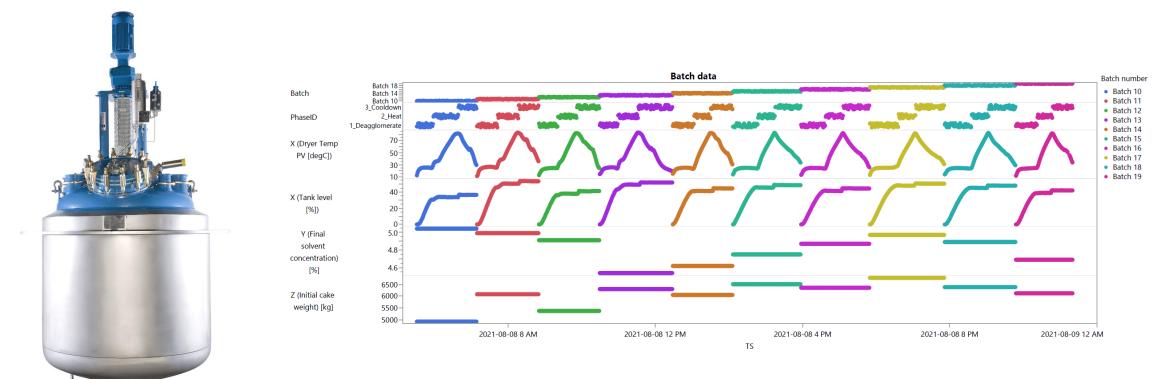
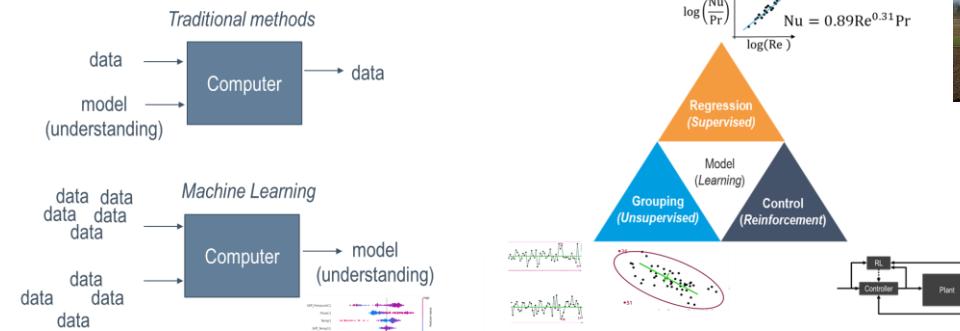
Industrial Data Scientist at Solvay and
Optimization Specialist [[in](#)]

Dr. Francisco Navarro

Sr. Data Science Training Lead at IFF and
Visiting Researcher at Imperial College London [[in](#)]

Agenda

- Day 1 – Industrial data science
 - Machine learning for process engineers
 - Industrial data and process control
 - Hands-on exercises
- Day 2 – Monitoring and screening
 - Identify relevant process changes
 - ML applied to batch processes
 - Hands-on exercises
- Day 3 – Modeling and understanding
 - Predictive modeling
 - Industrial examples
 - Hands-on exercises



Industrial Data Science
Workshop | Sept. 2023

Day 1

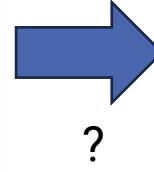
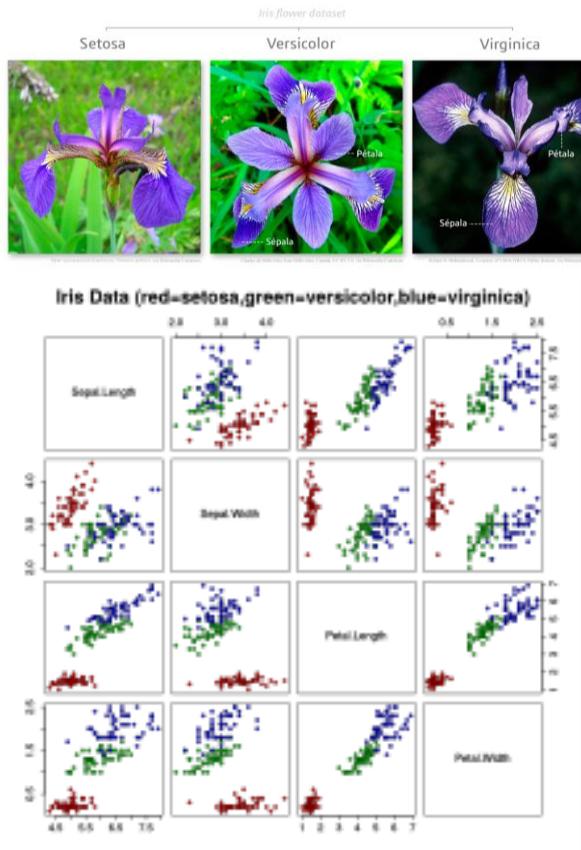
Industrial Data Science

*Manufacturing needs still are safety,
efficiency, product quality & process reliability.*

*(and not video recommendations
or predicting types of flowers)*

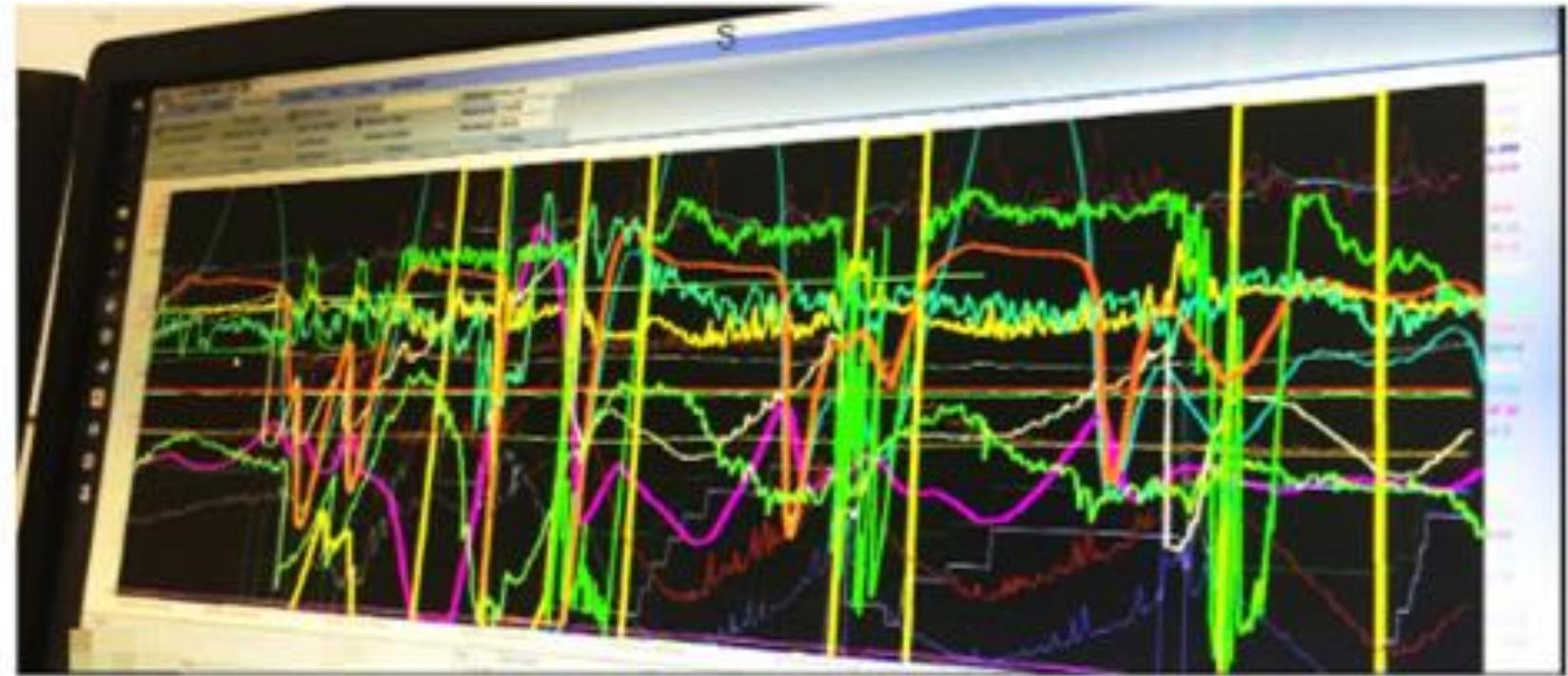


Machine learning is (often) taught with irrelevant datasets



?

101 ML example:
Types of flowers

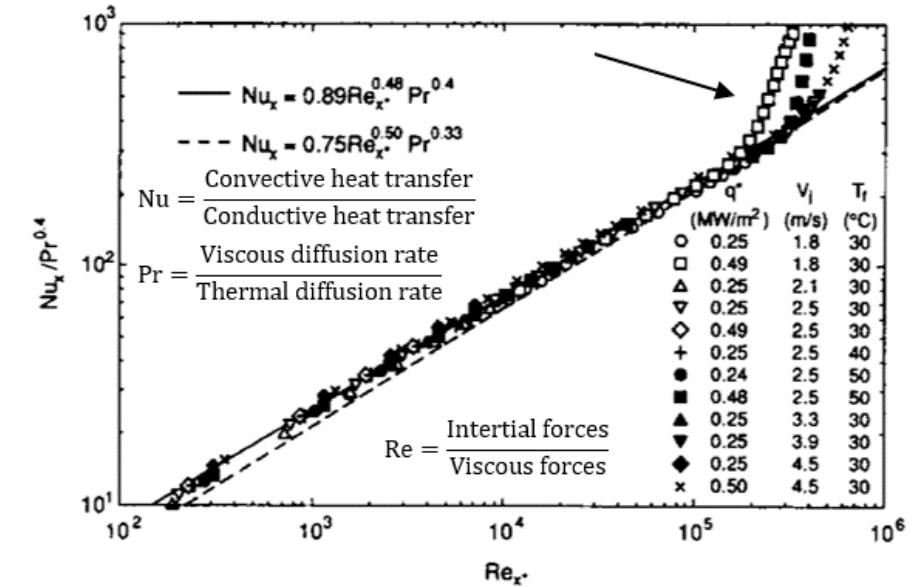
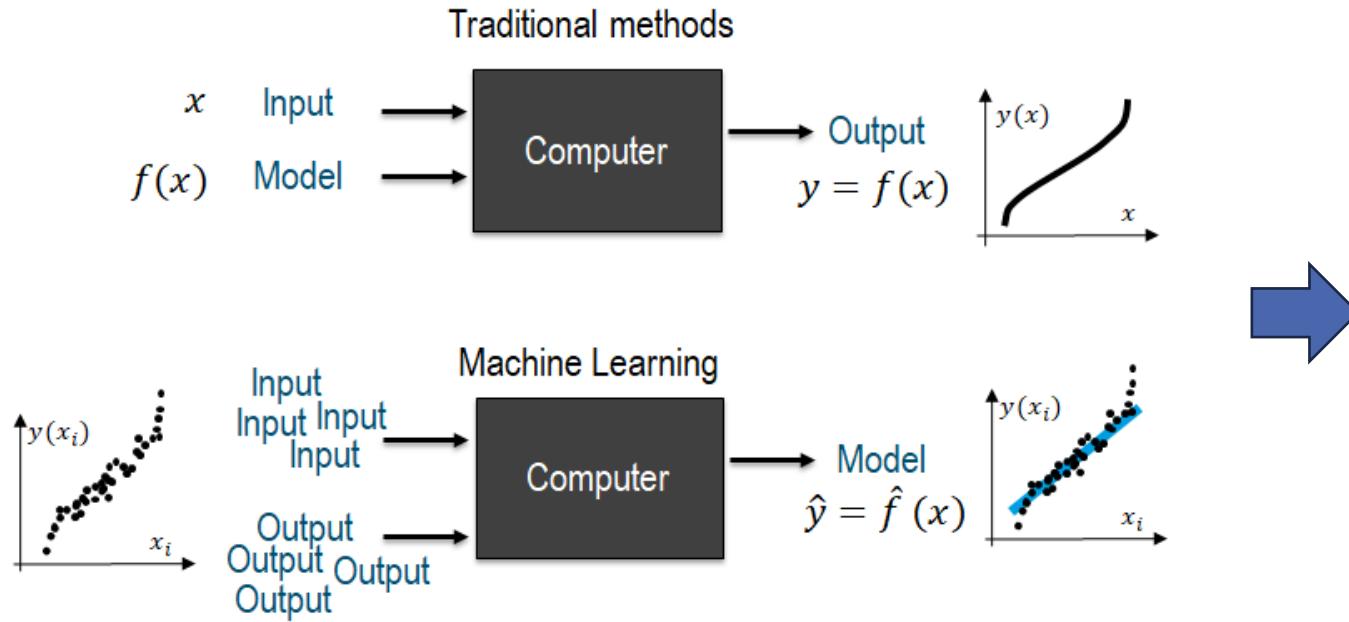


Tags in process historian (sensor data)

Image credit:

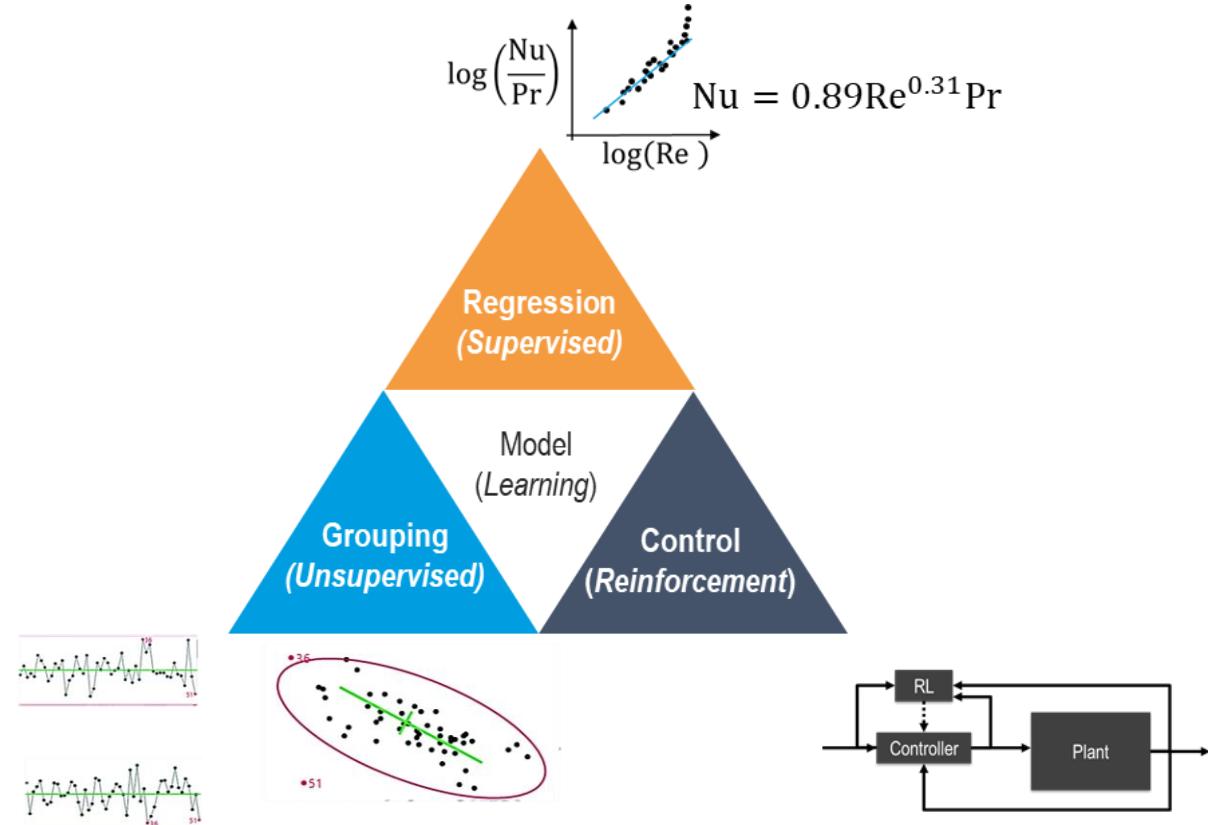
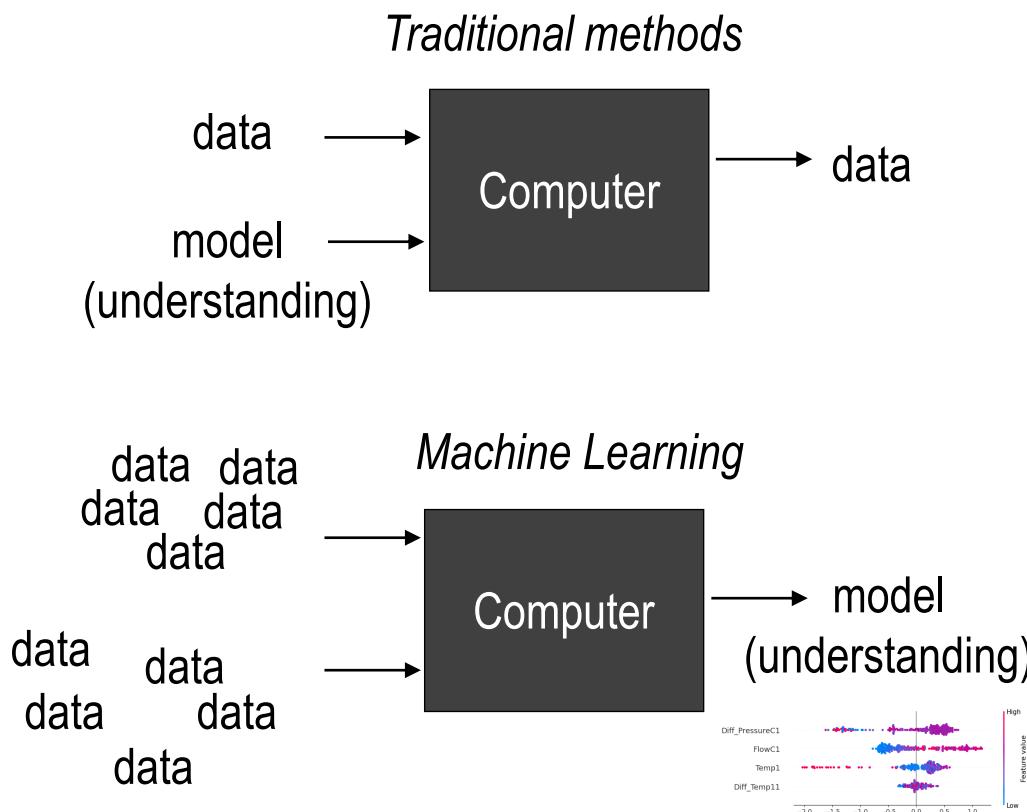
Scaling up the Use of Machine Learning in Chemical Process Industries ([2023-EU-30MP-1346](#))

Good news, process engineers already know ‘machine learning’



Traditionally, dimensionless numbers (Re, Pr, Nu, etc.,) and non-linear models were used to estimate heat or mass transfer coefficients from process data. In machine learning, these steps can be named feature engineering, feature selection, dimensionality reduction and regression (supervised learning). The risk of extrapolation is known, as models are specific to similar systems and operating conditions

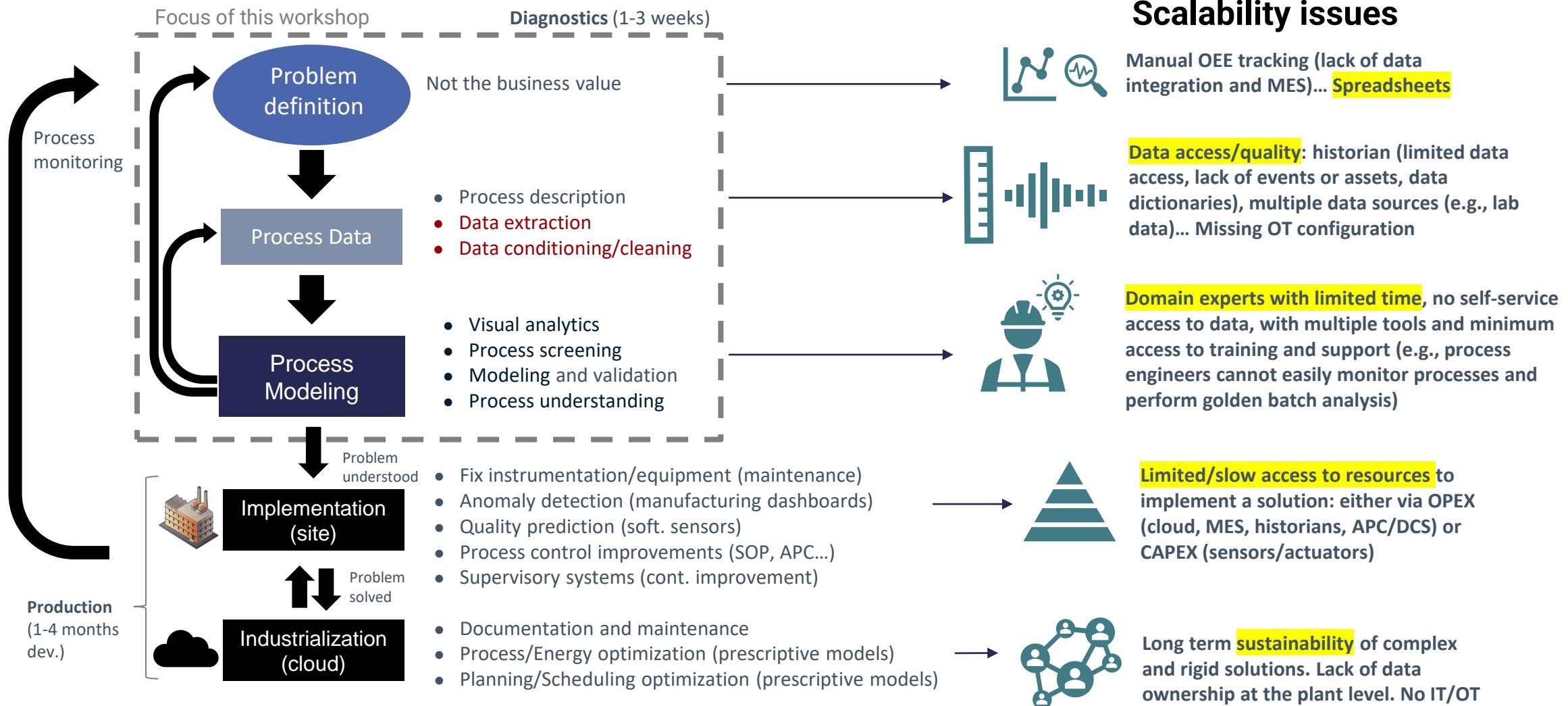
Process engineering and industrial data science in chemical manufacturing



Source:

[Industrial data science – a review of machine learning applications for chemical and process industries](#)
React. Chem. Eng., 2022, 7, 1471–1509

Data-driven decisions require process knowledge (and time)



ChatGPT agrees

Jan. 2023



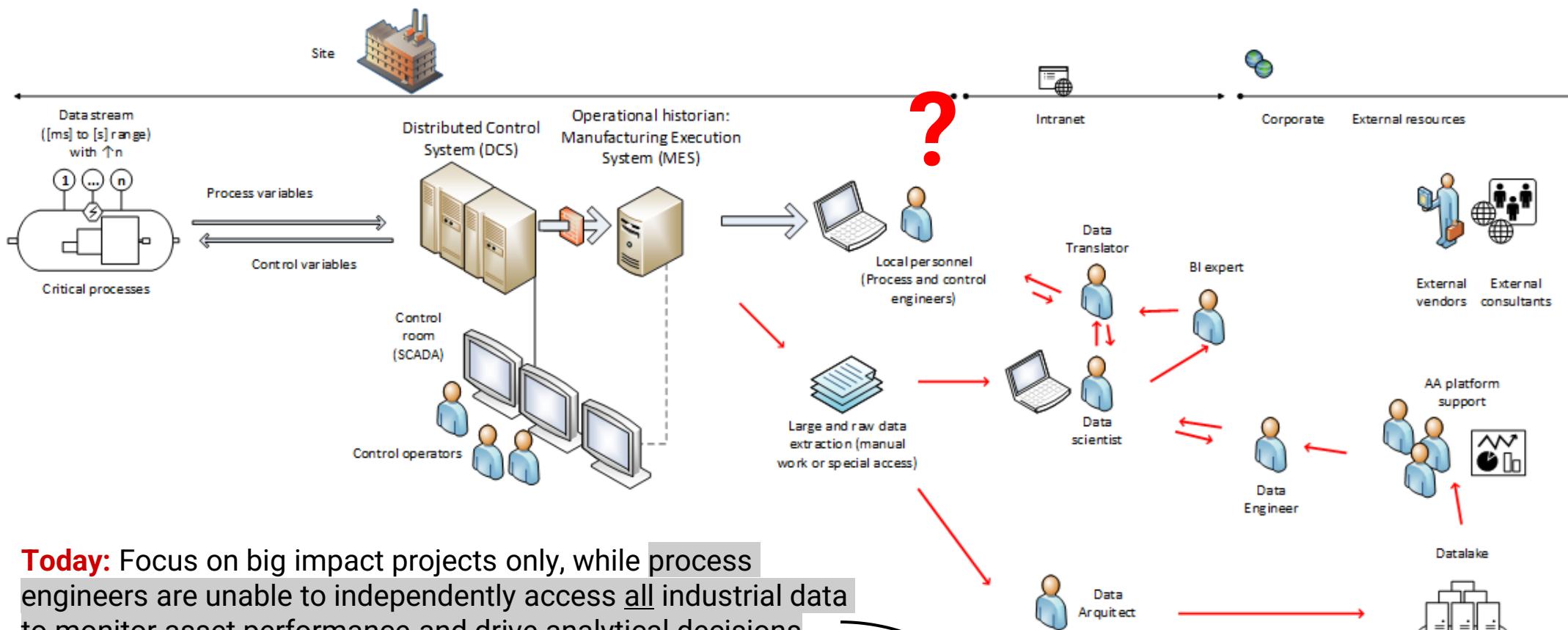
Could you describe, in one paragraph, how to use machine learning effectively in process and manufacturing industries so business impact is guaranteed?



To use machine learning effectively in process and manufacturing industries to guarantee business impact, it is important to start by clearly defining the problem or challenge that the business is facing, using relevant and high-quality data to train the model, selecting appropriate algorithms, monitoring performance, continuously improving, collaborating with experts, and communicating results to the relevant stakeholders. Additionally, process engineers should be involved in all stages of the process, from problem definition and data collection, to algorithm selection, model implementation, and performance monitoring, to ensure that the model is tailored to meet the specific needs of the manufacturing process and is delivering the desired results.



How not to scale a data-driven culture (number up approach)



Today: Focus on big impact projects only, while process engineers are unable to independently access all industrial data to monitor asset performance and drive analytical decisions.

Digital and central teams delayed by:

- 1) domain expert availability
- 2) poor data quality
- 3) limited data access

Risk: Lack of data ownership at the plant level resulting in poor data quality: missing tags, unrecorded batch events, unstructured data (assets), siloed quality data...

Small fraction of initiatives are addressed (big impact projects only by specialized teams).

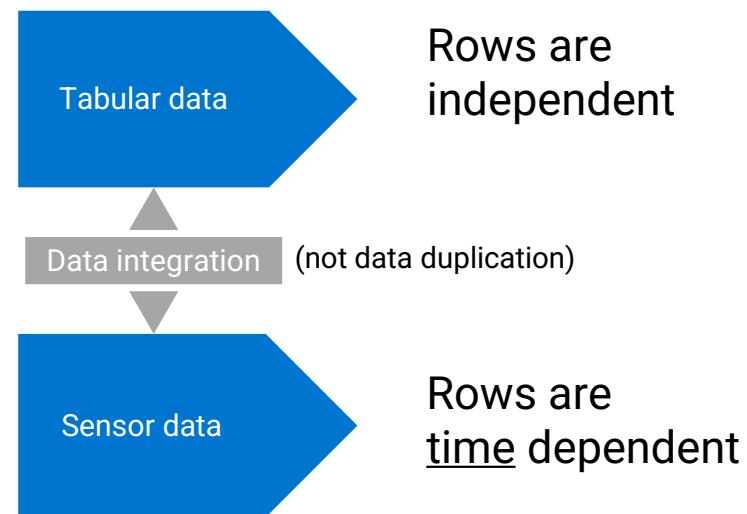
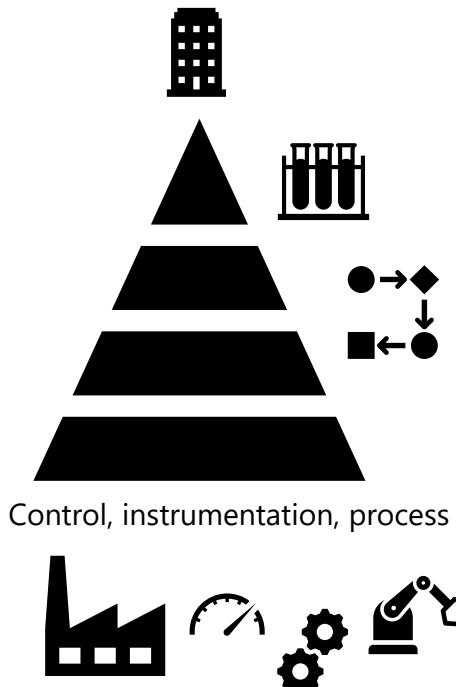
Complexity

**SELECTING
DUPLICATING
MAINTAINING {industrial data}**

through multiple tools, systems and personas



There are two* types of process data



BatchID	LotID (PO number)	Grade (P1, P2, P3)	Quality (pass/no pass)	Yield [%]
...
...

(e.g., KPIs, LIMS and SAP data)

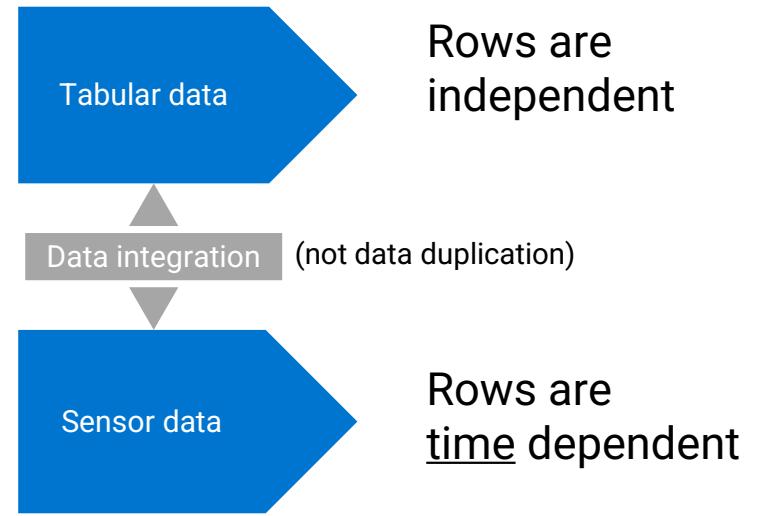
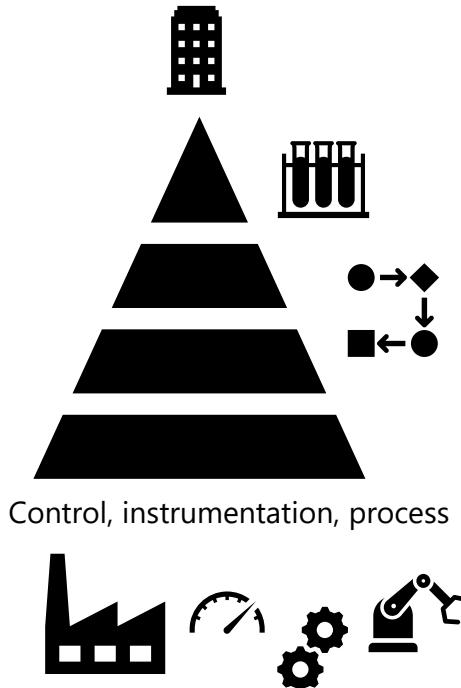
Datetime	PI_505.PV (Pressure) [bar]	TI_203.PV (Temperature) [degC]	FIC_101. MV (Valve position) [%]	FI_101.PV (Flowrate) [m ³ /min]
...
...

(e.g., tags in historians)

(* Not really true, of course ¹¹

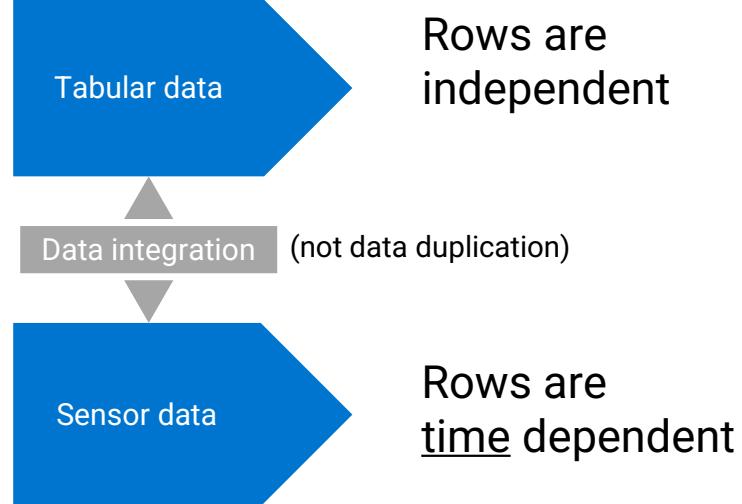
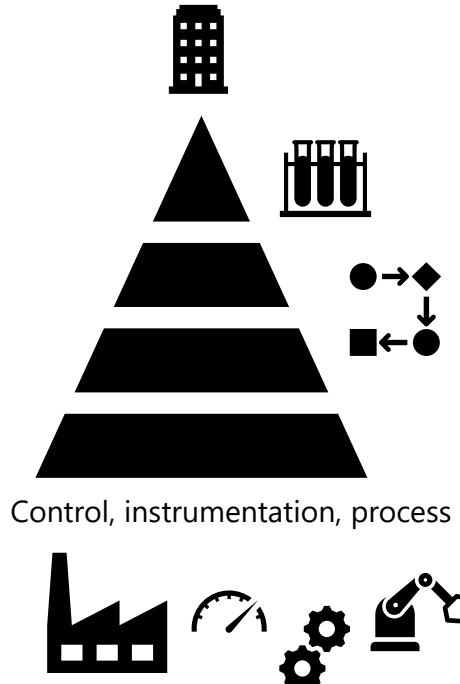
There are many, many types of industrial analytics software

+info: g2.com, trustradius.com, capterra.com, gartner.com, Insresearch.com

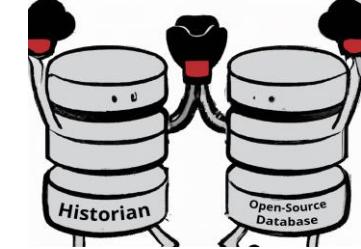


And many, many types of industrial databases

+info: g2.com, trustradius.com, capterra.com, gartner.com, Inresearch.com

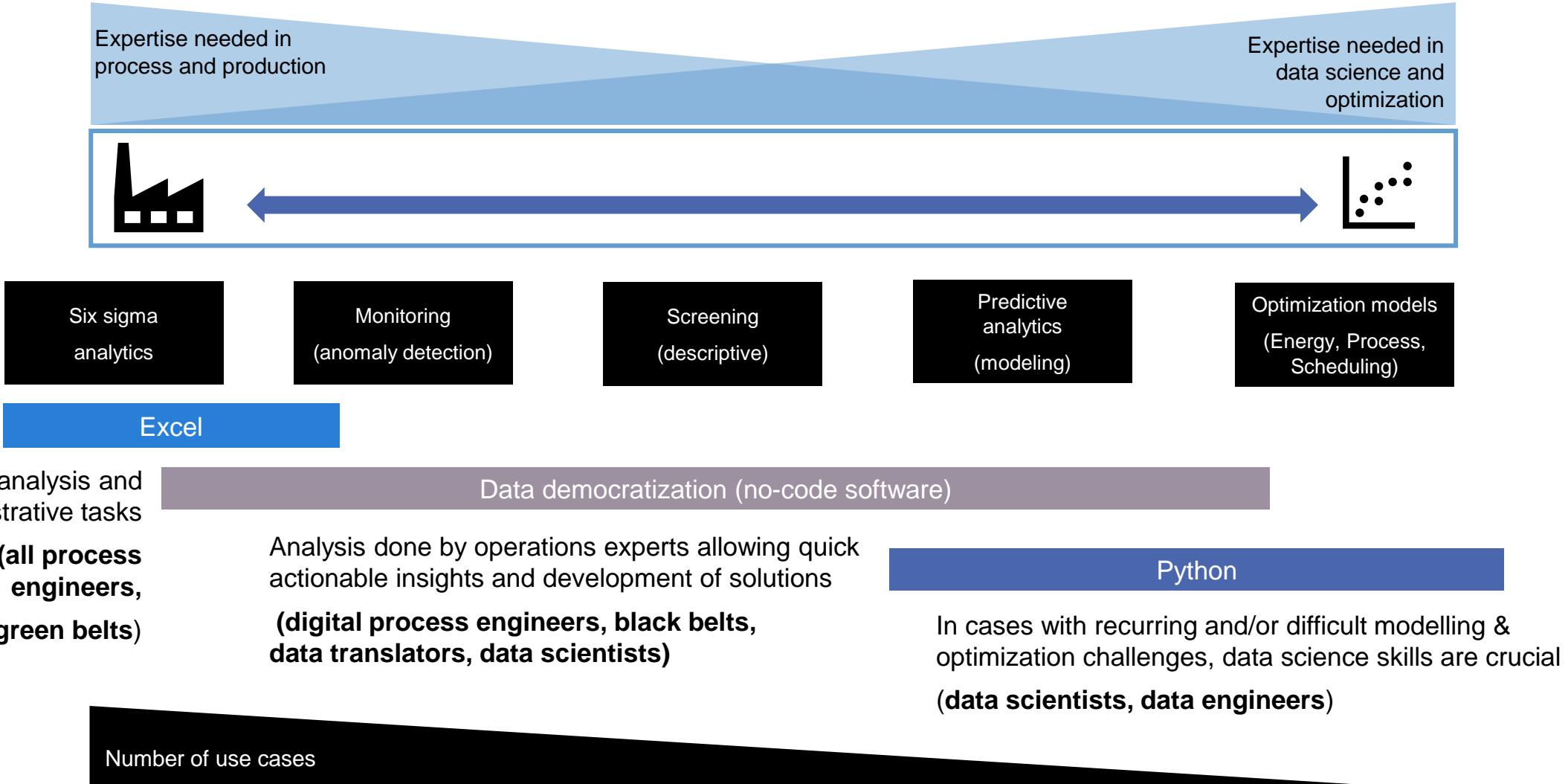


Historians (OT)
vs
(IT) open-source
databases



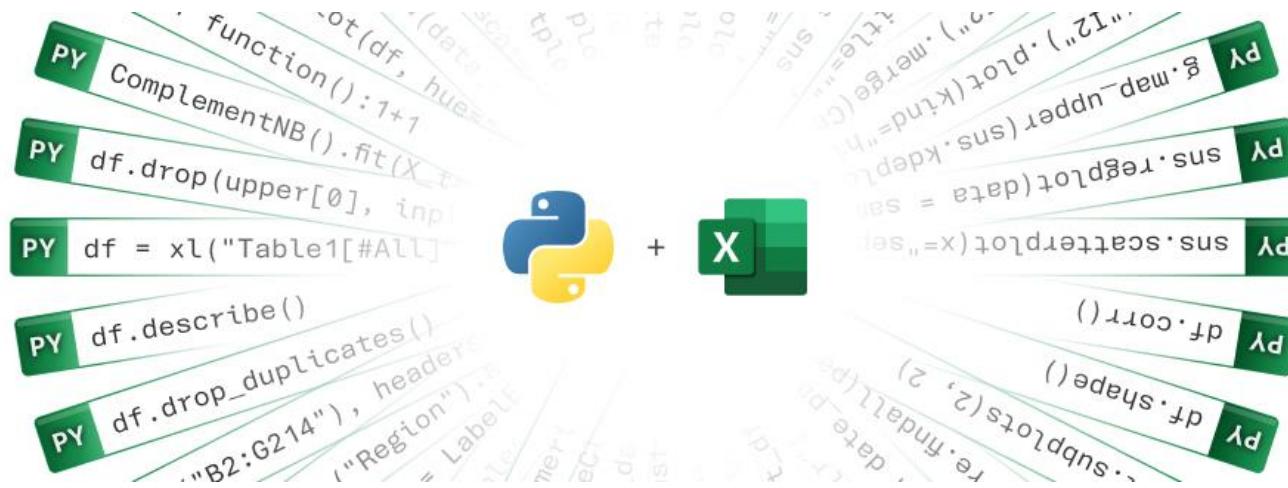
...

Data democratization tries to close the gap for industrial experts while covering areas with a majority of use cases

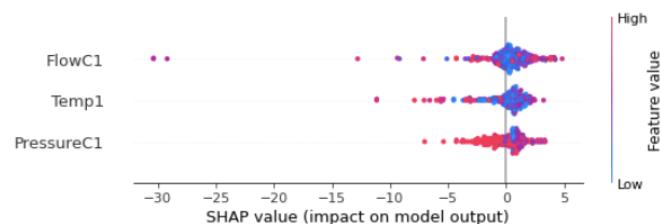


Python in Excel

- The most popular data analysis tool integrated with the most popular programming language



<https://techcommunity.microsoft.com/t5/excel-blog/announcing-python-in-excel-combining-the-power-of-python-and-the/ba-p/3893439>



ML and Explainable AI integrated within Excel

Demo: screening processes using machine learning



Reaction
Chemistry &
Engineering

REVIEW

Check for updates

Cite this: *React. Chem. Eng.*, 2022, 7, 1471

Received 1st December 2021,
Accepted 21st February 2022

DOI: 10.1039/d1re00541c

rsc.li/reaction-engineering

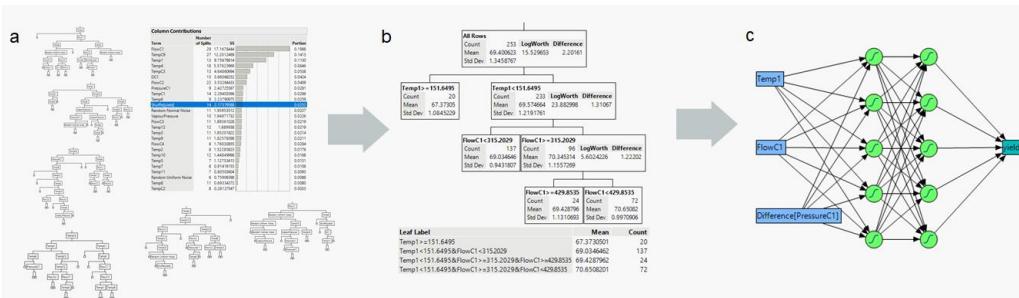


View Article Online
View Journal | View Issue

Industrial data science – a review of machine learning applications for chemical and process industries†

Max Mowbray, ^a Mattia Vallerio, ^b Carlos Perez-Galvan, ^{ab} Dongda Zhang, ^{ac} Antonio Del Rio Chanona ^{bc} and Francisco J. Navarro-Bruil ^{bc}

In the literature, machine learning (ML) and artificial intelligence (AI) applications tend to start with examples that are irrelevant to process engineers (e.g. classification of images between cats and dogs, house pricing, types of flowers, etc.). However, process engineering principles are also based on pseudo-empirical correlations and heuristics, which are a form of ML. In this work, industrial data science fundamentals will be explained and linked with commonly-known examples in process engineering, followed by a review of industrial applications using state-of-art ML techniques.



[Industrial data science – a review of machine learning applications for chemical and process industries](#)

CHEMISTRY WORLD



NEWS

RESEARCH

OPINION

FEATURES

CULTURE

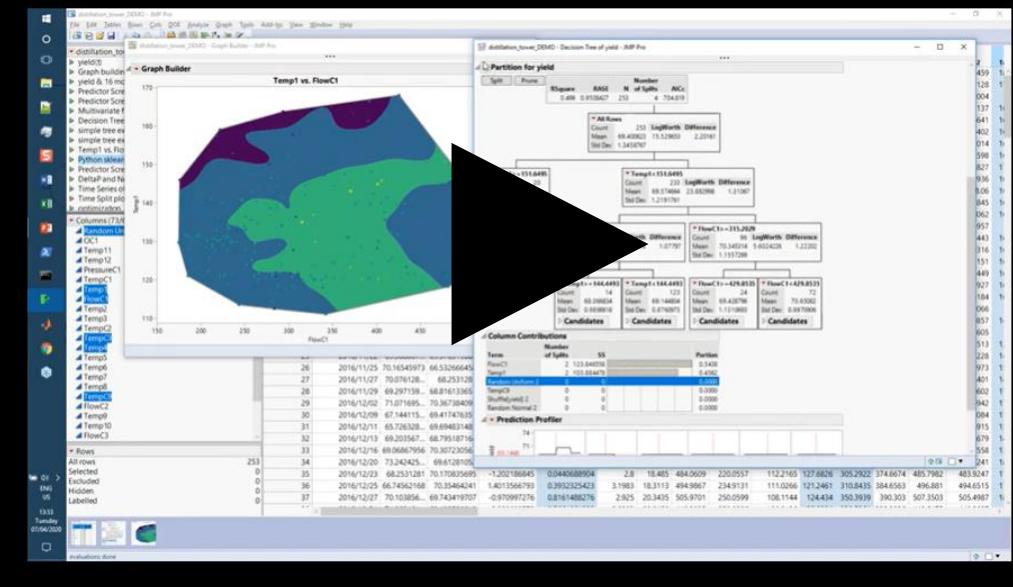
CAREERS

PODCASTS

WEBINARS

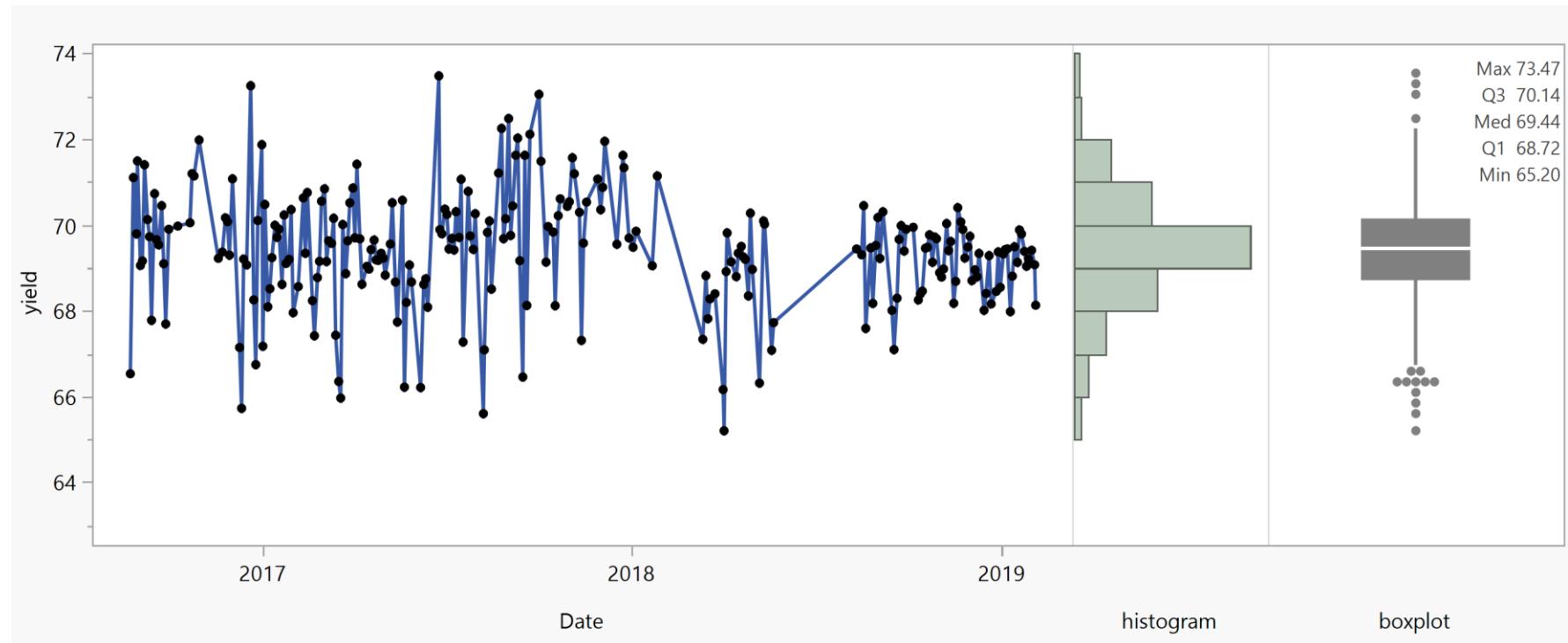
EXTRAS | Jobs | Reading room | Sign up to Re:action

WEBINARS



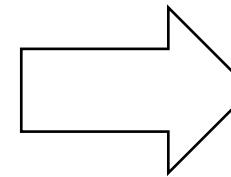
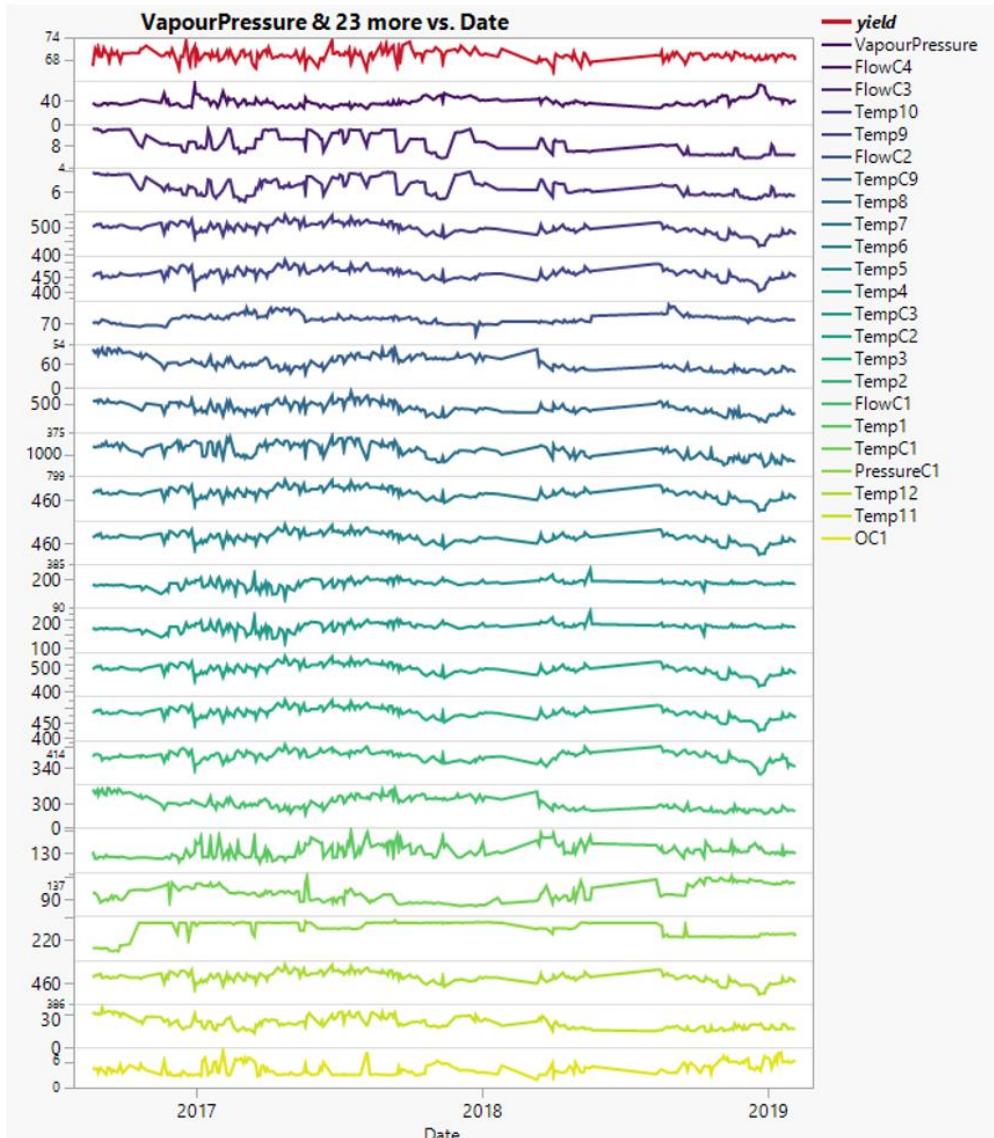
[Industrial data science for chemical process perfection | Webinar | Chemistry World](#)

Distillation tower demo



Process engineers from the site were able to stabilize the yield of a distillation tower but the baseline is now lower than before.

Distillation tower example

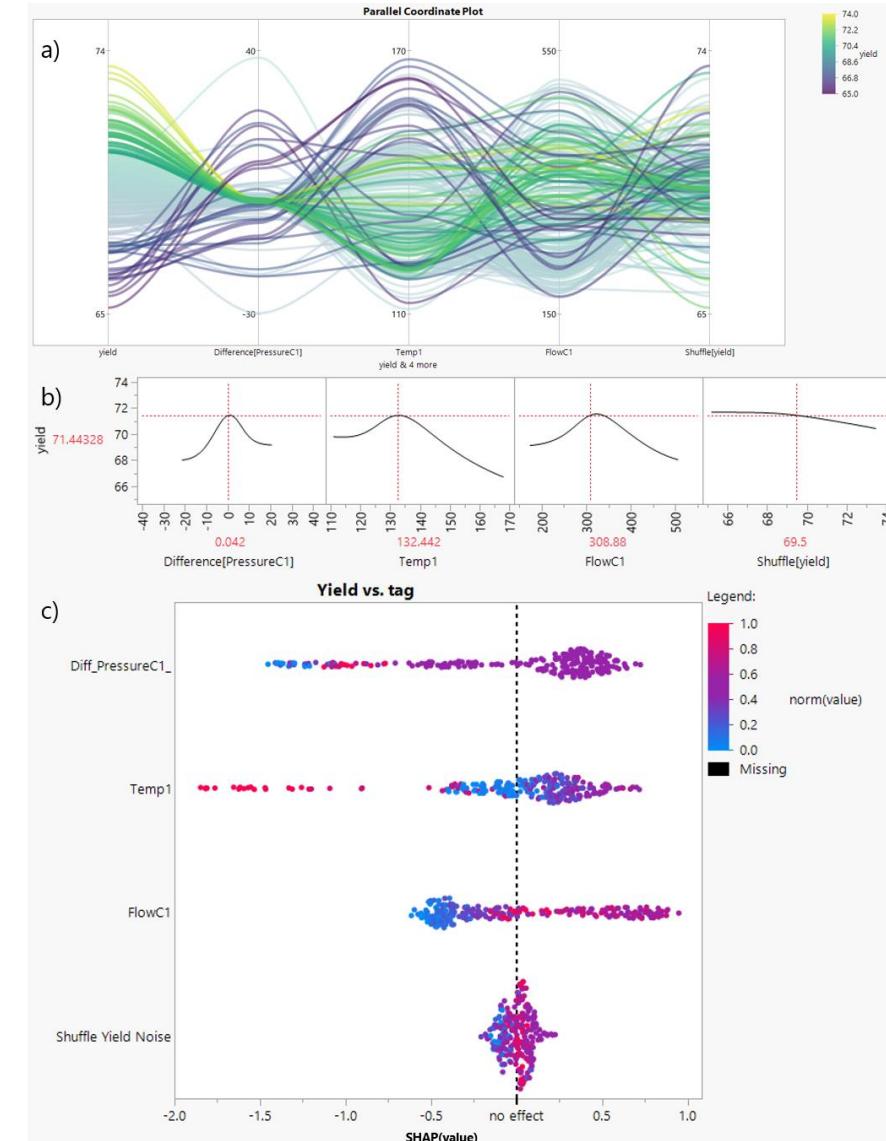


Automated screening of tags via correlation analysis (supervised learning)



Live demo

<https://www.chemistryworld.com/webinars/industrial-data-science-for-chemical-process-perfection/4011242.article>



Src: <https://pubs.rsc.org/en/content/articlelanding/2022/re/d1re00541c>

Industrial Data Science

Industrial data and process control



What happens with all that data?

- Data is generated by sensors, automation systems, events, incidents, people, transactions...
- Several dedicated systems exist to handle this data
- A plant engineer needs to understand which data and systems are relevant for process performance improvements
- Data is not just numbers (P&ID, DCS, PFD, SOP, OEE, instruction manuals...)



10-100 staff

1k-10k sensors and
actuators

process control
and automation



VGB Power Tech GmbH Germany

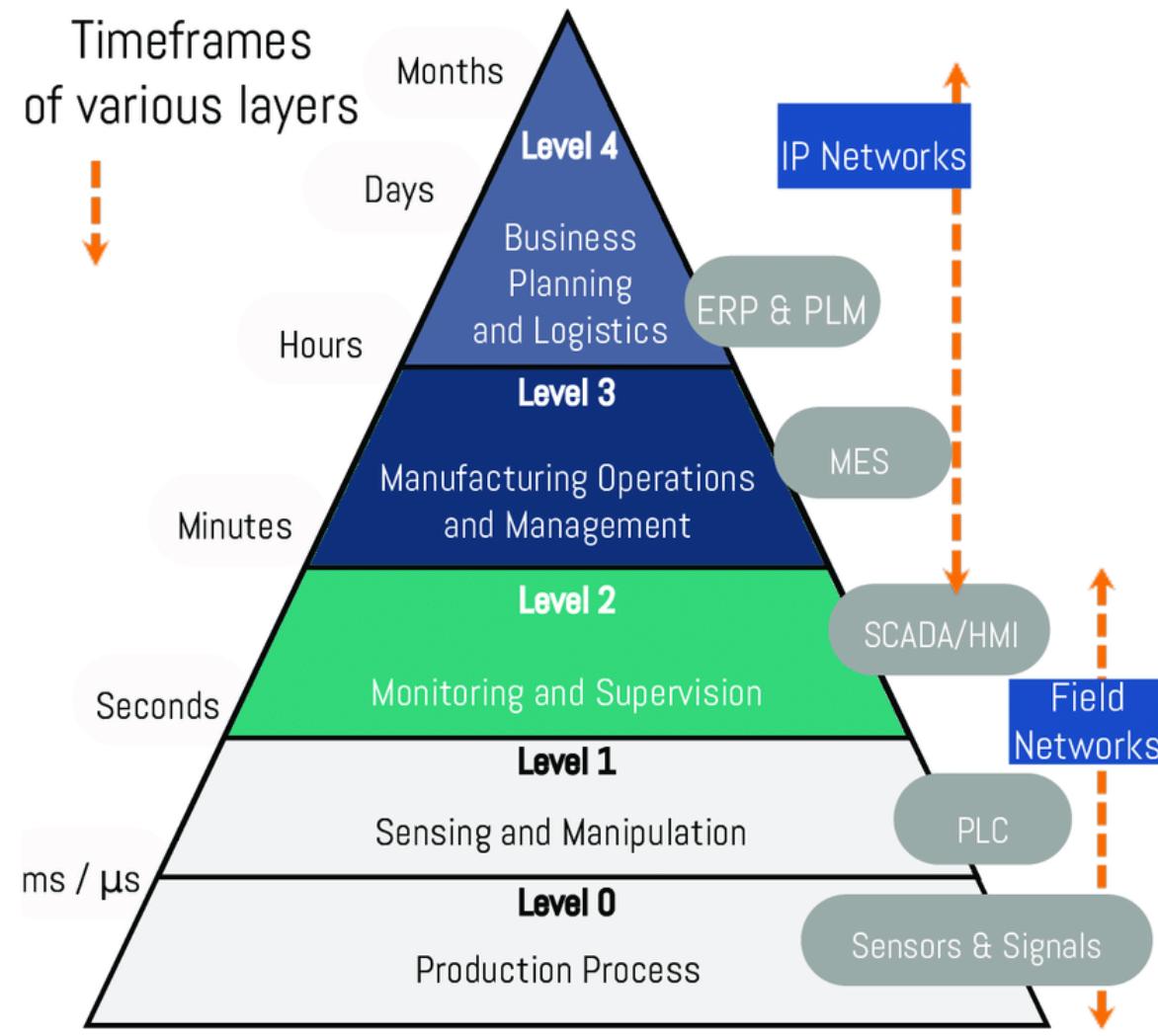
The Purdue model (Automation pyramid)



IT –
Informational
Technology



OT –
Operational
Technology

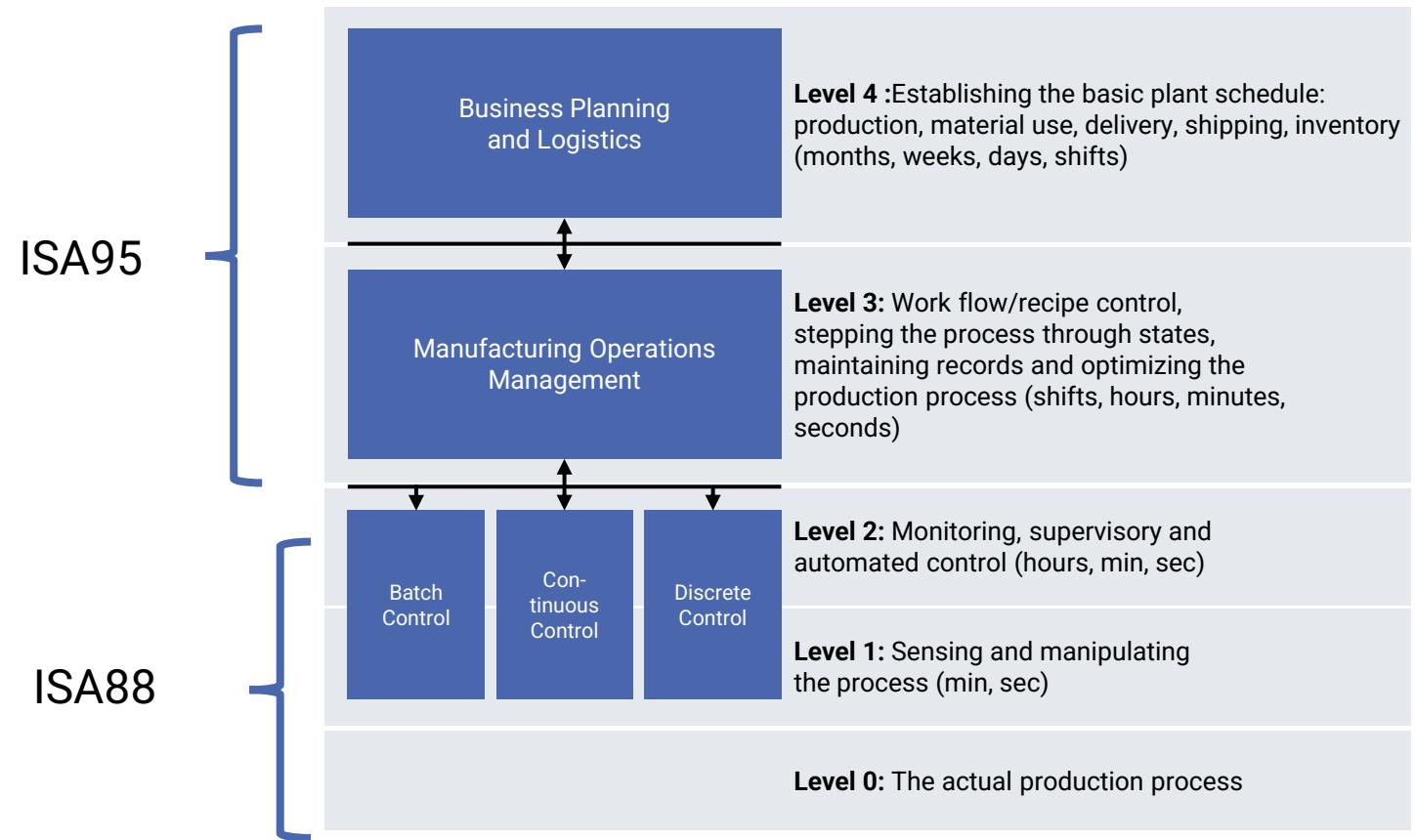


+info:

<https://learn.umh.app/lesson/introduction-into-it-ot-automation-pyramid/>

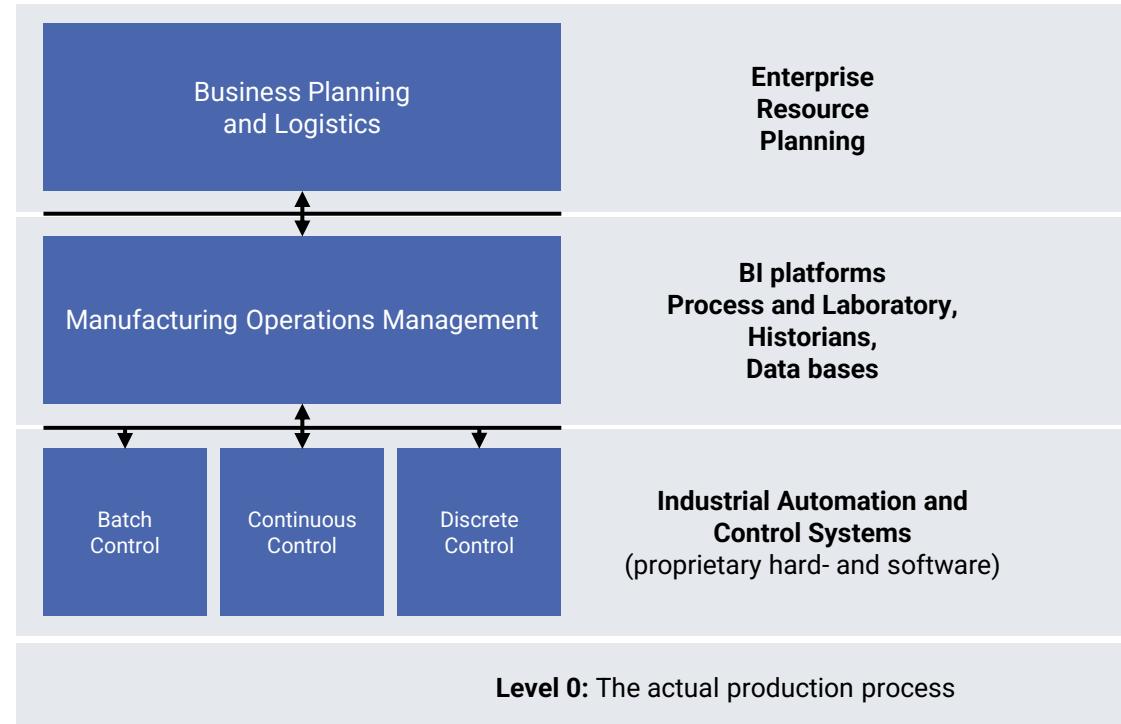
Generic structures of process data

- Data for process improvement often involves sensor and actuator data, their processed or aggregated forms, and information from transactions based on this data
- Industry standards describe the relationships and uses (ISA95/ ISA88/IEC62264)
- The higher the level, the longer the time horizon and the coarser the data



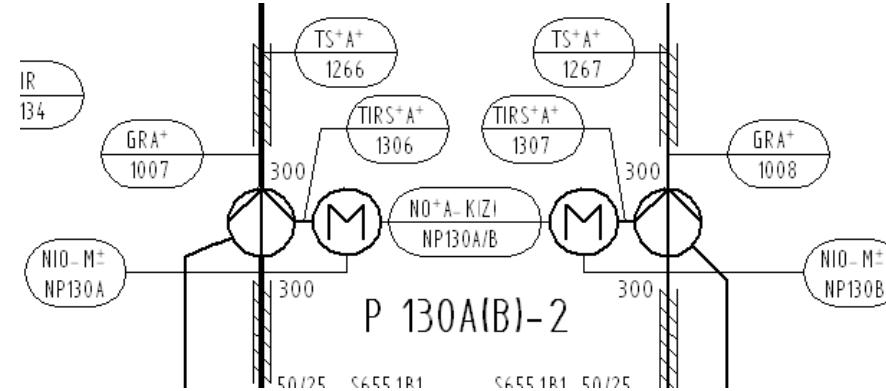
Data Systems Overview

- Depending on the function, proprietary or common IT solutions exist
- Particularly data in level 3 is relevant for Data Analytics, but data from other layers is used as well
- Data in level 3 is commonly stored in the form of times series data (trend data) with a sample frequency of ~minutes



Data identification level 1,2,3 – tag names

- Data in level 1,2 and 3 is often identified by a tag name: a combination of letters (that denote the function) number (that identify the particular sensor, actuator or control loop)

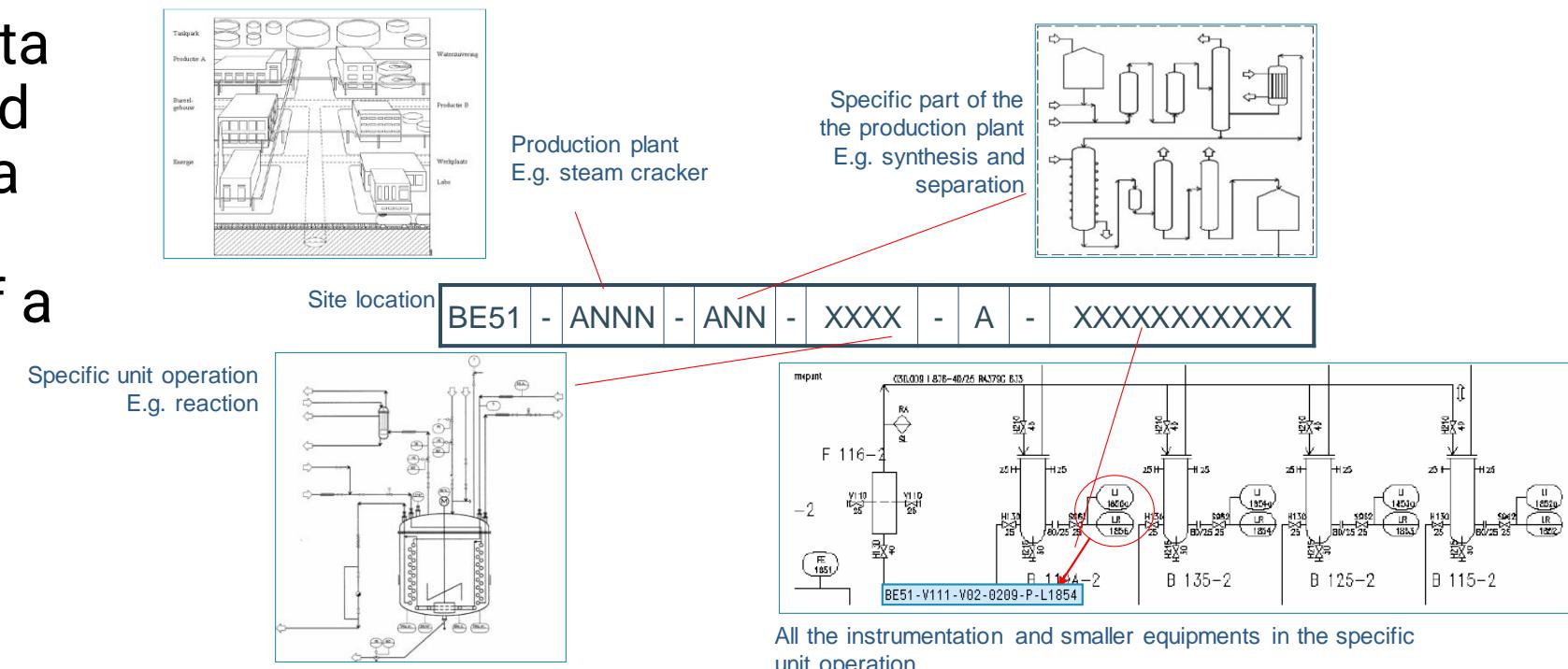


Letter	Measurement		Action
	1e letter	2e letter	
A			Alarm
C			Control
D	Density	Difference	
E	Electrical		
F	Flow	Ratio	
G	Position		
H	Manual action		
I			Indication
K	Time		
L	Level		
M			
N	Electrical consumer		

Letter	Measurement		Action
	1e letter	2e letter	
O			
P	Pressure		
Q	Quality	Integral	
R			Registration
S	Speed, frequency		Switch
T	Temperature		
U	Calculated property		
V	Viscosity		
W	Weight		
+			High level
-			Low level

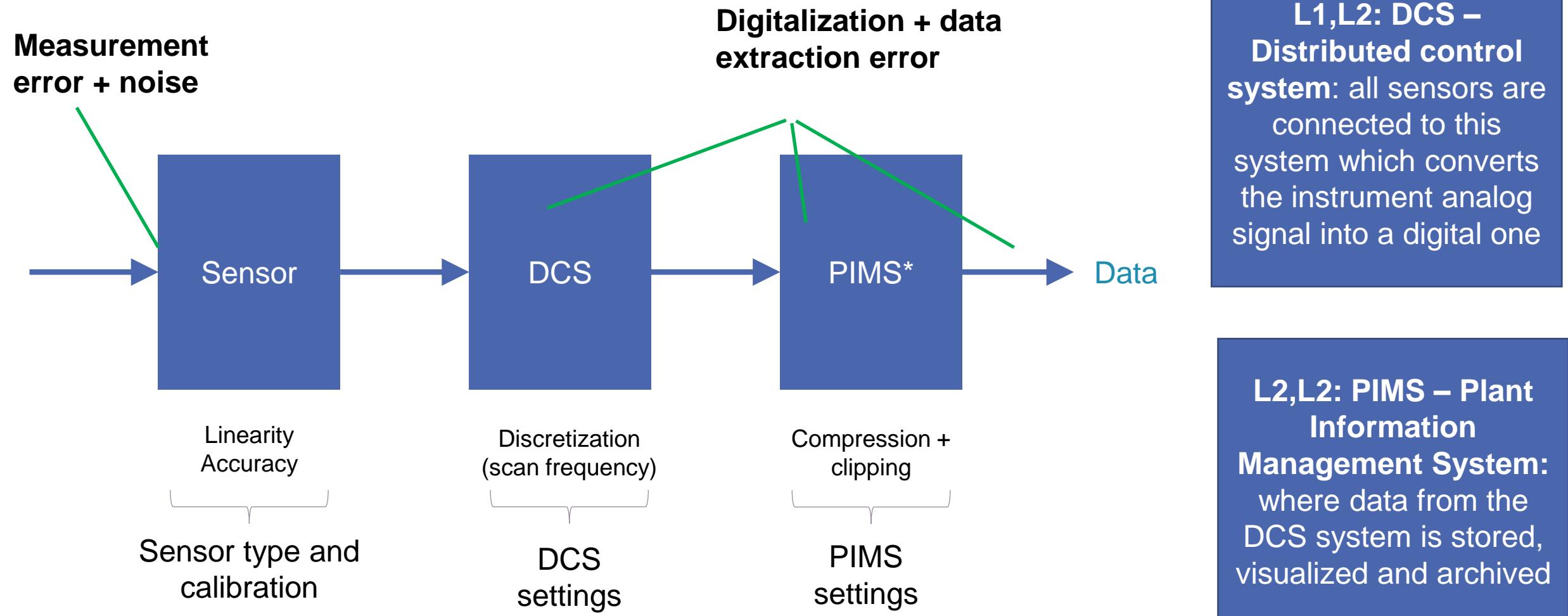
Data identification level 3,4 – functional location

- Data in level 4 (and to some extend in level 3) is identified by data base field names and functional location: a code that uniquely identifies all parts of a chemical plant



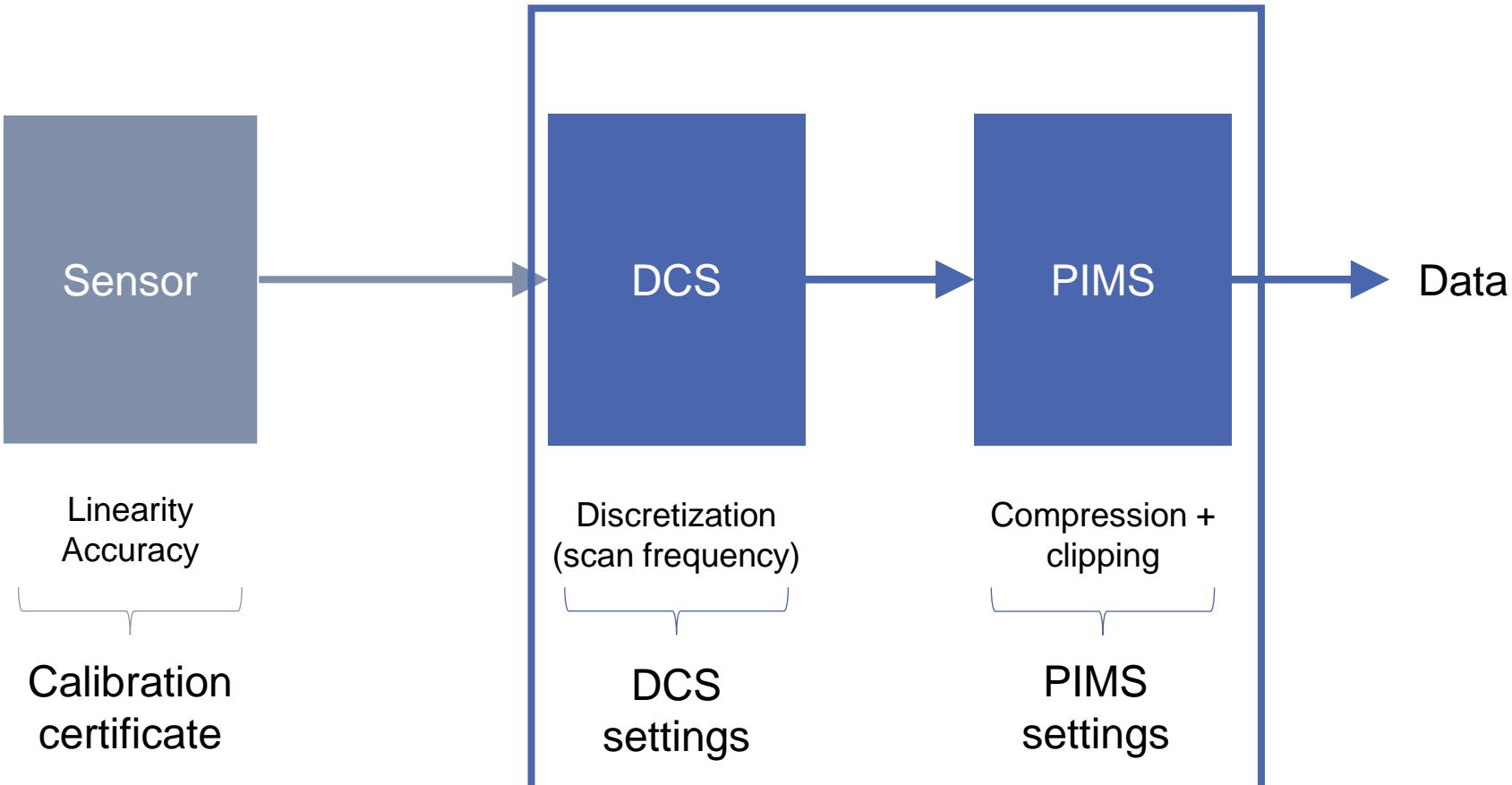
Data quality assessment for chemical production data

Where are the errors coming from?



Data quality assessment for PIMS data

Determine the digitalization error

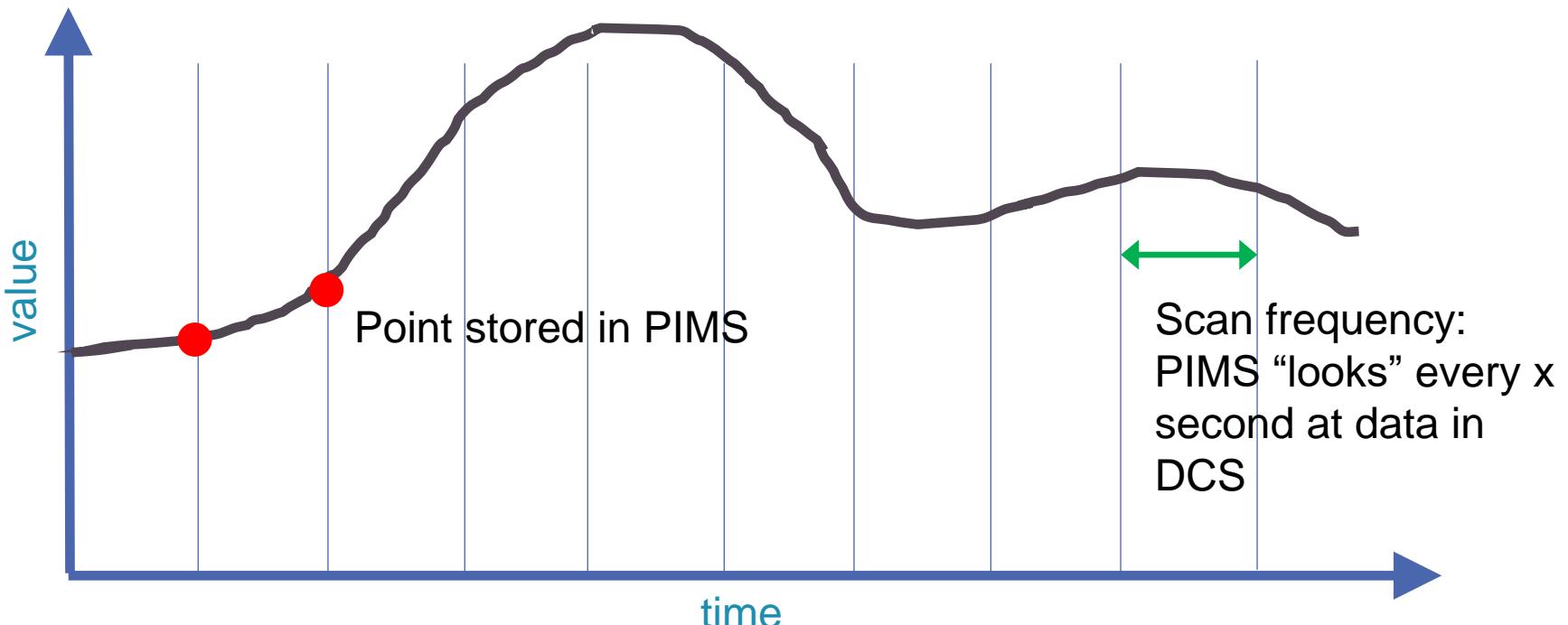


Data quality assessment for PIMS data

Digitalization error

Measurement error
Sampling error

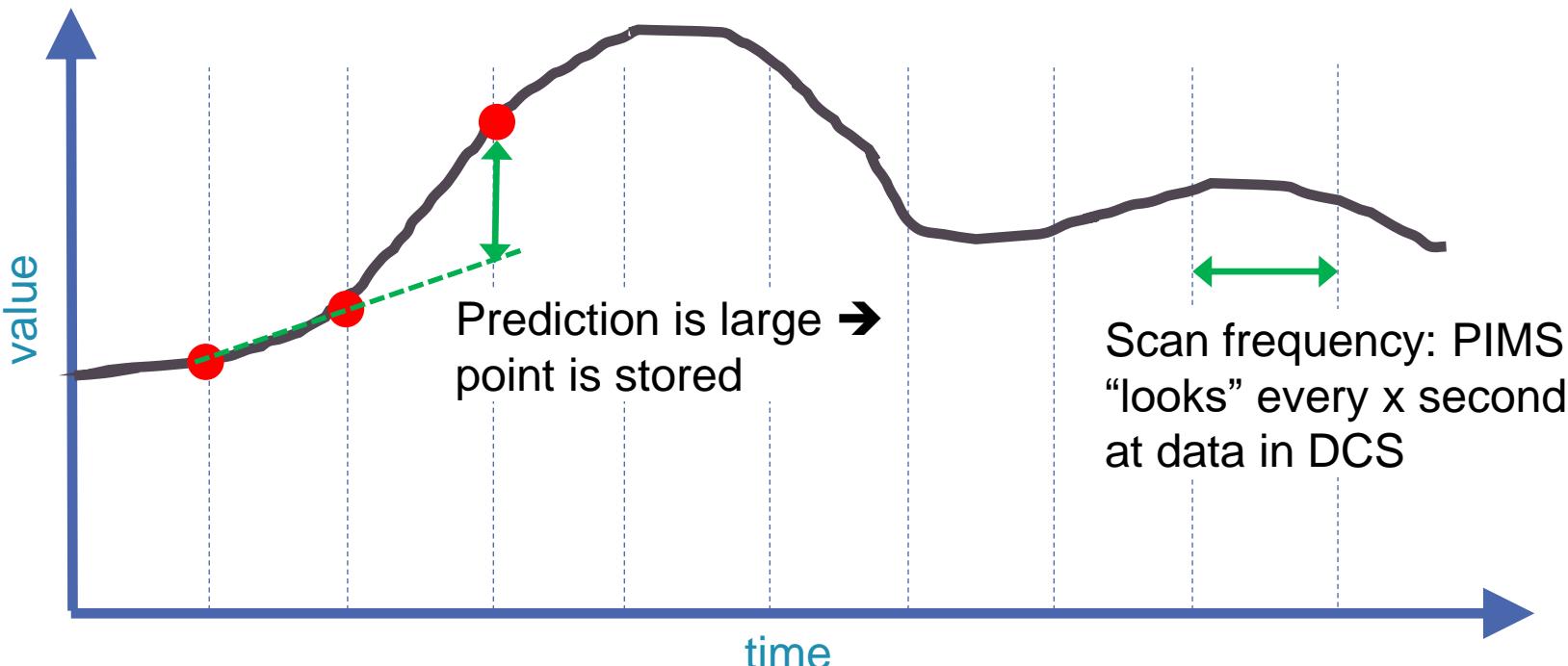
PIMS uses the scan frequency setting to determine
when to look at the DCS data



Data quality assessment for PIMS data

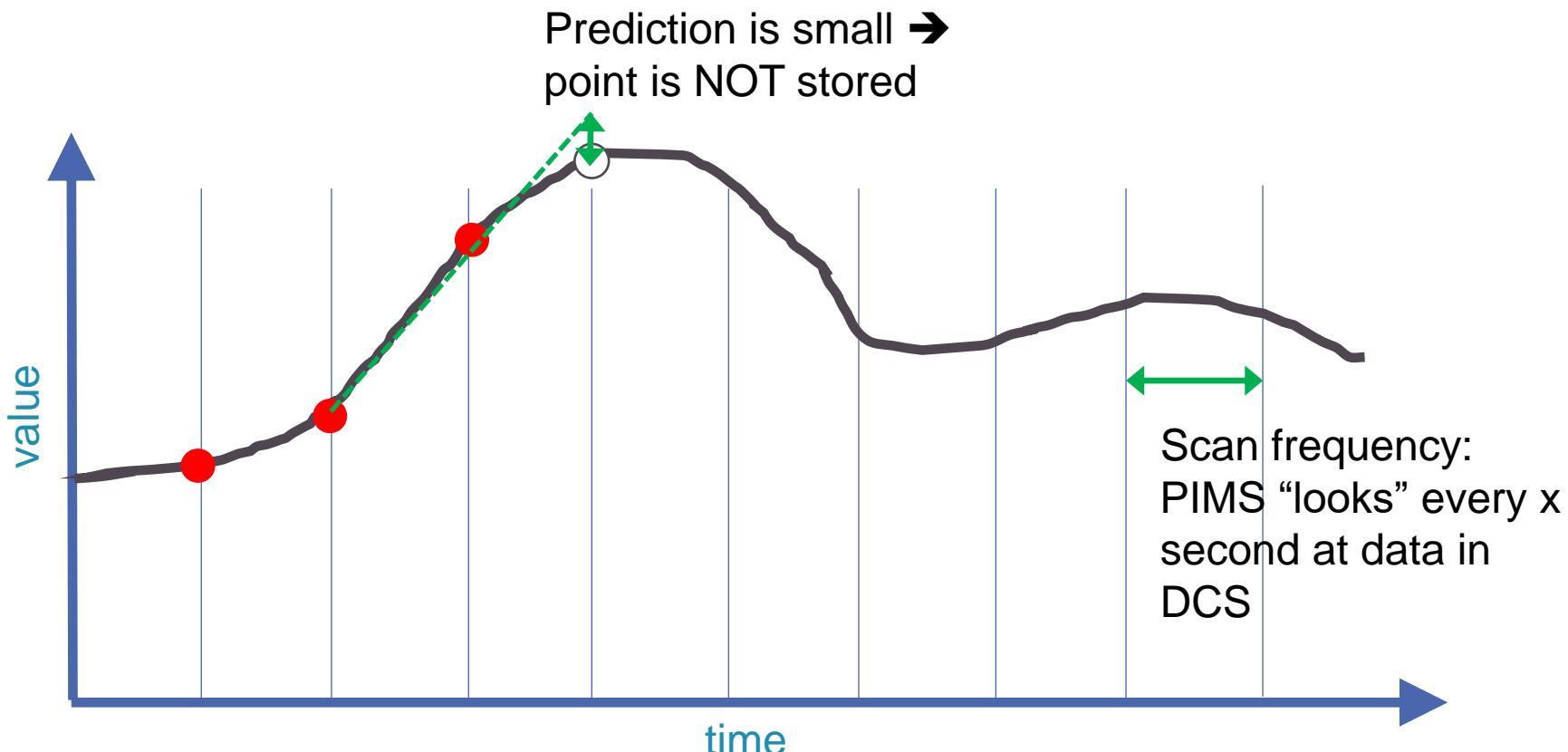
Digitalization error

For each new point, PIMS tries to predict the new point based on older points → if the prediction error is too large (depending on compression settings) the new point is stored



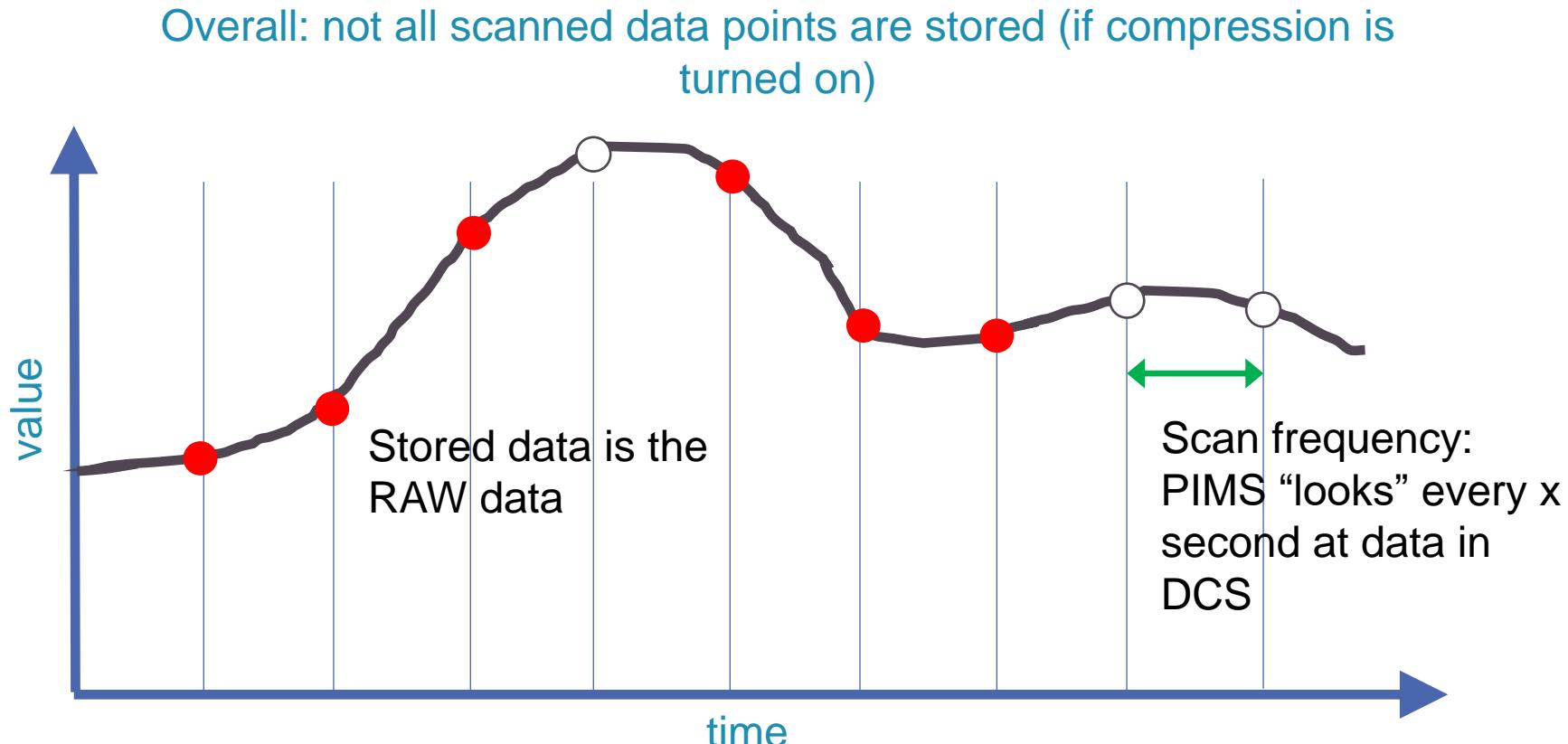
Data quality assessment for PIMS data

Digitalization error



Data quality assessment for PIMS data

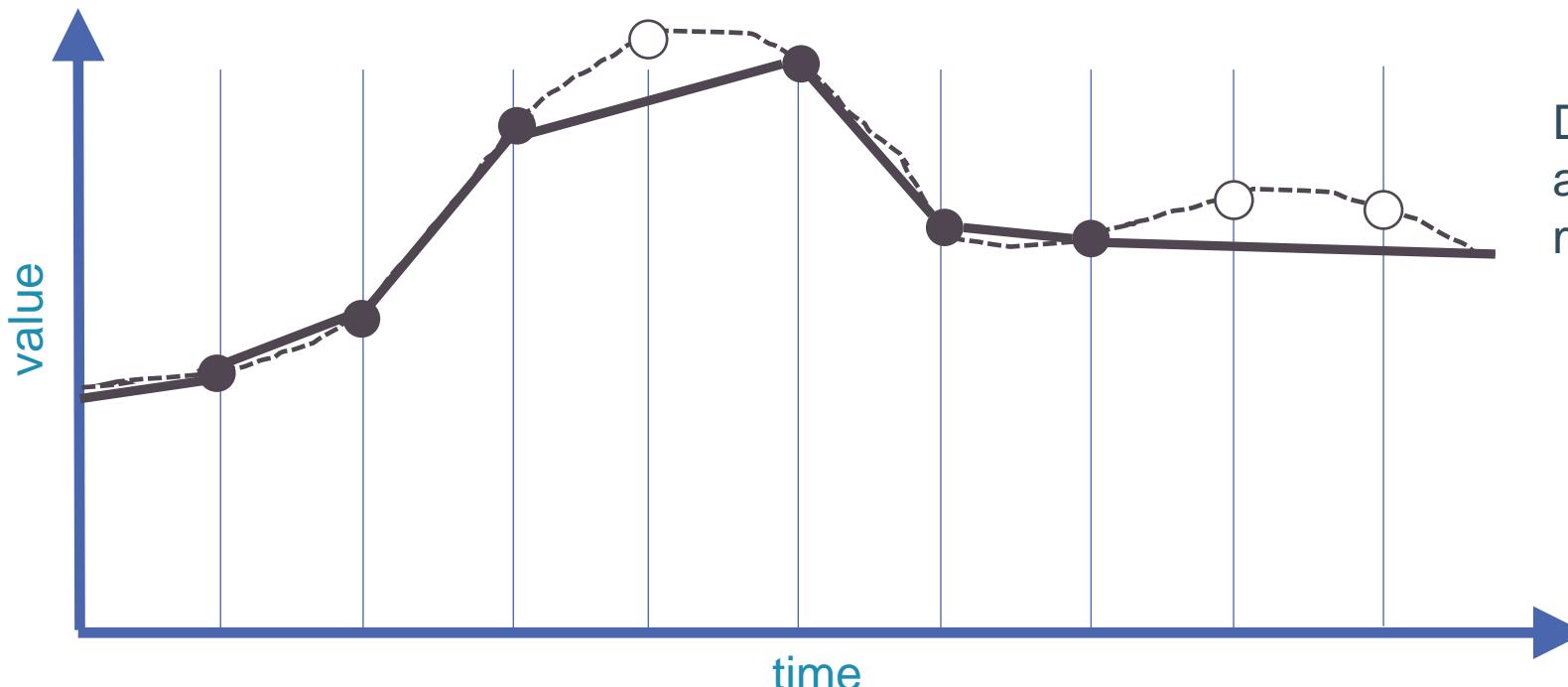
Digitalization error



Data quality assessment for PIMS data

Digitalization error

When retrieving the data (e.g. for data analysis) RAW data is interpolated to estimated missing data. Overall some detail is lost!



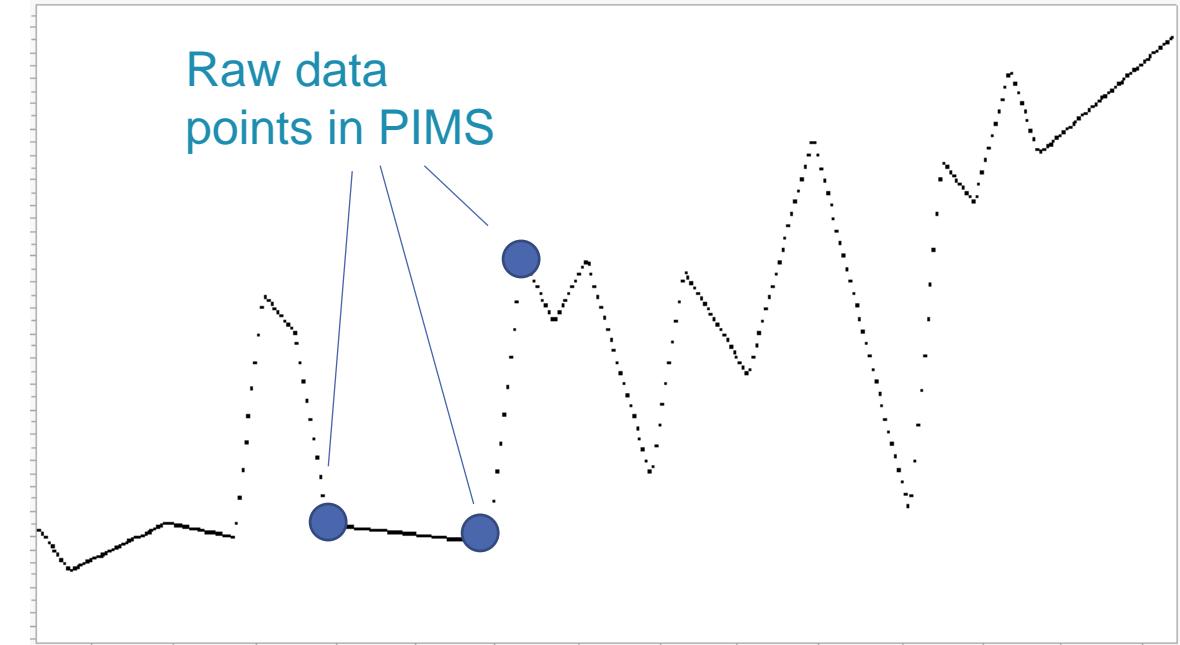
Different PIMS systems apply slightly different methods to store data.

Data quality assessment for PIMS data

Digitalization error – what if I ignore this?

Risks

- ▶ A lot of interpolated data that takes time to download from the system: data extraction takes much longer than needed.
- ▶ Variations that happen faster than what is captured by PIMS are not recorded and cannot be analyzed → the information is not in the data, but that does not mean it does not exist.



Industrial Data Science

Industrial metadata: asset hierarchy,
automation/batch sequences

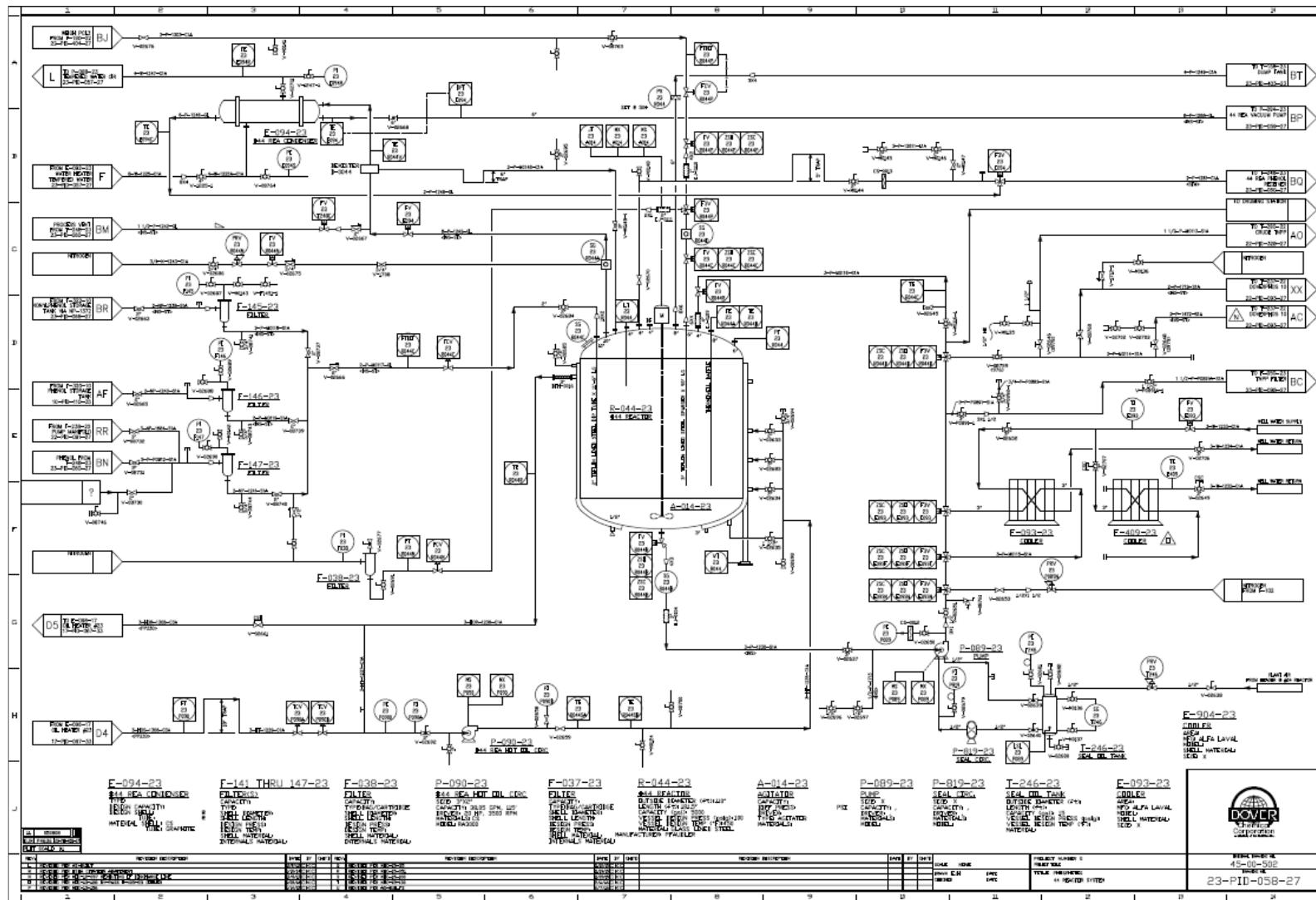


Diagrams

- Block Flow Diagram (BFD)
 - Unit operations (blocks) connected by lines (streams)
- Process Flow Diagram (PFD)
 - Unit operations (icons) connected by numbered lines (streams)
 - Temps, Pressures, compositions, flows specified (e.g. Aspen)
- Piping and Instrumentation Diagram (P&ID)
 - Units, pipes, valves, pumps, exchangers, instruments, pipe specs
 - Control strategies (loops, interlocks, alarms, etc.)
- Isometric Diagram – scale drawing of process in 3D space
- Equipment Drawing – scale drawing of equipment + details

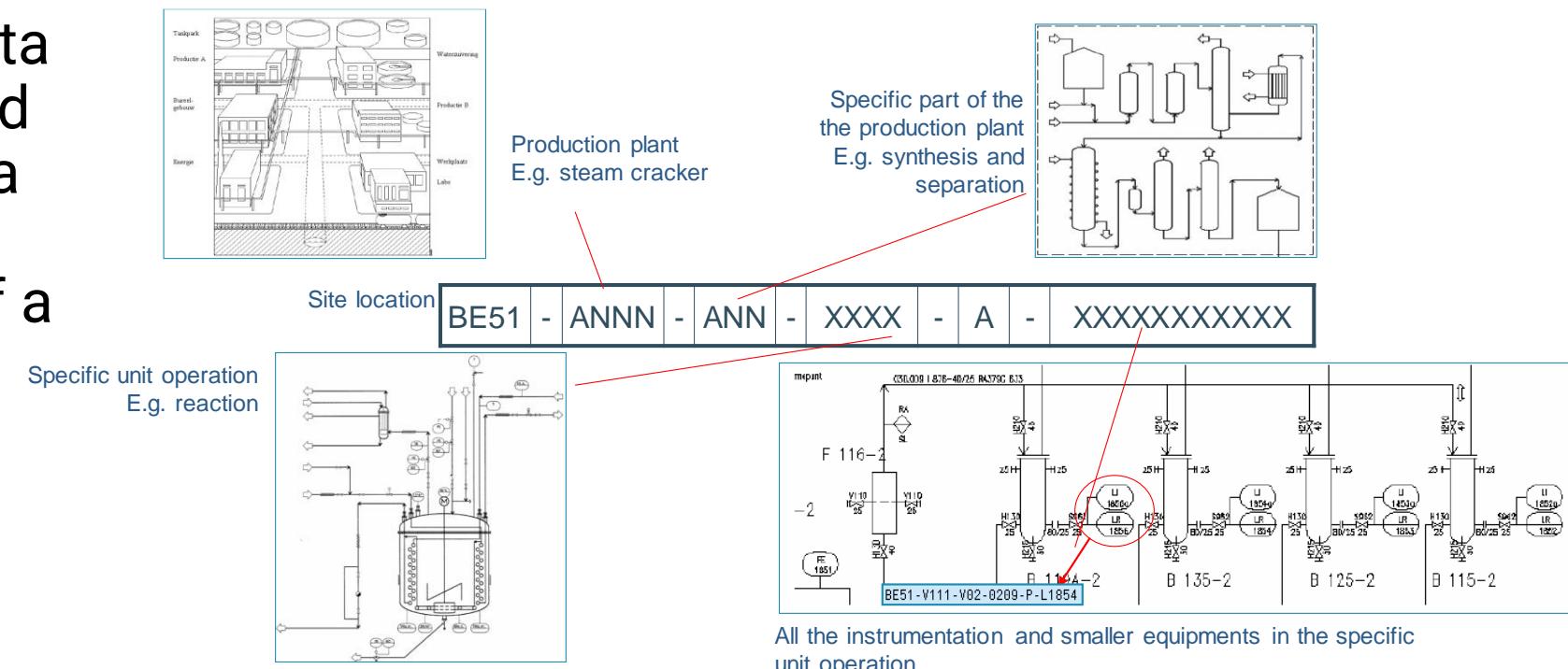
Source:
Industrial Data Science for
Batch Manufacturing
Processes
[arXiv:2209.09660](https://arxiv.org/abs/2209.09660)

Actual Batch Reactor P&ID



Data identification level 3,4 – functional location

- Data in level 4 (and to some extend in level 3) is identified by data base field names and functional location: a code that uniquely identifies all parts of a chemical plant



Asset hierarchy

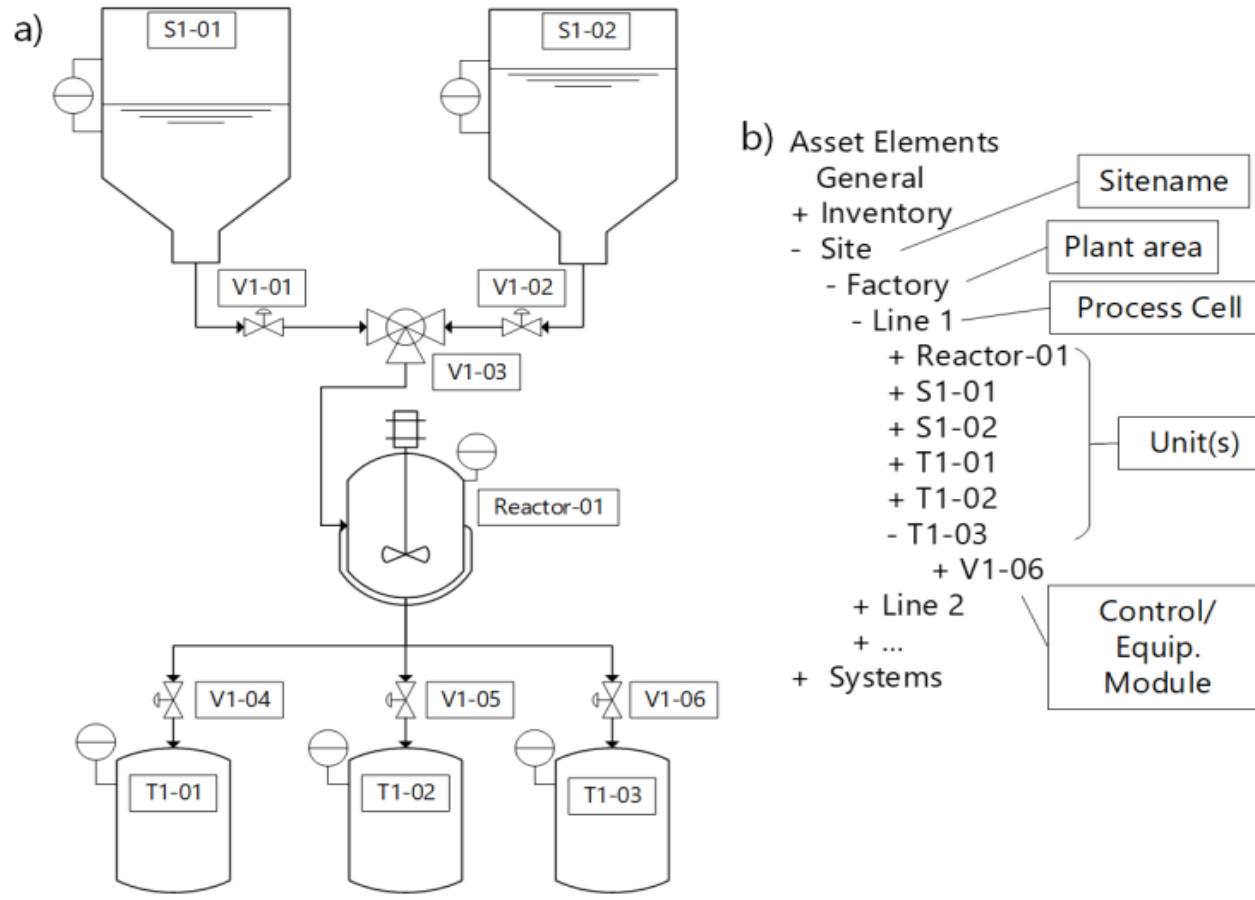


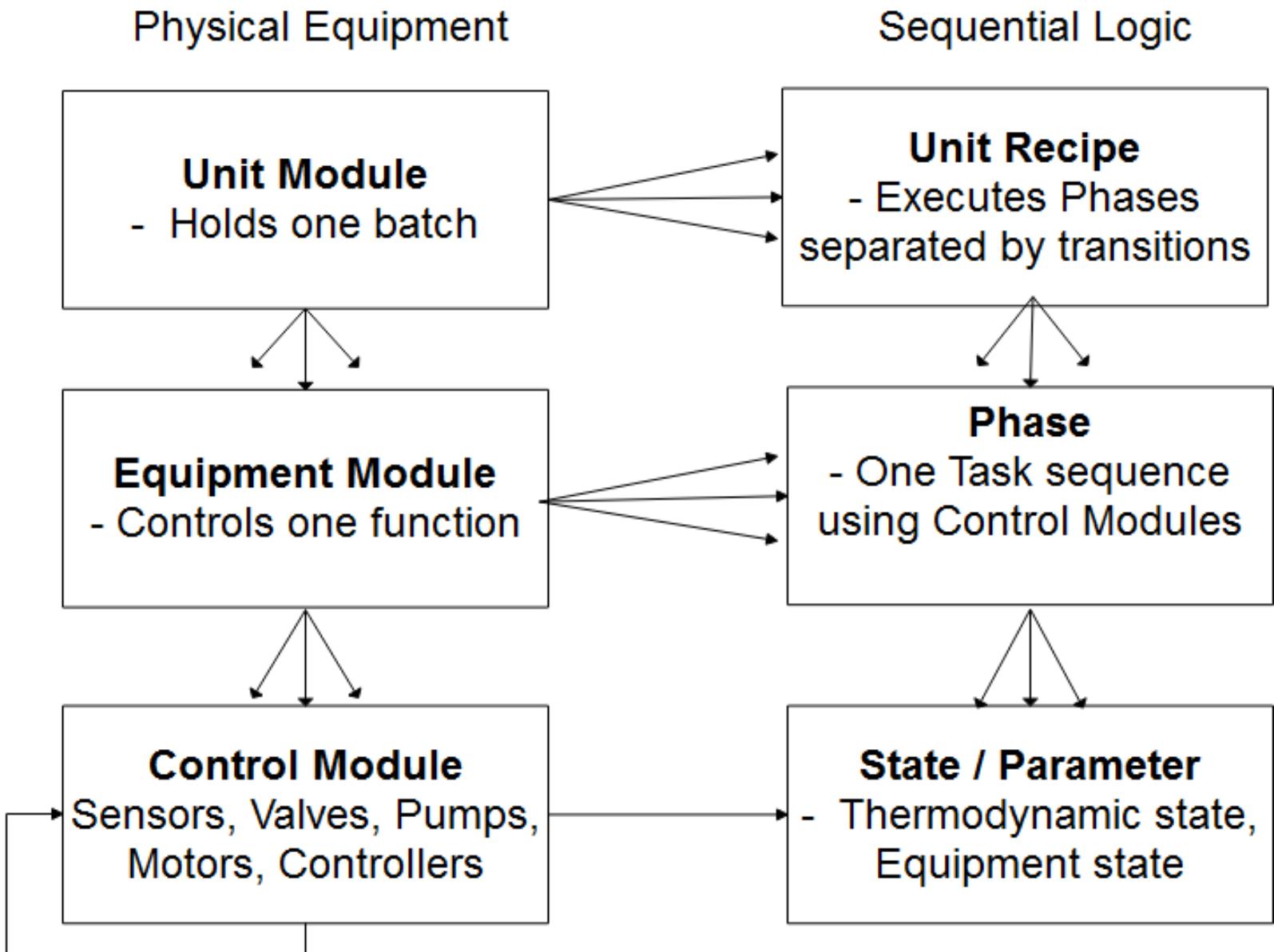
Figure 13: Simplified process (a) following ISA-88 to standardized asset hierarchy (b) in the plant.

<https://try.seeq.dev/workbooks?t=TVlfRk9MREVS>

Source:
Industrial Data Science for
Batch Manufacturing
Processes
[arXiv:2209.09660](https://arxiv.org/abs/2209.09660)

S88 Model

- We will focus on bottom 3 layers
- Model also includes Process, Area, Plant, and Business layers
- Unit Recipes and Phases are Sequences: SFCs



S88 Layers – Control Modules / States

- Control Modules are the lowest level of the physical model
 - Most detailed – As shown on a P&ID, e.g.:

Control Module	Command	Feedback
Solenoid Valve	Open/Close	Open/Closed/In Transition
RTD	(alarm levels)	Process Temperature, alarm status
Variable Speed Pump	On/Off, Target Speed	Running/Stopped, Actual Speed
Level Control Loop	Auto/Manual, Target Level (setpoint)	Measured level, control valve position, target deviation alarm

- Control Modules measure the **Thermodynamic State**
- Control Modules manipulate the **Equipment State**
 - We also measure the Equipment State (e.g. motor/valve status)

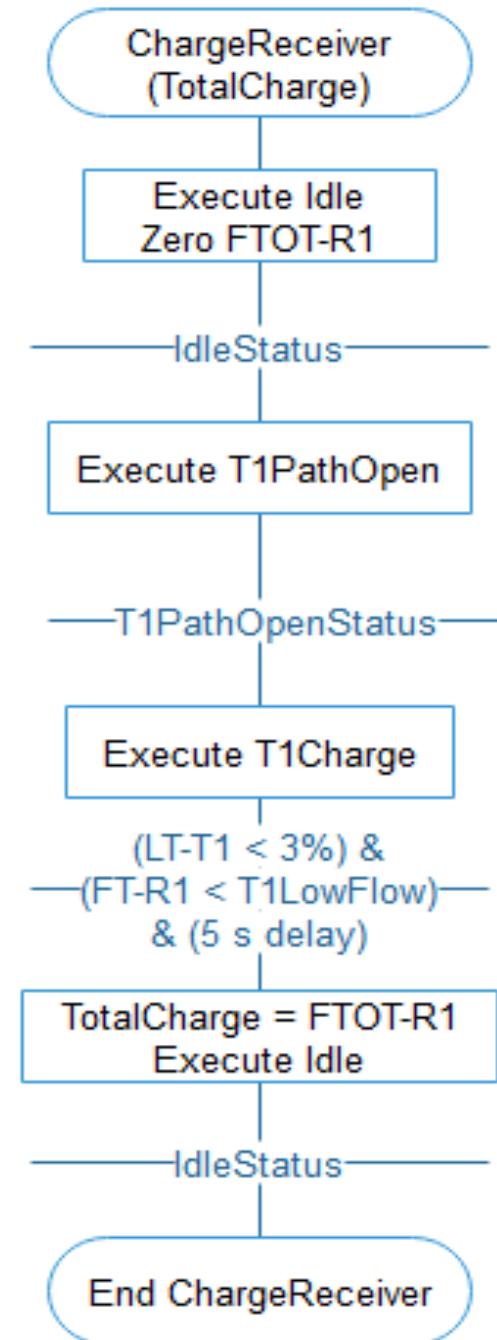
S88 Layers – Equipment Modules / Phases

- Equipment Modules are a collection of Control Modules associated with one Function, e.g.
 - **Temperature Control:** Temperature sensors, control valves and block valves for steam/hot oil/cooling water, Temp. controllers
 - **Batch Charging:** Flow meters, pumps, valves, level meters
 - **Discharge/Recirculate:** Pumps and valves for directing flow out of and back to the reactor and storage
- Equipment Modules have one or more **Phases**, Sequences that manipulate control modules to carry out one task
 - **Phases** are the vocabulary of the **Unit Recipes**
 - **Control Module States** are the vocabulary of the **Phases**

Phases

- Phases manipulate Control Modules to Complete a Task

Device	Equipment States		
	Idle	T1PathOpen	T1Charge
SV-T1	Closed	Open	Open
SV-P1	Closed	Open	Open
P-1	Stopped	Stopped	Running
SV-T2	Closed	Closed	Closed
SV-P2	Closed	Closed	Closed
P-2	Stopped	Stopped	Stopped
FCV-R1	0%	100%	100%



Unit Modules / Unit Recipes

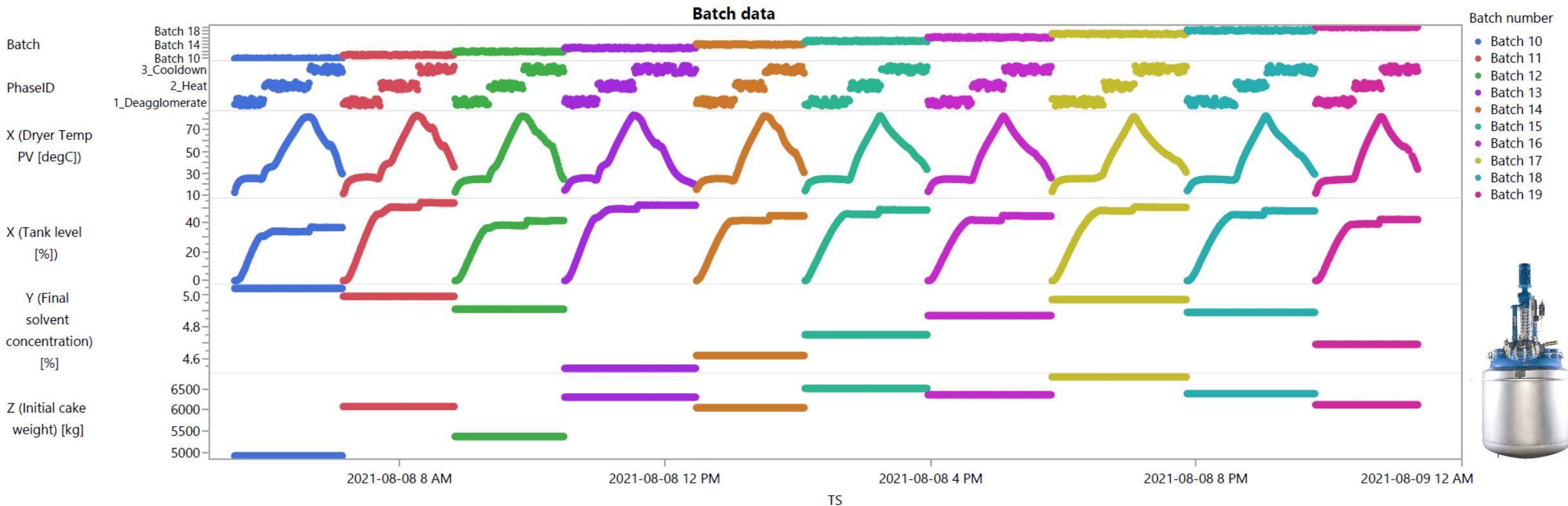
- A Unit Modules holds one batch
 - It is made up of many Equipment Modules
- It can execute one to many Unit Recipes
 - Unit Recipes are written with Phases as their vocabulary
 - A batch reactor would have one recipe for each product
 - A mixer or filter unit might only have one recipe
 - A unit recipe can also be a start up/shutdown sequence of a continuous plant
- Unit Recipes resemble Chemist Recipes
 - Good Communication tool for R&D / Production / Engineering
 - Operators can track batches with dynamic displays of the SFC

Master Recipes – Generic Unit Recipes

- Unit recipes and phases are specific to **one unit**
- **Master recipes** apply at any scale from lab to pilot to various production vessels
 - Phase parameters are **scaled** (usually to mass) or **unscaled**
 - Similar to extrinsic vs. intrinsic thermodynamic variables.

Variable	Description	Scaled/Unscaled	Master Recipe	Unit Recipe (R-3)
Batch Size	Product mass	Scaled	1000	5500
A_Charge (kg)	Mass A Charged	Scaled	364	2000
B_Charge (kg)	Mass B Charged	Scaled	636	3500
React_Temp (°C)	Rx Temperature	Unscaled	180	180
React_Wait (min)	Wait time	Unscaled	180	180
Strip_Press (mmHg)	Target pressure for Strip	Unscaled	50	50
Strip_Temp (°C)	Target Strip tem[.	Unscaled	160	160

Industrial data and batch processes



<https://try.seeq.dev/workbooks?t=TVlfRk9MREVS>

Source:
Industrial Data Science for
Batch Manufacturing
Processes
arXiv:2209.09660

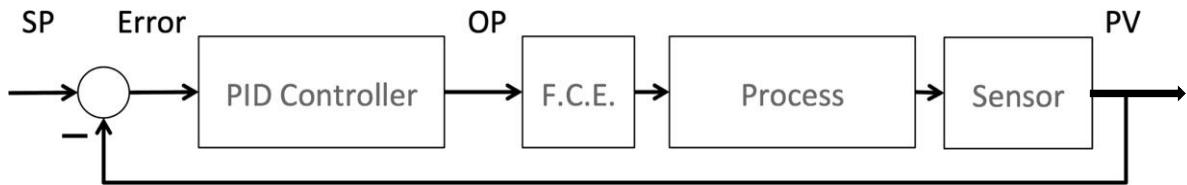
Industrial Data Science

Advanced Process Control



Control performance monitoring

– Industrial datasets and applications

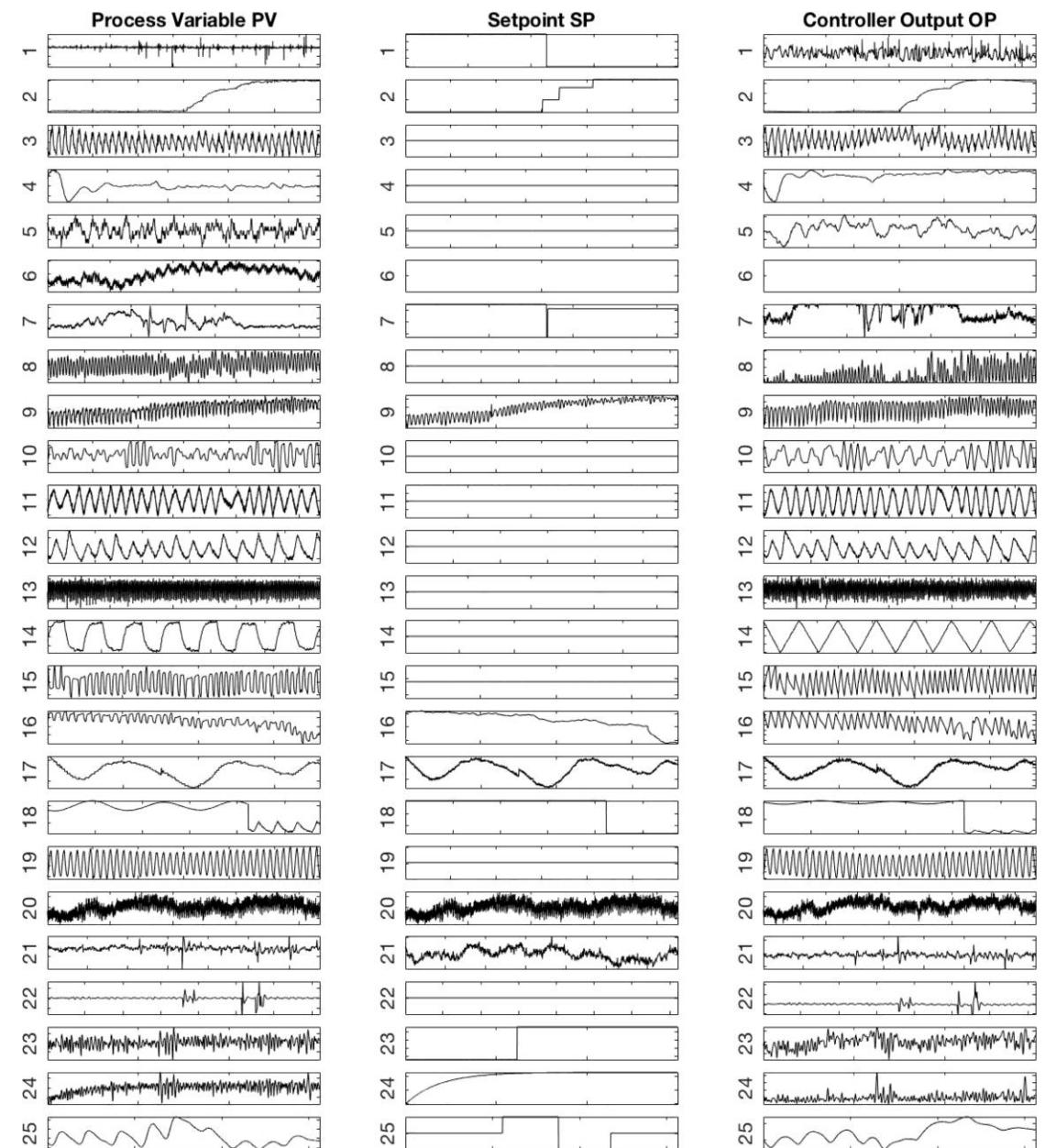


- incorrect tuning settings ([SysID add-on](#))
 - valve stiction ([add-on](#))
 - saturation ([add-on](#))
 - other actuator faults
 - sensor faults
 - other
 - unknown
- Source: [\[Open Access\] Industrial PID Control Loop Data Repository and Comparison of Fault Detection Methods | Industrial & Engineering Chemistry Research \(acs.org\)](#)

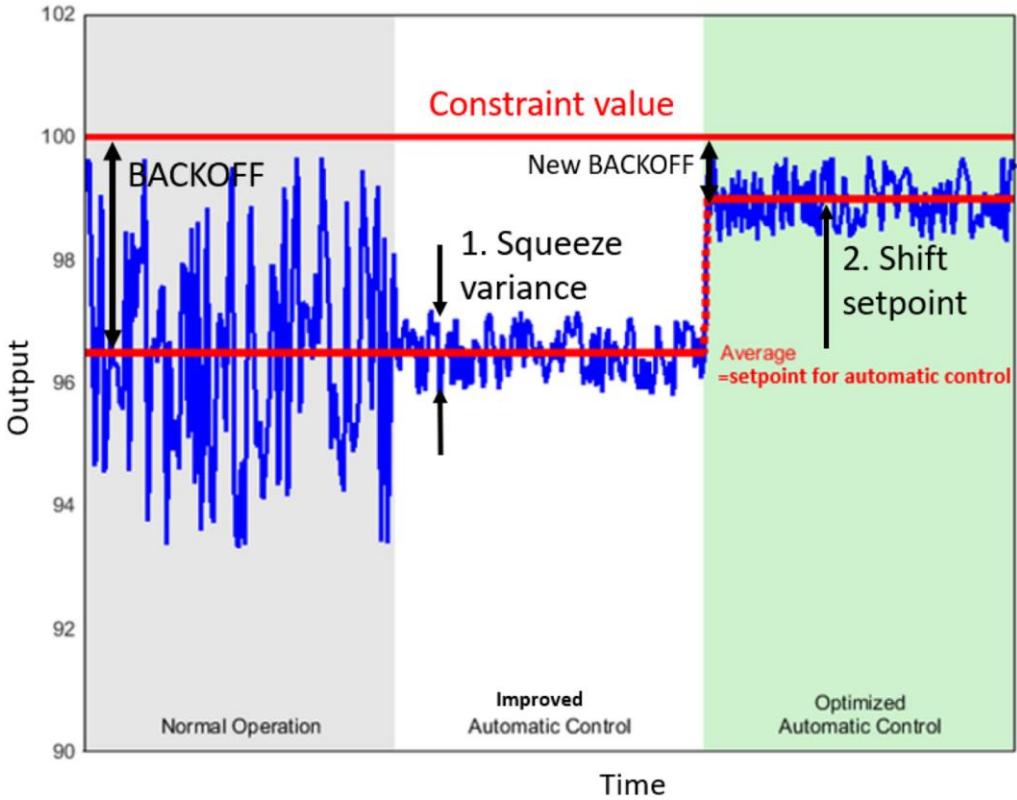
Other references

[Industrial datasets and a tool for SISO control loops data visualization and analysis – ScienceDirect](#)

[Detection and Diagnosis of Stiction in Control Loops \(ualberta.ca\)](#)

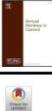
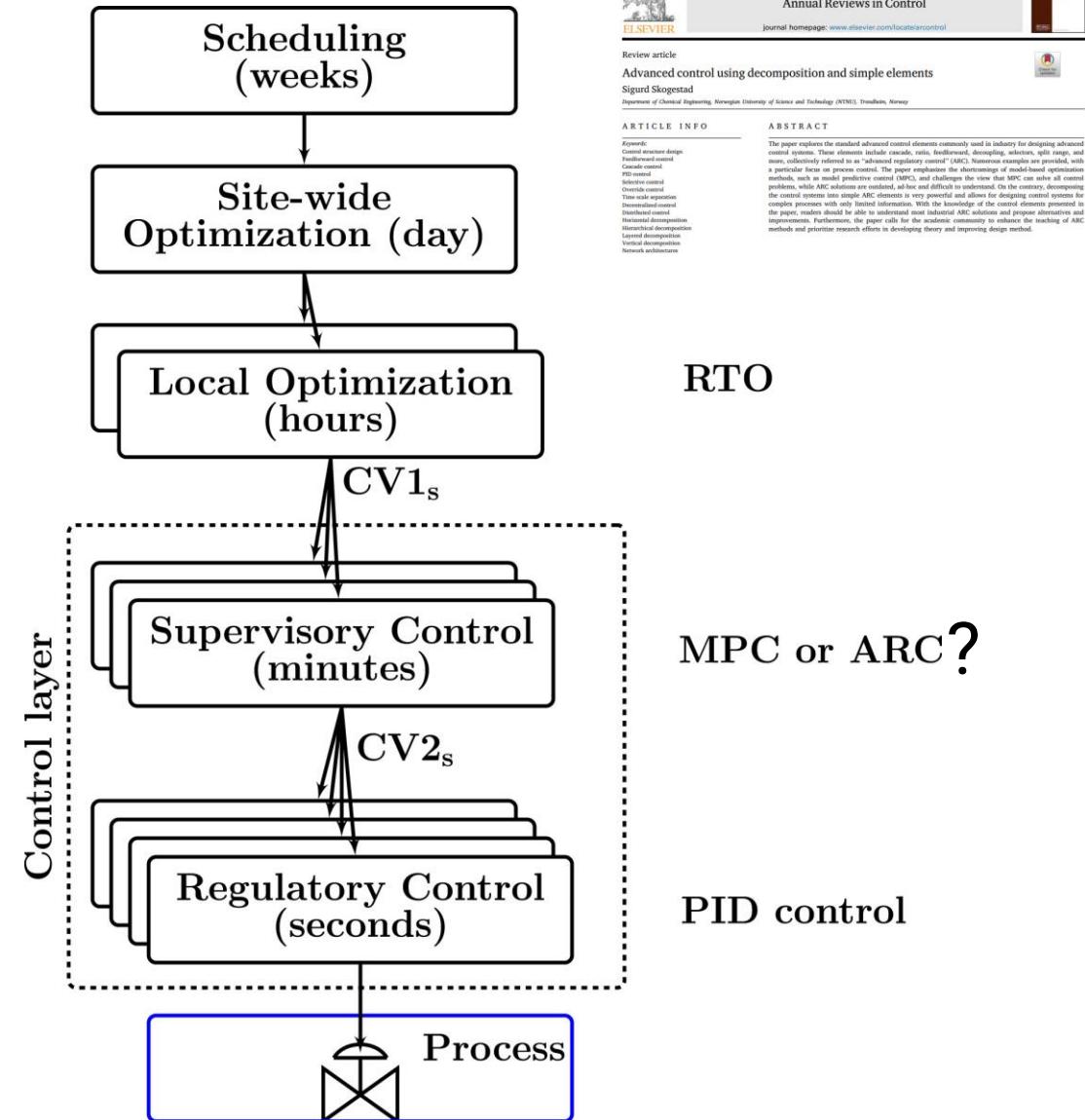


Advanced Process Control



Advanced control using decomposition and simple elements S. Skogestad

Annual Reviews in Control 56 (2023) 100903
<https://www.sciencedirect.com/science/article/pii/S1367578823000676>



RTO

MPC or ARC?

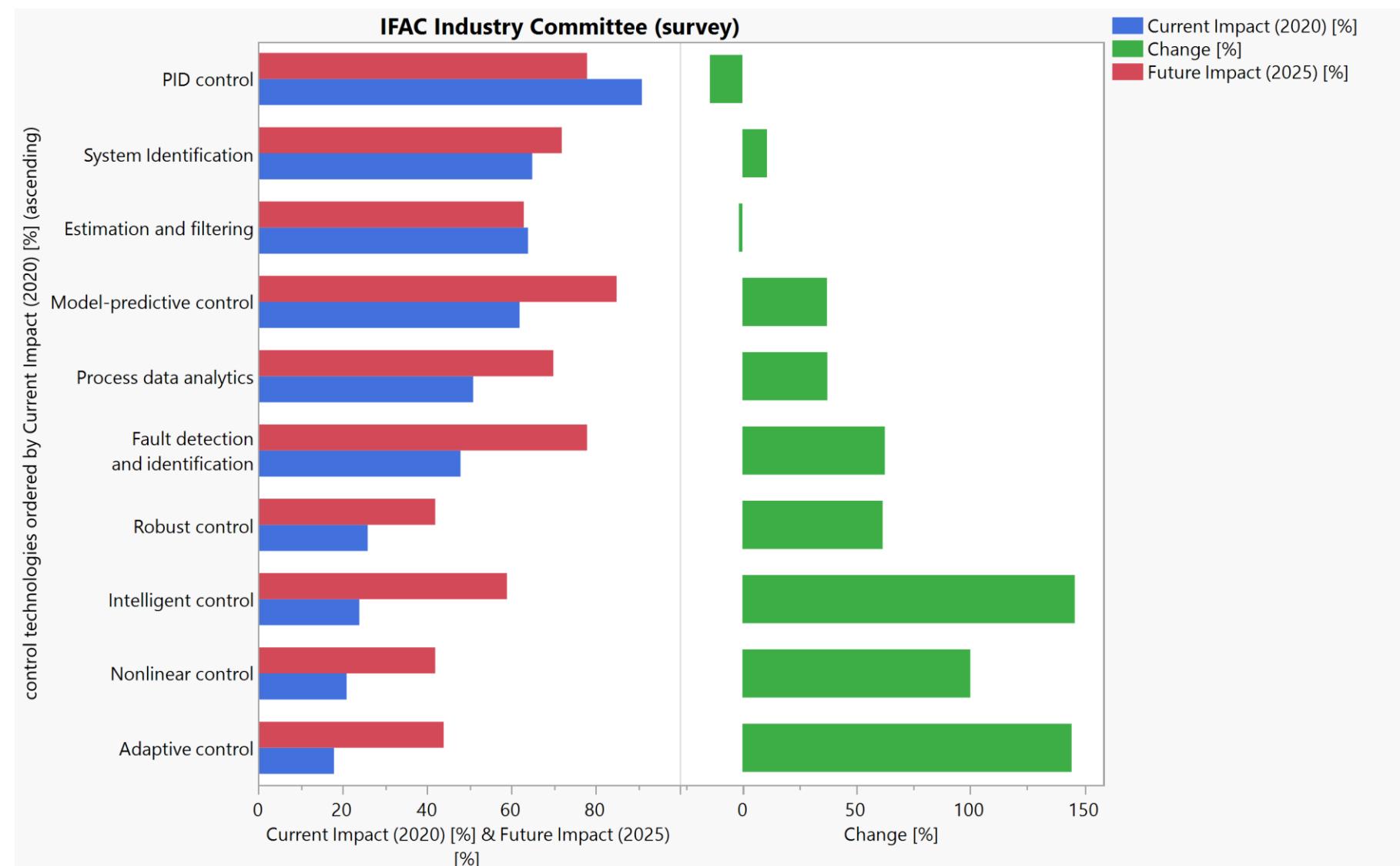
PID control

Control technologies and their impact (present and future)

- Industry engagement with control research: Perspective and messages

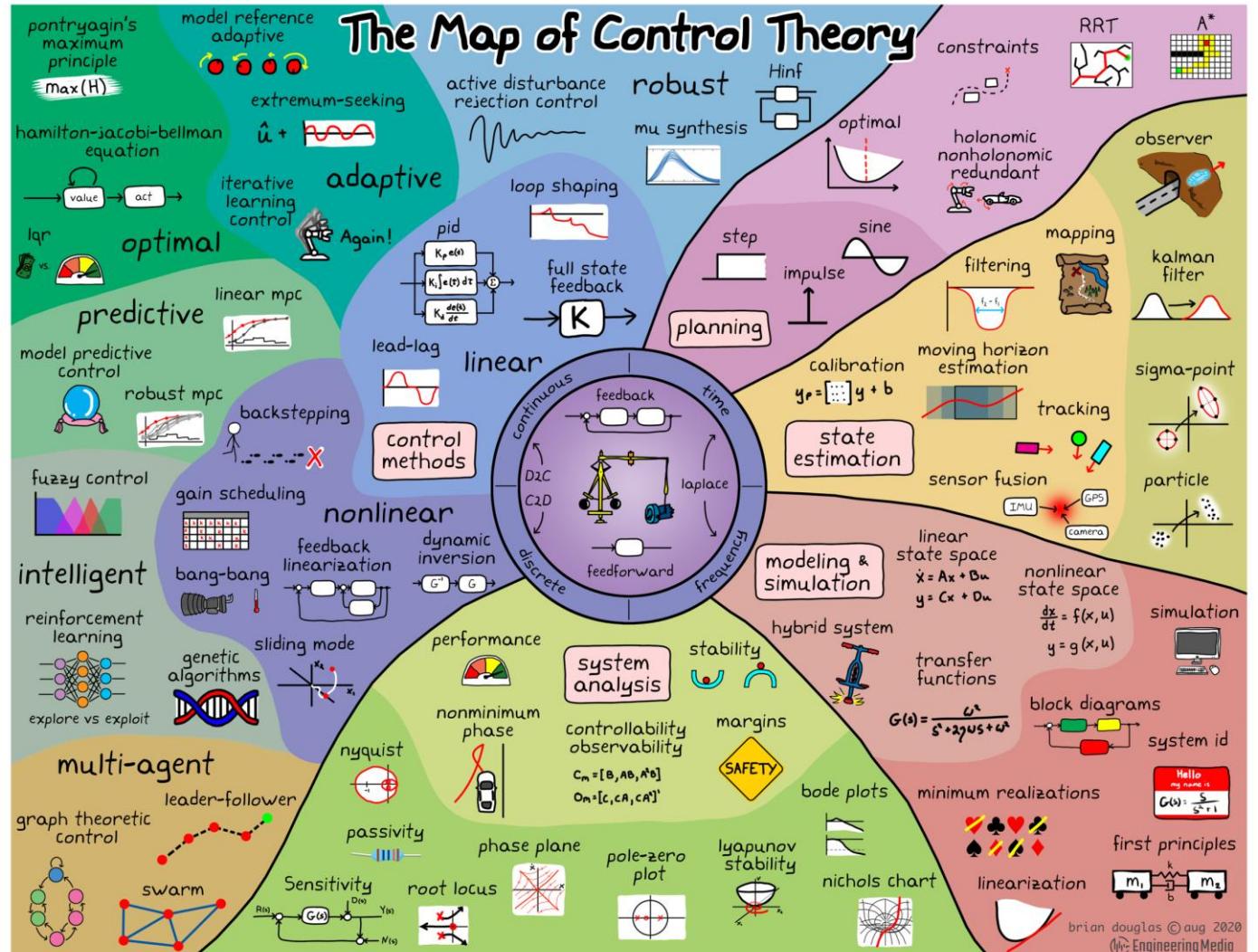
<https://www.sciencedirect.com/science/article/pii/S1367578820300080>

Percentage of survey respondents indicating whether a control technology had demonstrated (“Current Impact”) or was likely to demonstrate over the next five years (“Future Impact”) high impact in practice.



Online resources to learn process control

- <https://www.chemicalengineeringpractice.org/> - Prof.
- <https://engineeringmedia.com/> - Brian Douglas
- [https://apm.byu.edu/prism/index.php/Site/Online Courses](https://apm.byu.edu/prism/index.php/Site/Online%20Courses)
- <http://www.pc-education.mcmaster.ca/>
- Process Control. A Practical Approach. Myke King
- Applying S88. Batch Control from a User's Perspective. Jim Parshall and Larry Lamb



Industrial Data Science

Data integration



The Purdue model (Automation pyramid)



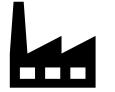
IT -
Informational
Technology

Business Planning & Logistics
Plant Production Scheduling, Business Management, etc.



4 - Establishing the basic plant schedule production, material use, delivery, and shipping
Determining inventory levels.

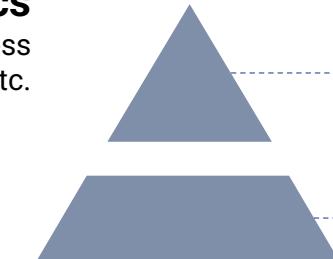
Time Frame: Months, weeks, days, shifts



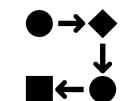
OT -
Operational
Technology

Manufacturing Operations Management

Dispatching production, detailed production, scheduling, reliability assurance



3 - Work flow / recipe control to produce the desired end products. Maintaining records and optimizing the production process.



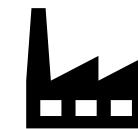
Time Frame: Shifts, hours, minutes, seconds



Manufacturing Control
Basic control, supervisory control, process sensing, process manipulation.



2 - Monitoring, supervisory control and automated control of the production process



1 - Sensing the production process, manipulating the production process

0 - The physical production process

LIMS, ERP & Sensor data

Enterprise Resource Planning (ERP)

- Purchase Orders
- Finished Products and Raw Materials
- Suppliers and Inventory Levels

Tabular data:				
	KPIs	LIMS	ERP	
BatchID	LotID (PO number)	Grade (P1, P2, P3)	Quality (pass/no pass)	Supplier
...
...
...

Rows are independent*

LIMS

- Quality of intermediate
- Specification Limits
- Control Limits

Sensor Data

- P, F, T
- P_{SP} , F_{SP} , T_{SP}
- Contextual data

Sensor data: tags in historians ...

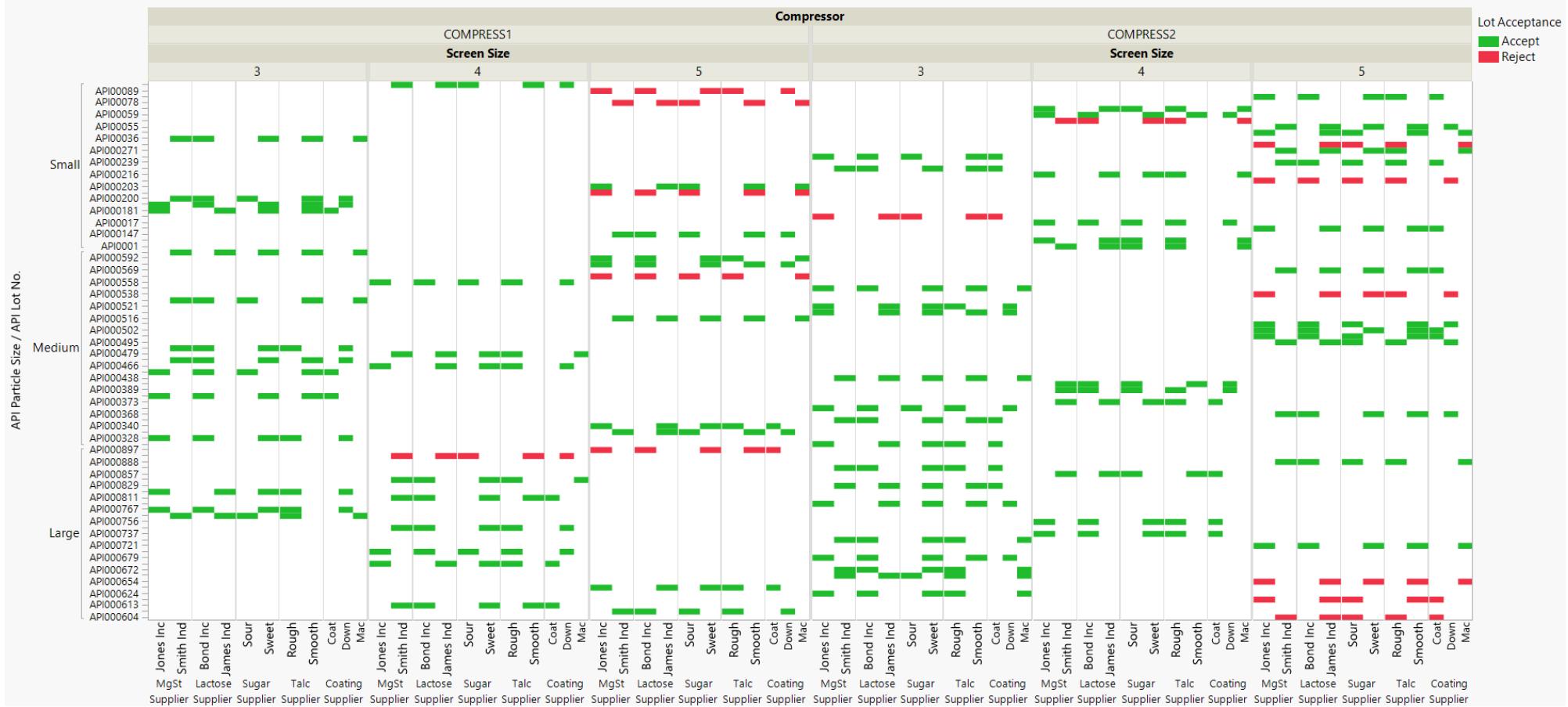
◀ Data integration ▶
(not data duplication)

Datetime	PI_505.PV (Pressure) [bar]	TI_203.PV (Temperature) [degC]	FIC_101.MV (Valve position) [%]	FI_101.PV (Flowrate) [m³/min]
...
...

Rows are not independent (sampling, dynamics...)

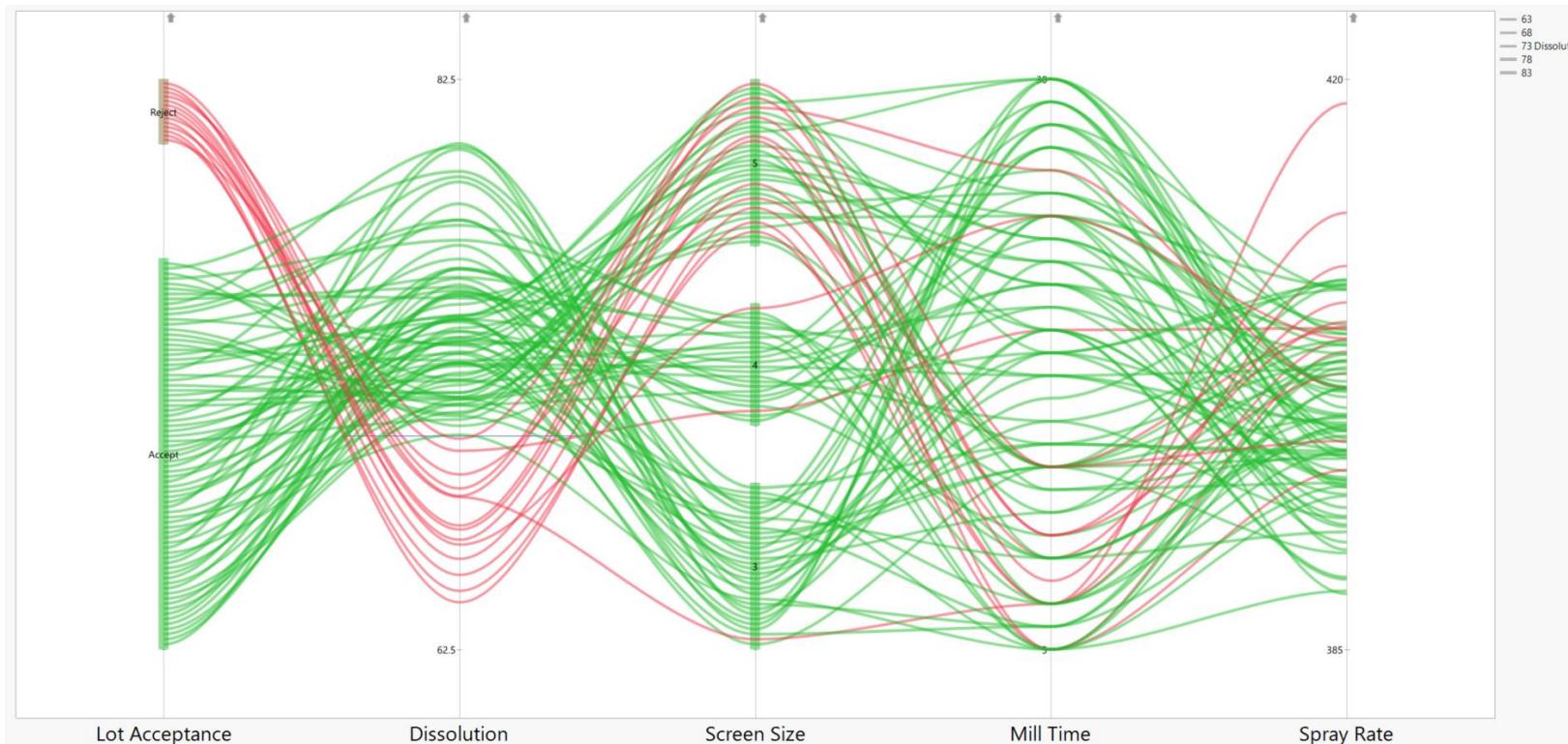
To leverage the full potential of the information captured by the business. The manufacturing data can be complemented with the operations management. This type of data integration can bring insights on information such as the supplier of the raw materials to investigate possible issues

ERP data example



ERP data contains multiple types of information. It is itself a consolidation of different types of data from different functions (controlling, operations, supply chain, quality, etc). Often, a good portion if the inputs to the system are manual. Among the challenges we have: finding out the latest up to date, name recoding (versioning), applying the right filters, different maturity depending on site.

Screening SAP-type data



Since multiple types of information are maintained in ERP, also multiple data formats and types are available. Our tools need to cope with continuous and discrete data as well as with tabular and time series formats.

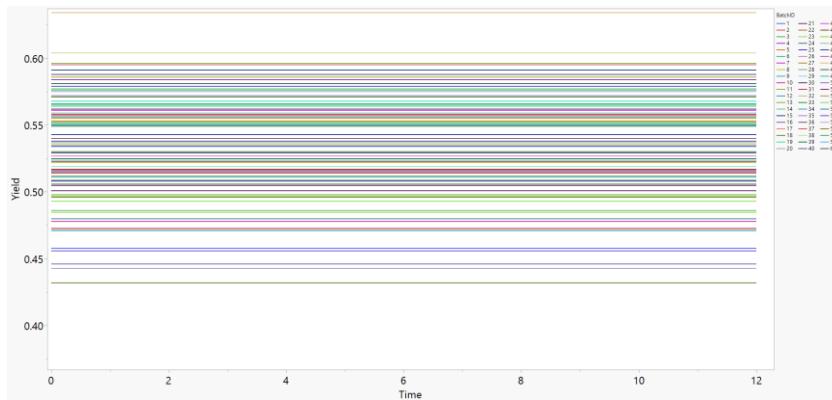
Predictor	Contribution	Portion	Dissolution	Rank ^
Screen Size	157.567	0.2760		1
Mill Time	125.876	0.2205		2
Spray Rate	61.262	0.1073		3
Blend Time	39.452	0.0691		4
Coating Viscosity	38.897	0.0681		5
Blend Speed	36.890	0.0646		6
Atom. Pressure	23.873	0.0418		7
Inlet Temp	15.500	0.0272		8
API Particle Size	12.864	0.0225		9
Force	10.962	0.0192		10
Lactose Supplier	10.167	0.0178		11
Exhaust Temp	9.853	0.0173		12
Talc Supplier	8.906	0.0156		13
Coating Supplier	5.972	0.0105		14
Sugar Supplier	5.556	0.0097		15
Compressor	4.585	0.0080		16
MgSt Supplier	2.629	0.0046		17

The right filtering conditions must be applied when using screening methods. Understanding of the discrete elements is key

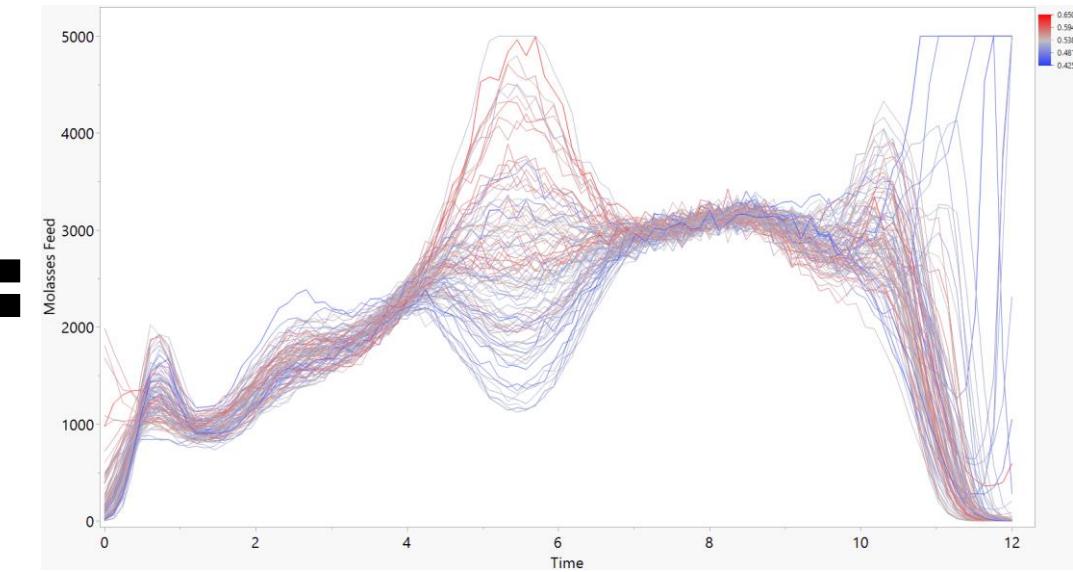
Fermentation example (quality + process)

- Quality data per batch
- Measurement is taken at the end of the batch
- Repeatability and real sampling time must be verified

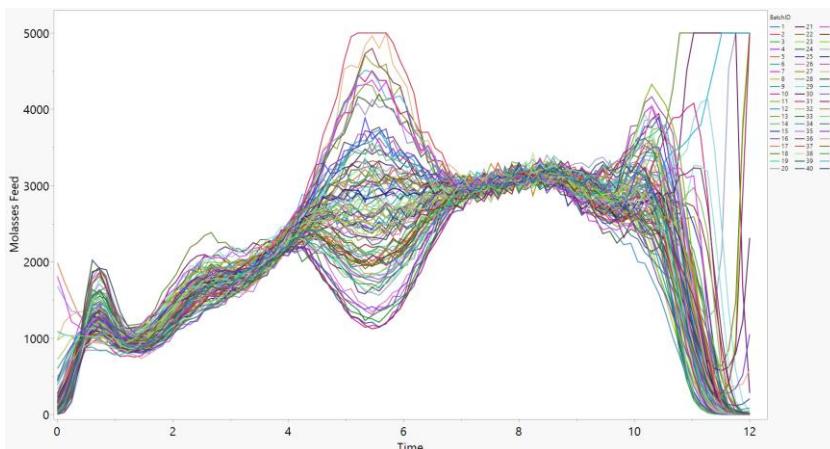
LIMS



Integrated Data

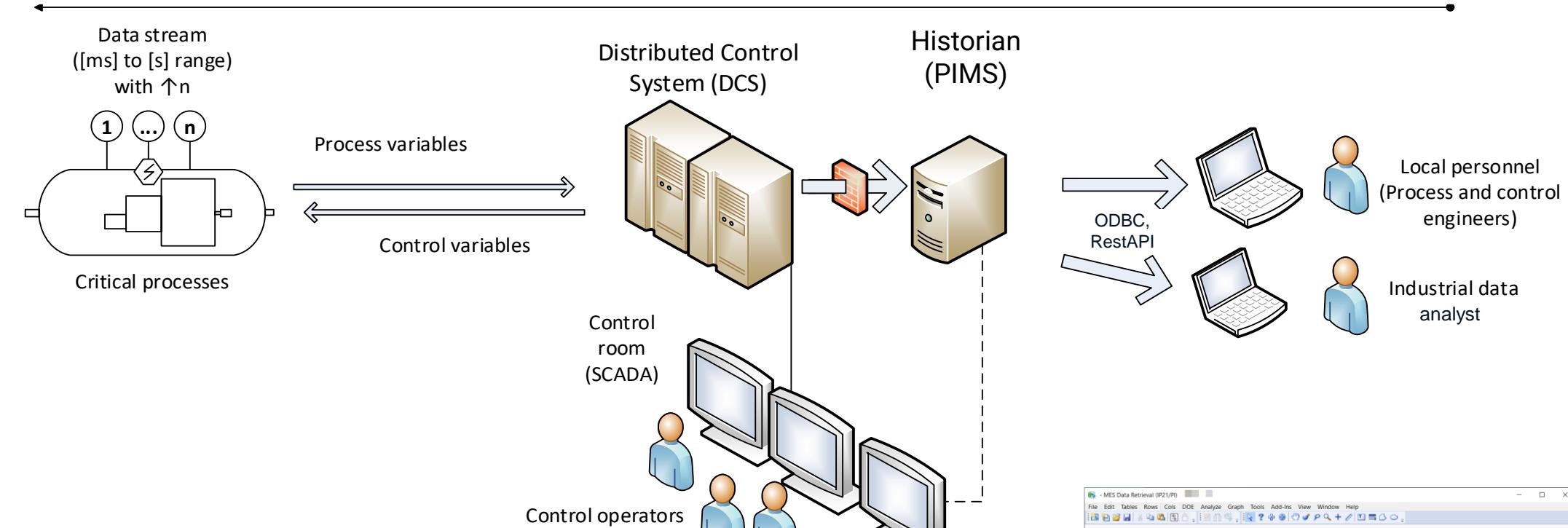


Sensor Data



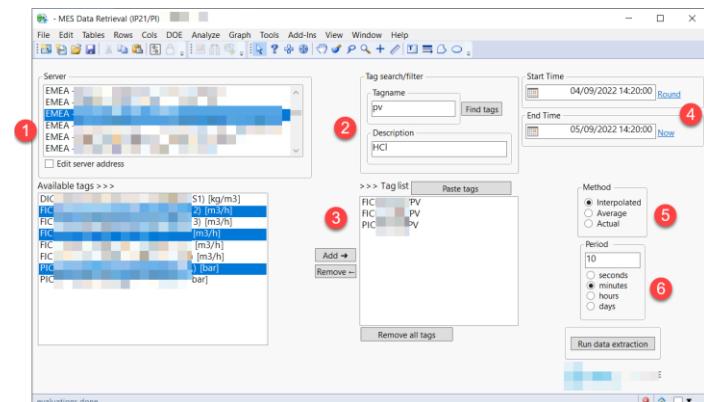
Join of process data and quality data. Color now belongs to quality. You need a unique batch ID and in-sample time

Industrial data (open-source) native connectors



tagreader 4.2.5

[https://github.com/equinor/
tagreader-python](https://github.com/equinor/tagreader-python)



<https://github.com/industrial-data/JMP-MES-connector>

Industrial Data Science

Hands-on



First steps with JMP

Getting data into JMP, several options:

★ [Directly from Excel \(add-in\)](#) [1.5 min]

★ [Excel file](#) [3 min]

- [Text file](#)
- [Multiple text files](#)
- Optional: [External databases](#)
- [IP21/PI data extraction](#)

Demos for visual analytics:

★ [Graph builder for manufacturing & SAP data](#) (28 min) [dataset [I](#), [II](#)]

★ [Process data and ML](#) (tags, sensors) (45 min) [dataset]

The screenshot shows the JMP Learning Library interface. At the top, it says "Learn JMP: The one place to access all JMP learning materials" and "Subscribe to Category". Below are three main sections: "STEP 1 Learn the Basics" (with a line graph icon), "STEP 2 Build Your Expertise" (with a person icon), and "STEP 3 Go For More" (with a graduation cap icon). "Go For More" contains a grid of resource categories: Statistical Thinking for..., Data Access & Exploration (12 resources), Design of Experiments (18 resources), Data Analysis & Modeling (25 resources), Predictive Modeling &... (24 resources), Quality, Process Engineering,... (23 resources), Consumer & Market Research (3 resources), Automation & Scripting (12 resources), and Content Organization &... (5 resources).

All learning material: <https://community.jmp.com/t5/Learn-JMP/ct-p/learn-jmp>

1-2 min tutorials: https://www.jmp.com/en_us/learning-library.html

Control charts: <https://community.jmp.com/t5/Mastering-JMP/Analyzing-Improving-and-Controlling-Process-Stability-and/ta-p/500753>

Training material for the Federal Government: www.jmp.com/fedgov

Online course (+certification): https://www.jmp.com/en_us/online-statistics-course.html

SPC training: <https://community.jmp.com/t5/Learning-Center/JMP-Statistical-Process-Control-Course/ta-p/575330>

Summary

- Machine learning concepts are not new for chemical engineers
 - Supervised learning → Regression
 - Unsupervised learning → Statistical Process Control charts
- Process, control and automation knowledge is essential to ask relevant questions and interpret results in process industries
- Democratization of machine learning requires self-service access to all process data and analytics
- One of the main applications of machine learning is to understand chemical processes faster (e.g., selecting relevant sensors (tags) amongst noise)
- If the answer is obvious, question should be redefined. Example: data analysis points to mass or heat balance, which are already known. Target variable should be redefined to situations where assumptions do not hold (i.e., discrepancy models).

*Business impact is measured in
tones, not lines of code*

Day 2

Monitoring and screening process data

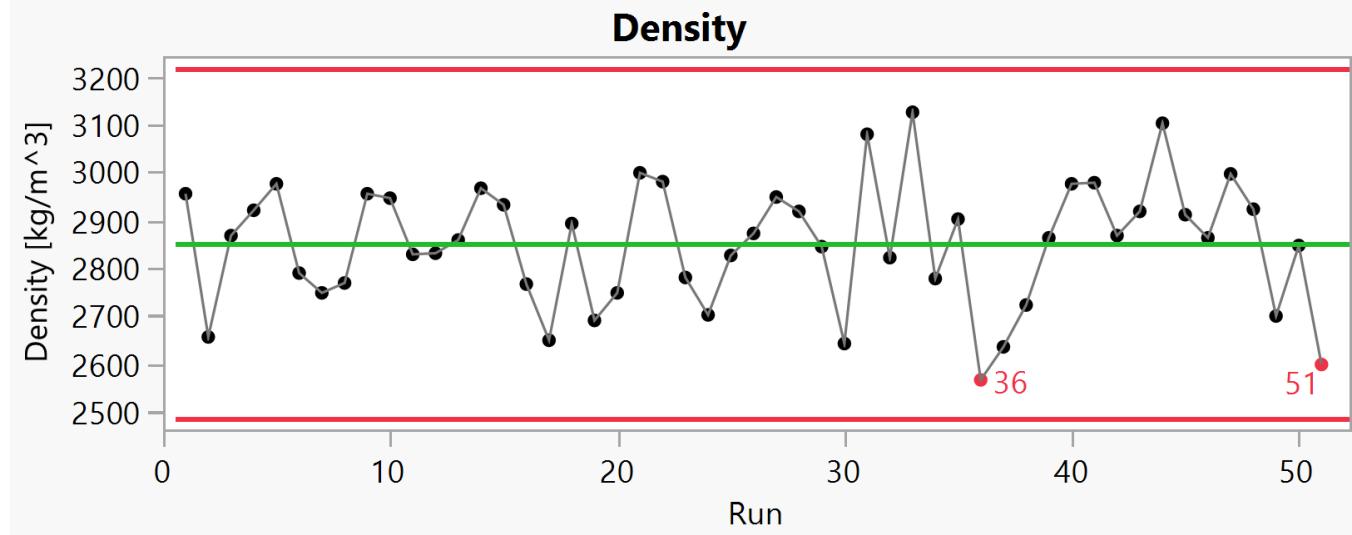
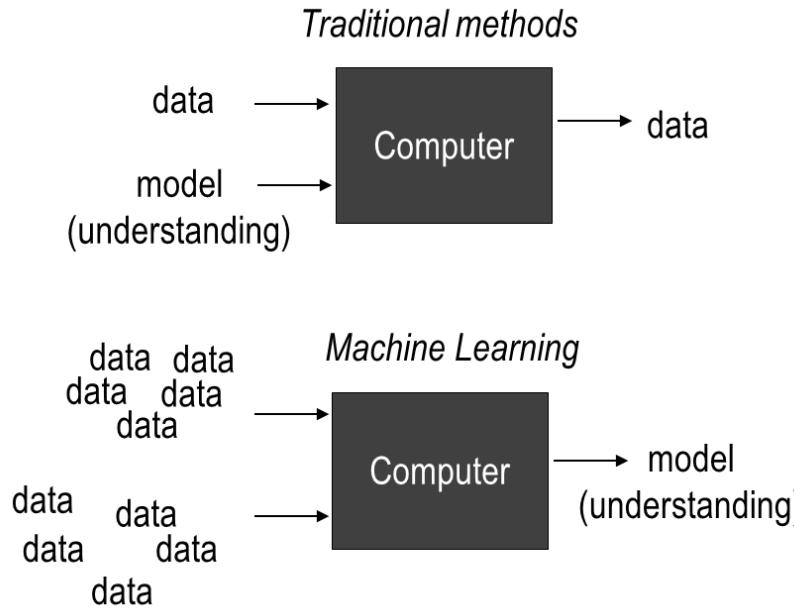


Monitoring and screening process data

Anomaly detection

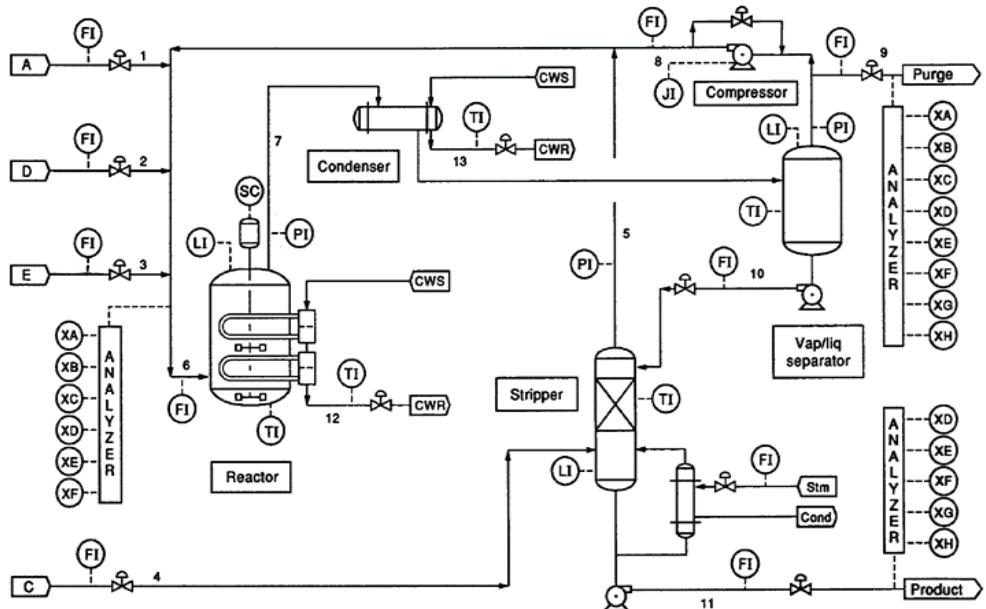


Control charts (Unsupervised learning)



Control charts are a form of unsupervised models where only input training data is given and a statistical model is built (mean and standard deviation, since controlled processes should follow a normal distribution). Finally, we can classify data points as in- or out-of-control.

Anomaly detection example: Tennessee Eastman Plant



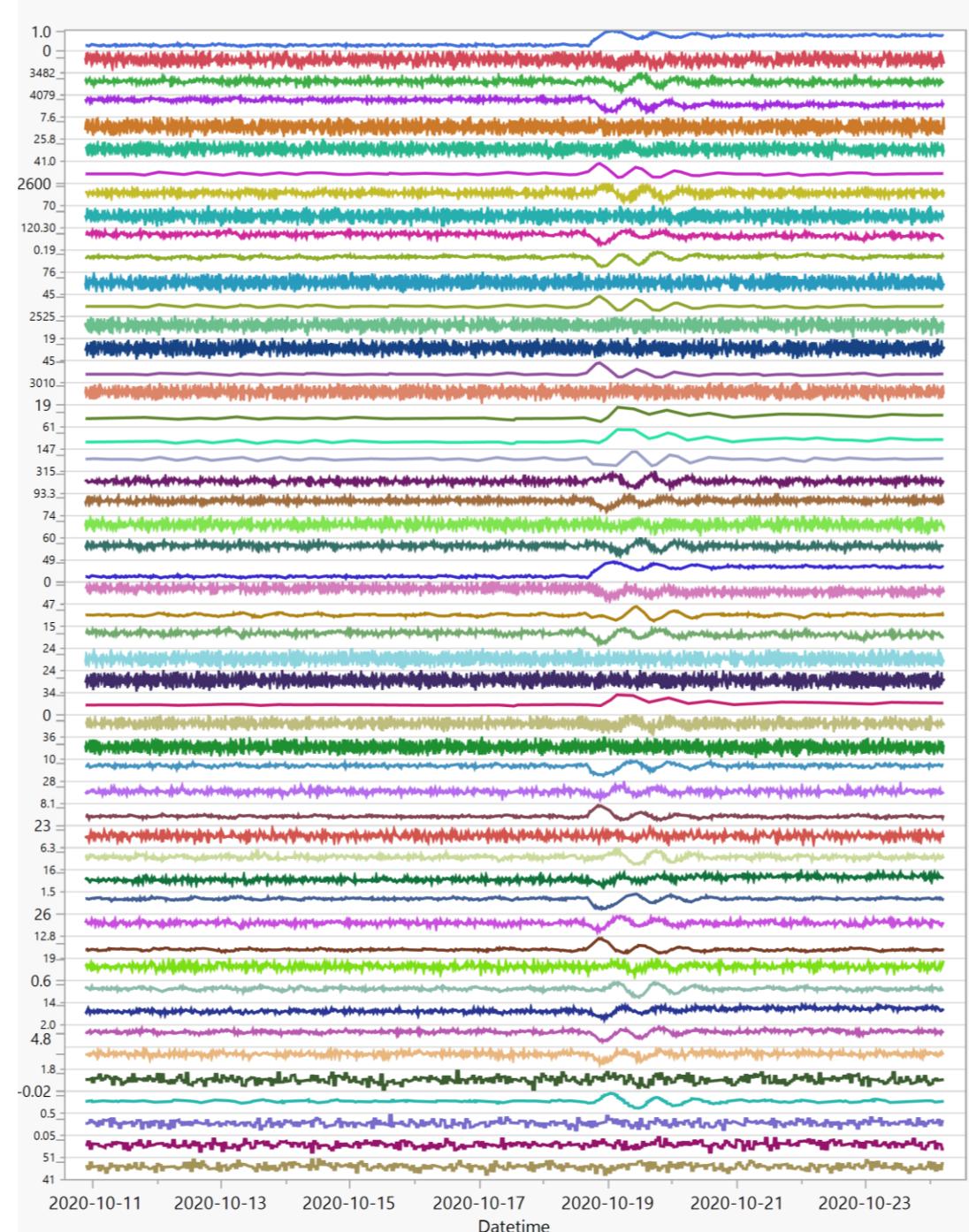
Tennessee Eastman plant has transitioned between two regimes, multiple measurements have changed.

Src: <https://pubs.rsc.org/en/content/articlelanding/2022/re/d1re00541c>

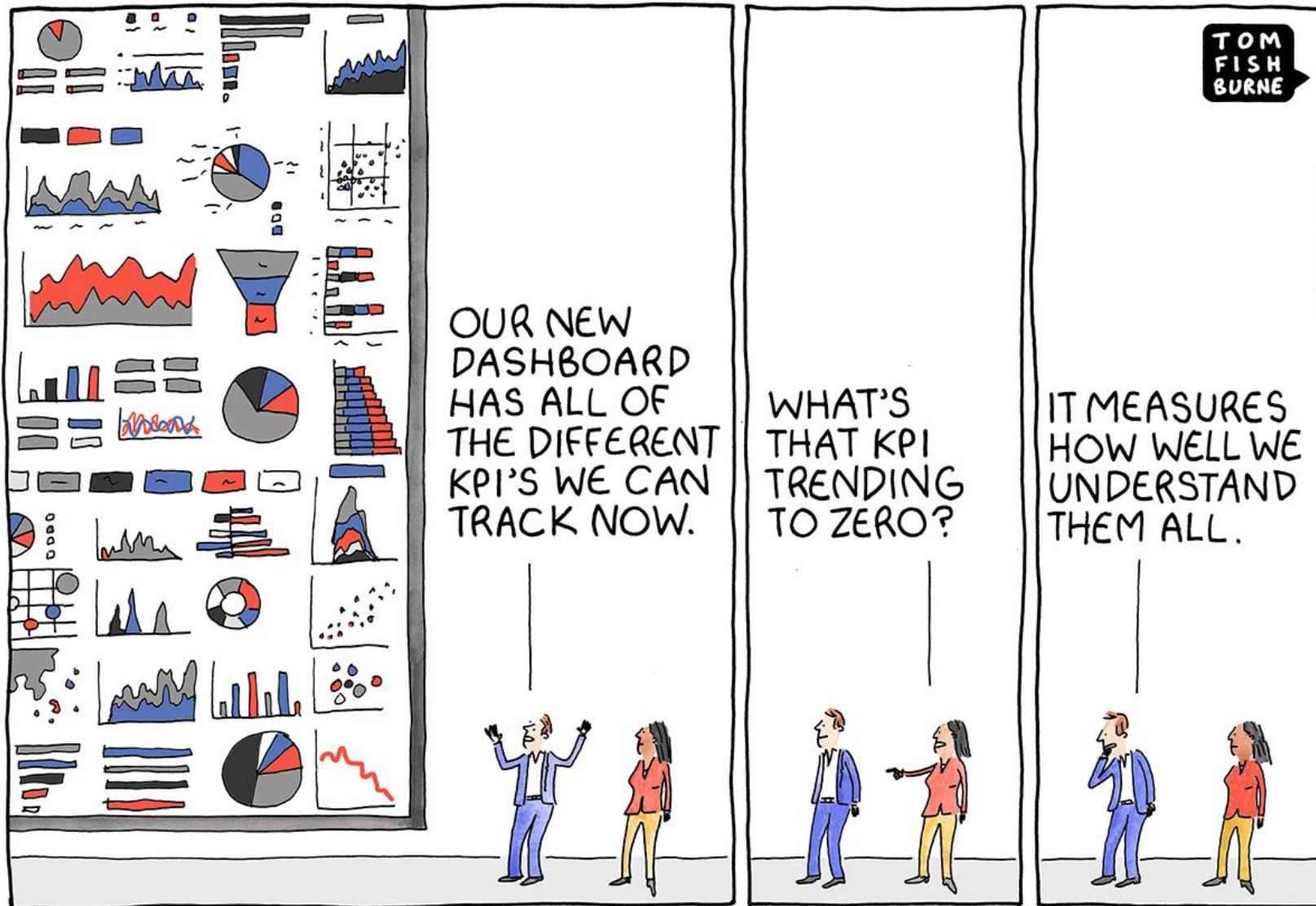
PV

MV

Quality
(online)

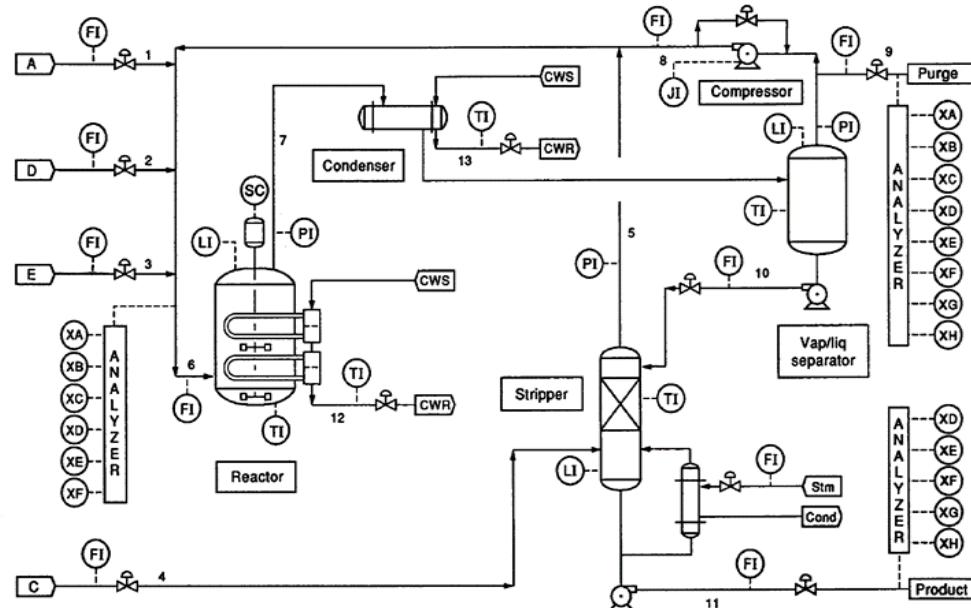


Let's check our KPIs

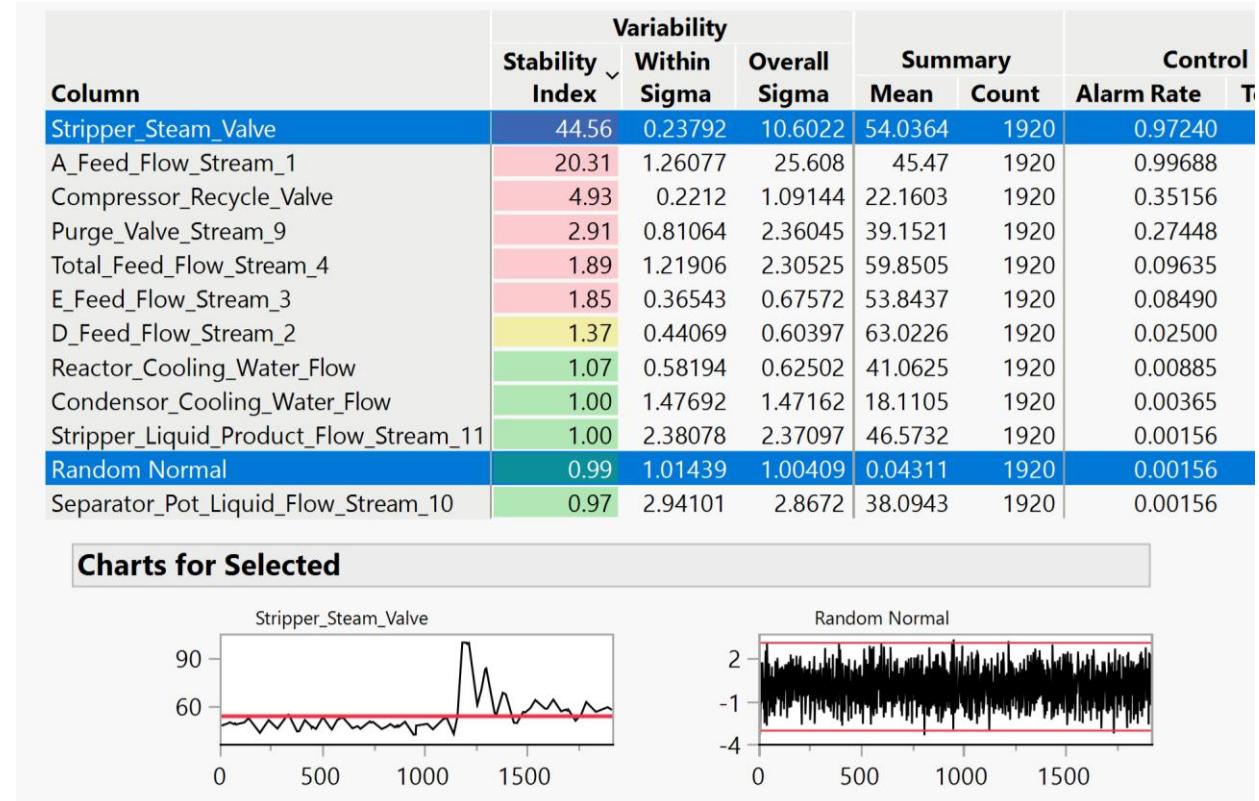


© marketoonist.com

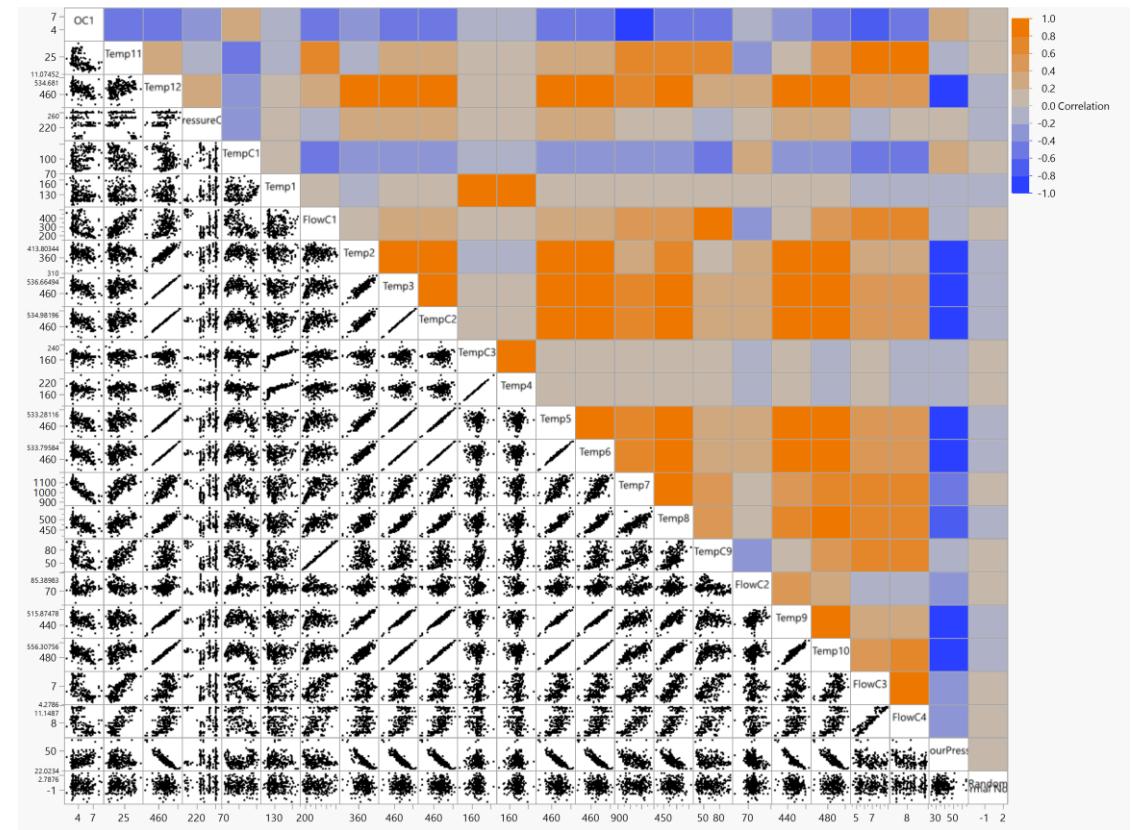
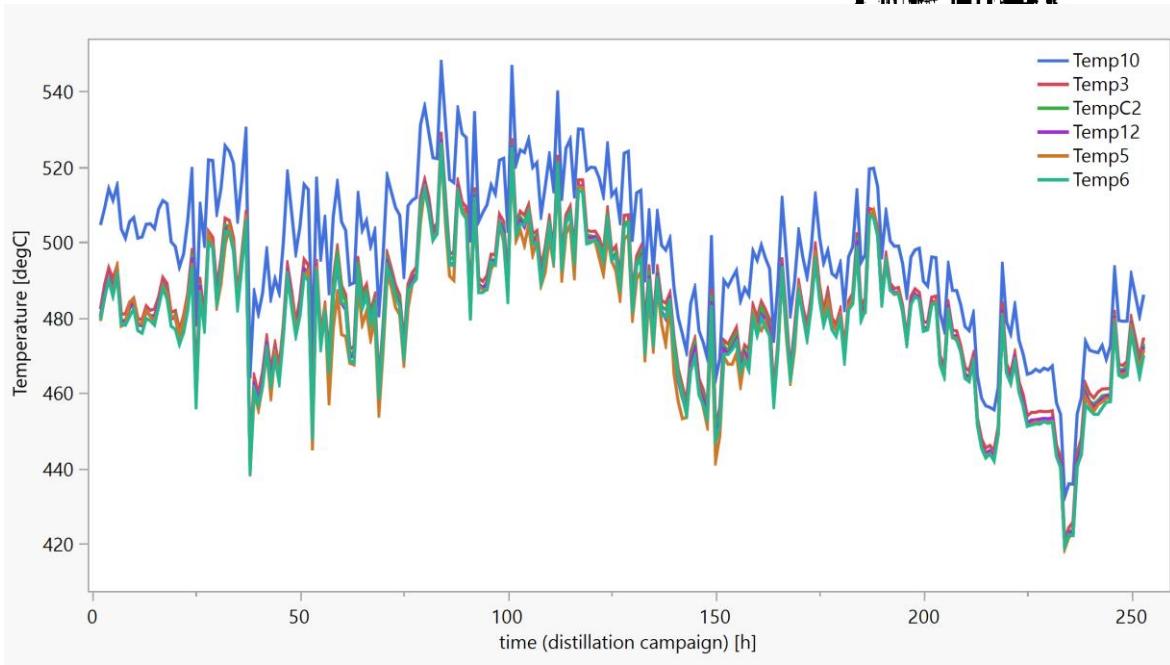
Anomaly detection (univariate analysis)



Univariate analysis can be used to identify major process changes, specially when interactive data filters are applied.

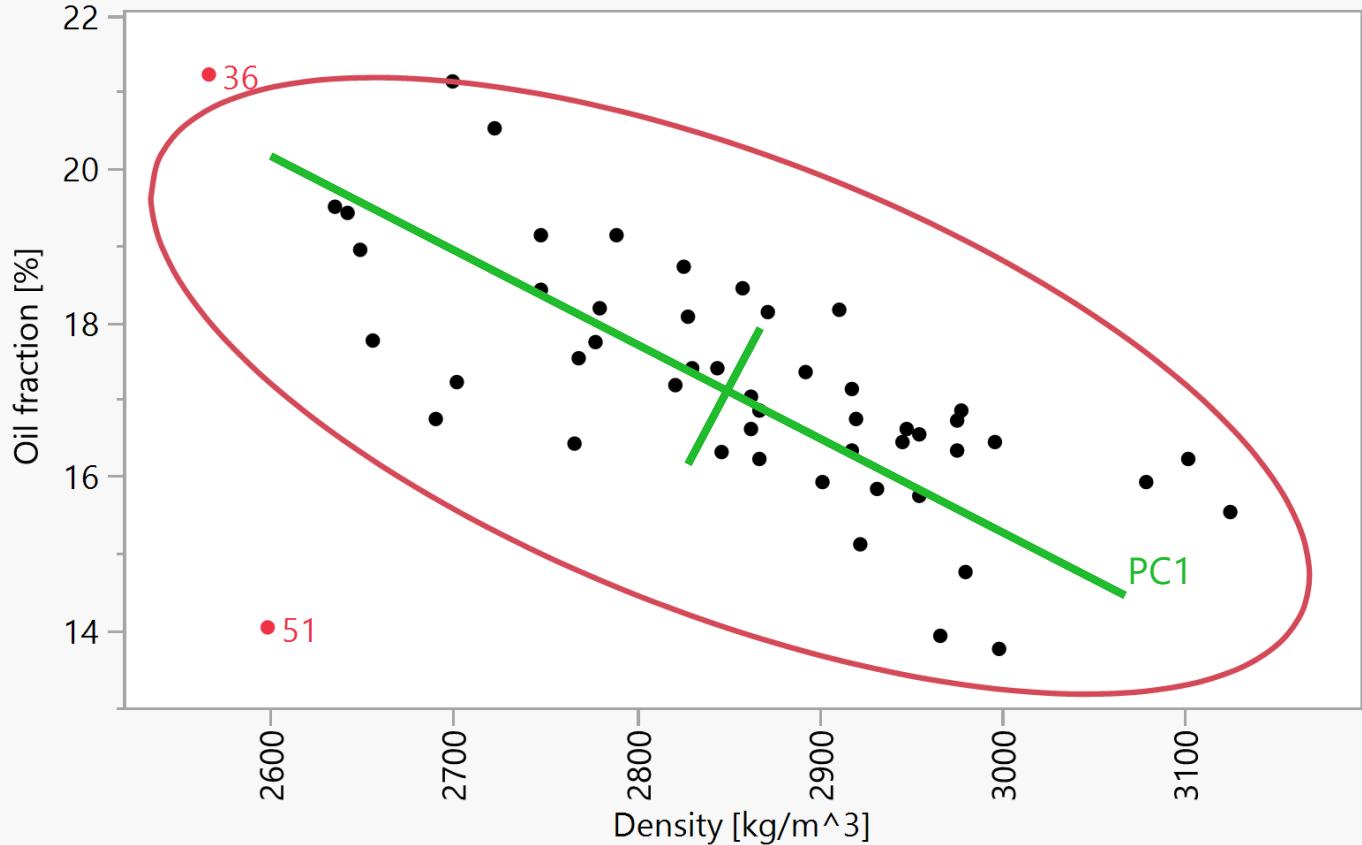
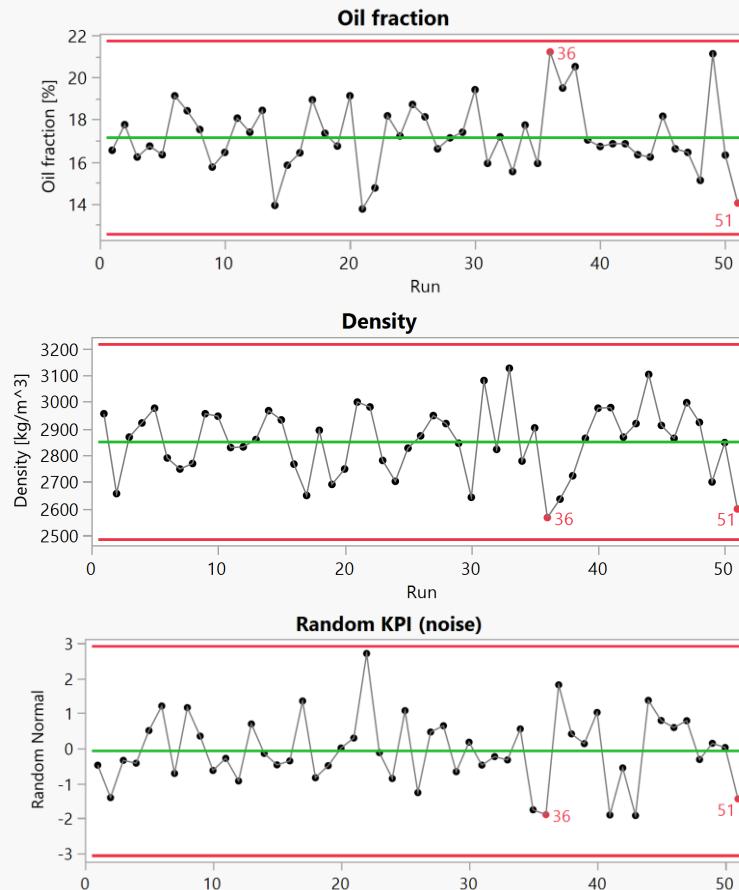


Redundant sensors



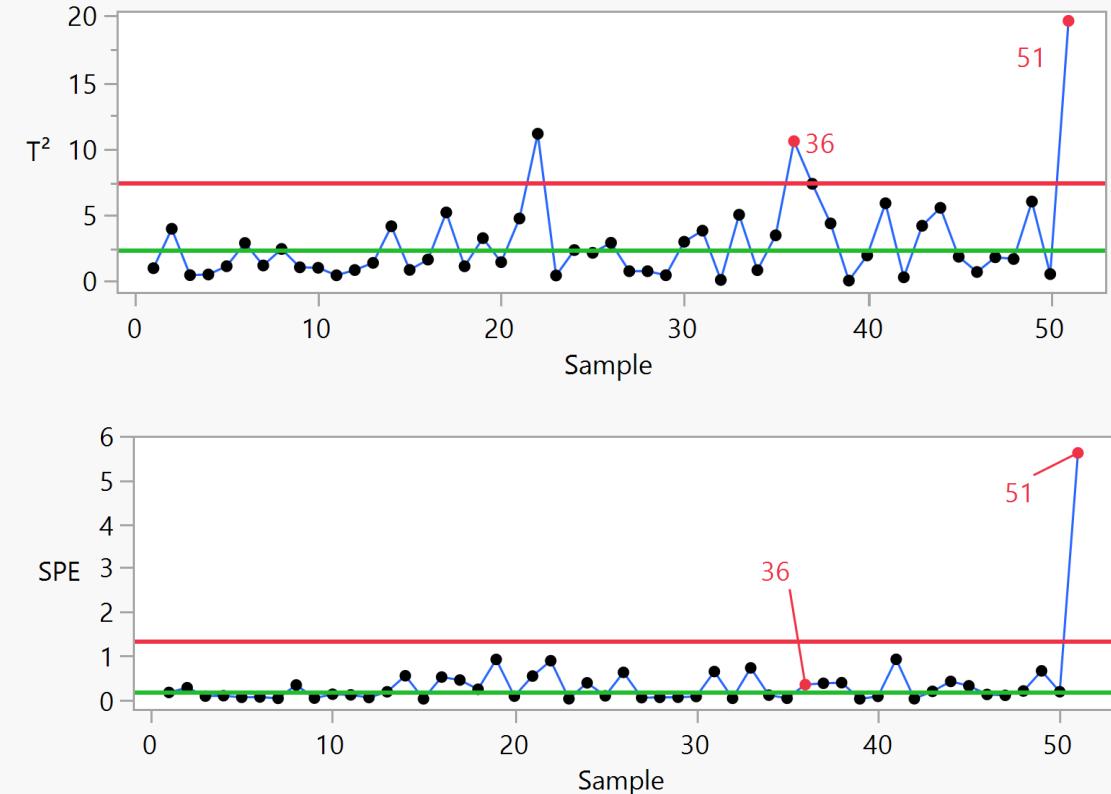
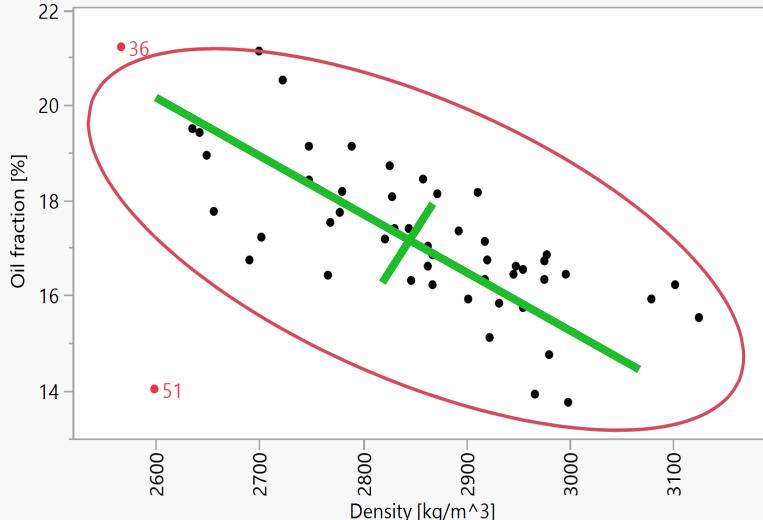
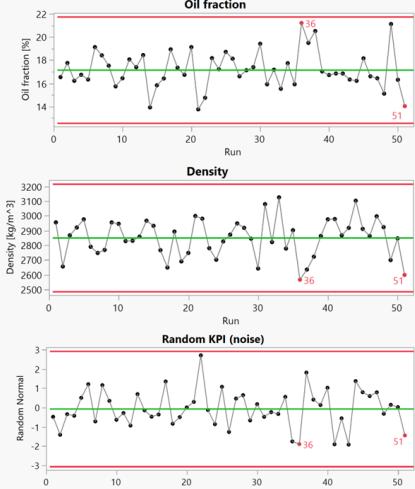
A common example of multivariate analysis is to monitor multiple redundant sensors. When several thermocouples are used to measure a critical temperature, these can be summarized by taking the average of all the sensor readings. The average is a linear combination of all these terms with equal weight. This way, the information is being reduced to one latent variable, the temperature we want to monitor. If a big variation exists between the average of the sensors and one thermocouple, in particular, an alert can be triggered. This reasoning is the same behind principal component analysis (PCA), and it has been widely used for multivariate process analysis, monitoring, and control.

Multivariate control charts (PCA)



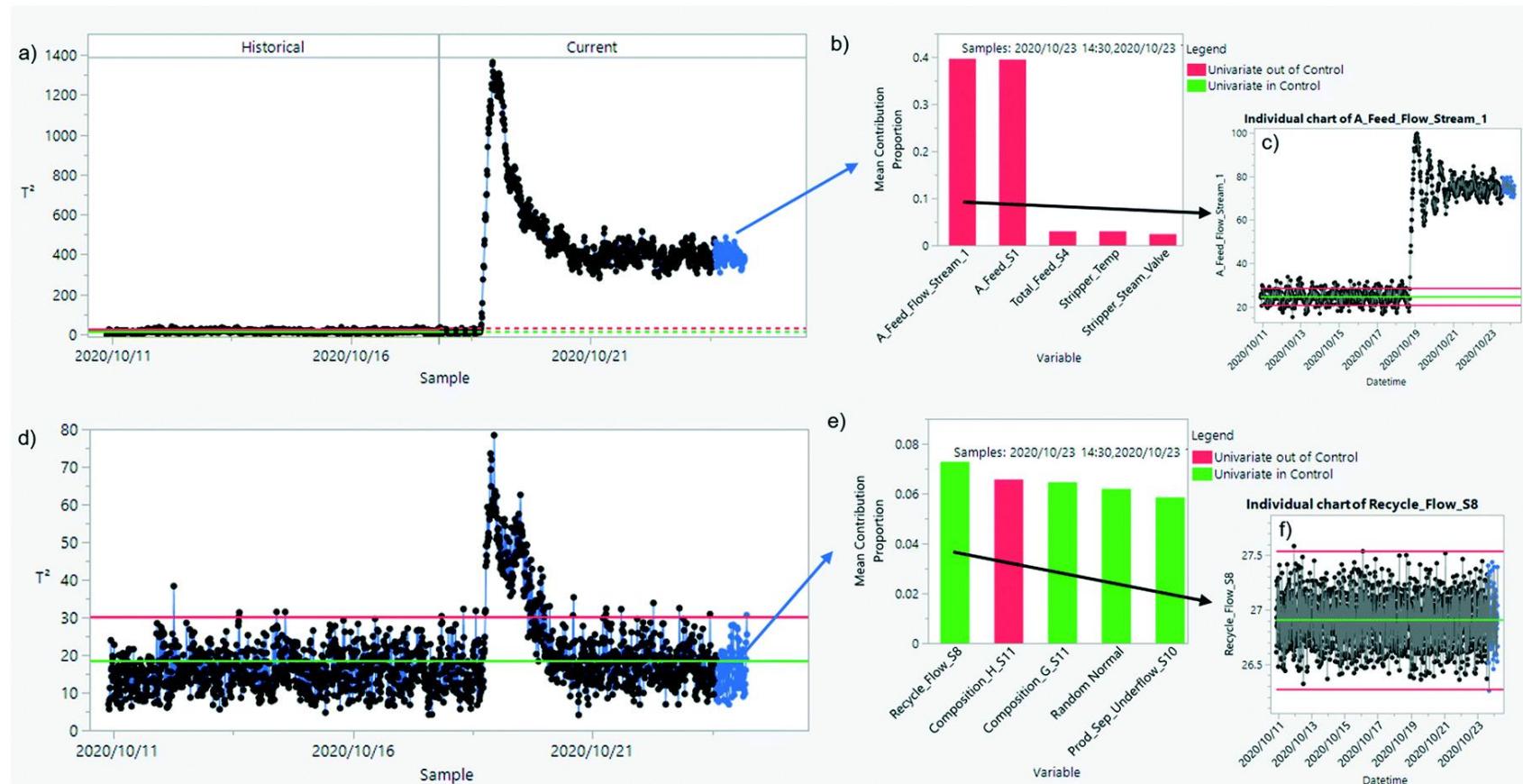
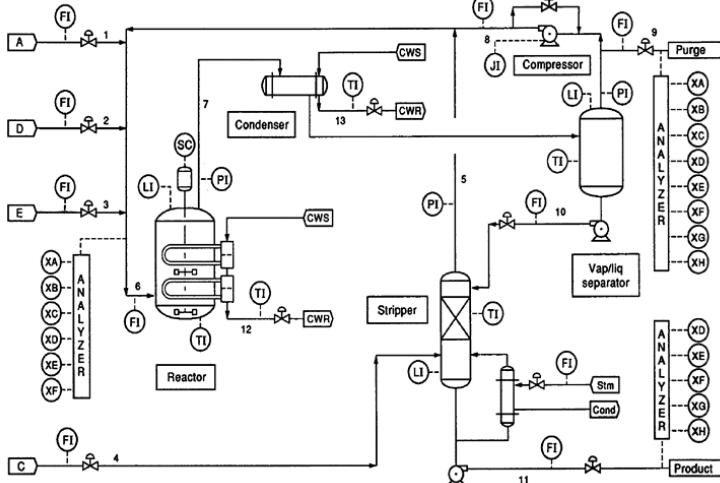
Highlighted samples appeared to be in control (individual control charts). Co-linearity among process variables can be used to detect multivariate anomalies.

Multivariate control charts (T² and SPE)



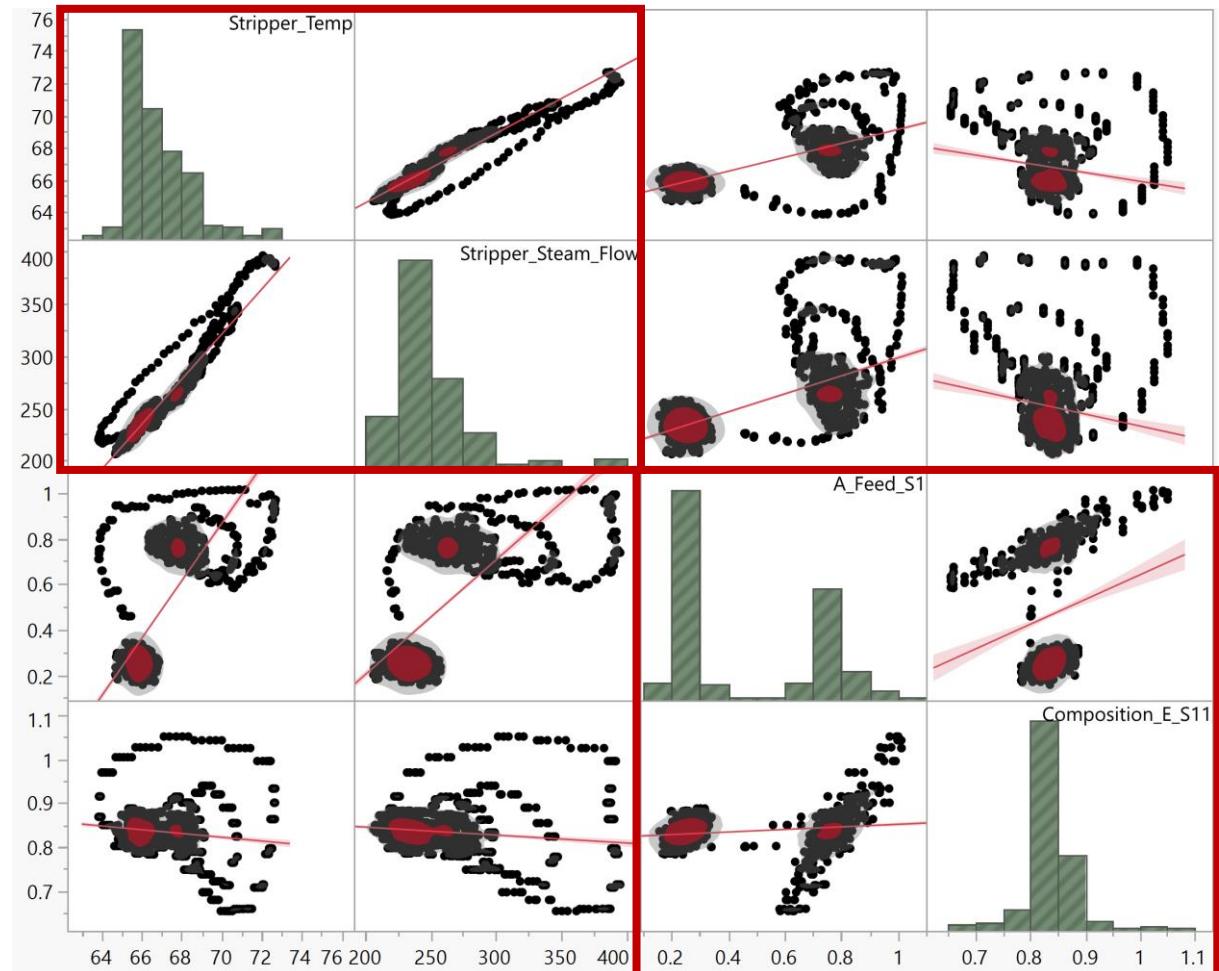
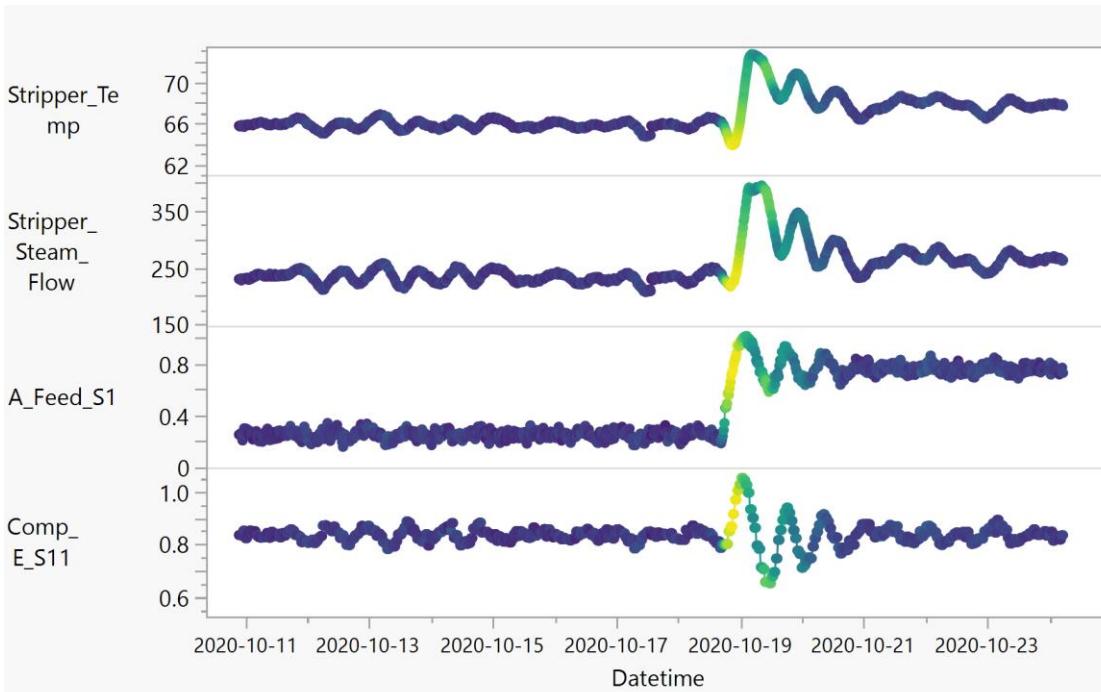
T² and SPE are two popular metrics to identify outliers. T² identifies points that are far from previous data while SPE focuses on those points that do not follow the co-linearity identified across redundant sensors

Anomaly detection (multivariate analysis)



A transition between two steady-state regimes is detected using a multivariate control (PCA). If the model is built using historical data before the perturbation (a) the step changes in the feed flow of chemical A (b and c) are found in the current dataset for the points highlighted in blue. If all historical data is used to build the model (d) the contribution of recent data points in blue (e) shows signals close to random noise. The plant wide control in the simulation stabilizes the control loops and anomalies are only seen in the transition period, even though the plant is operating in a different state for chemical A.

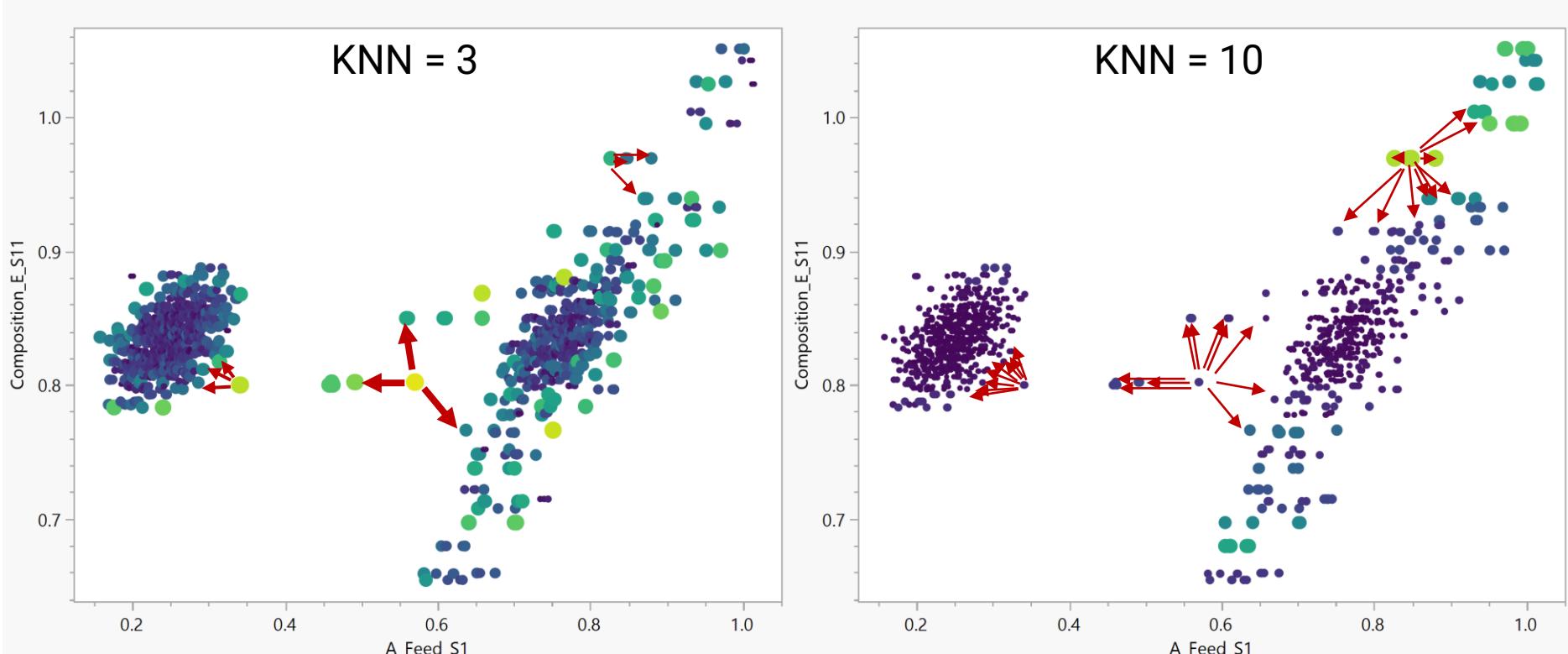
Limitations of multivariate analysis with process data



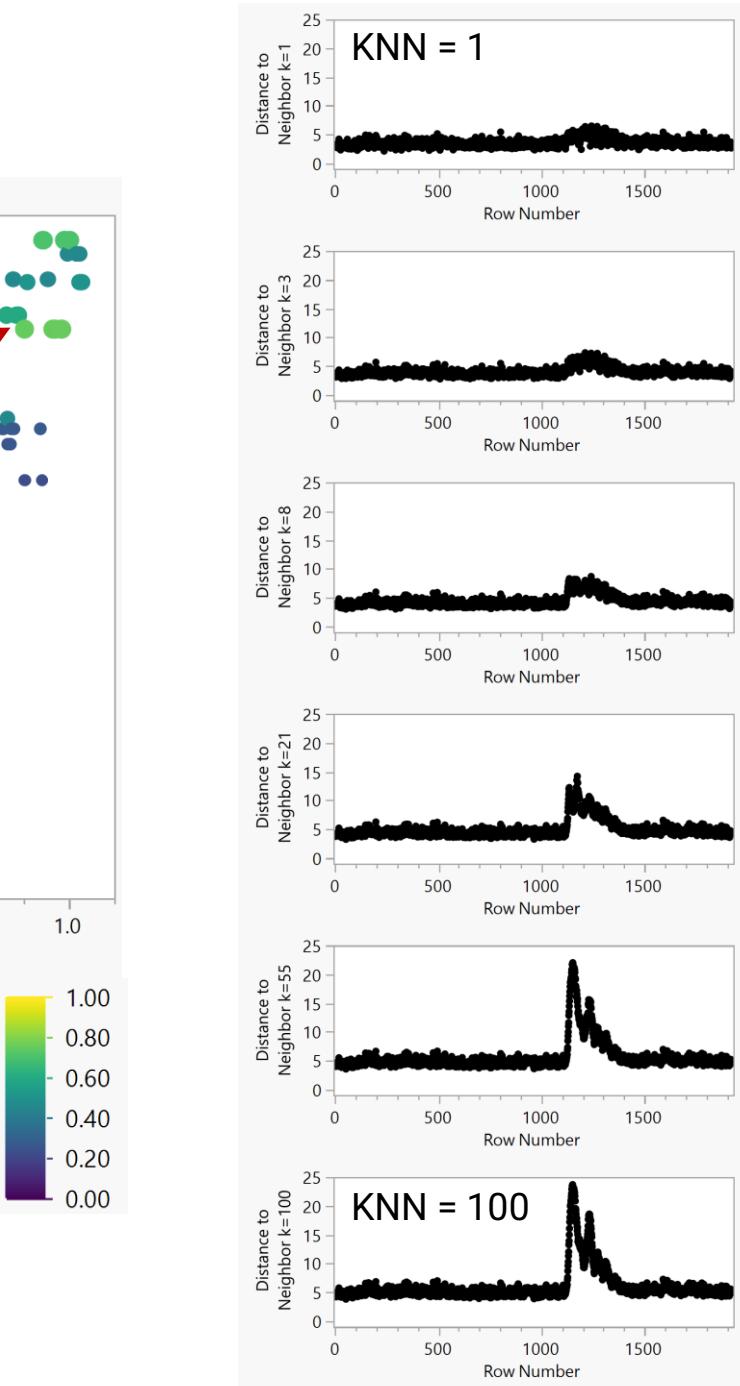
- System dynamics interfere with the analysis, requiring more robust statistics. KNN is another popular algorithm for anomaly detection, in the example above highlighting a process change.

- Steady state regimes can be identified by looking at the stability of main process variables (e.g., flows), or higher point-density areas (colored in red).

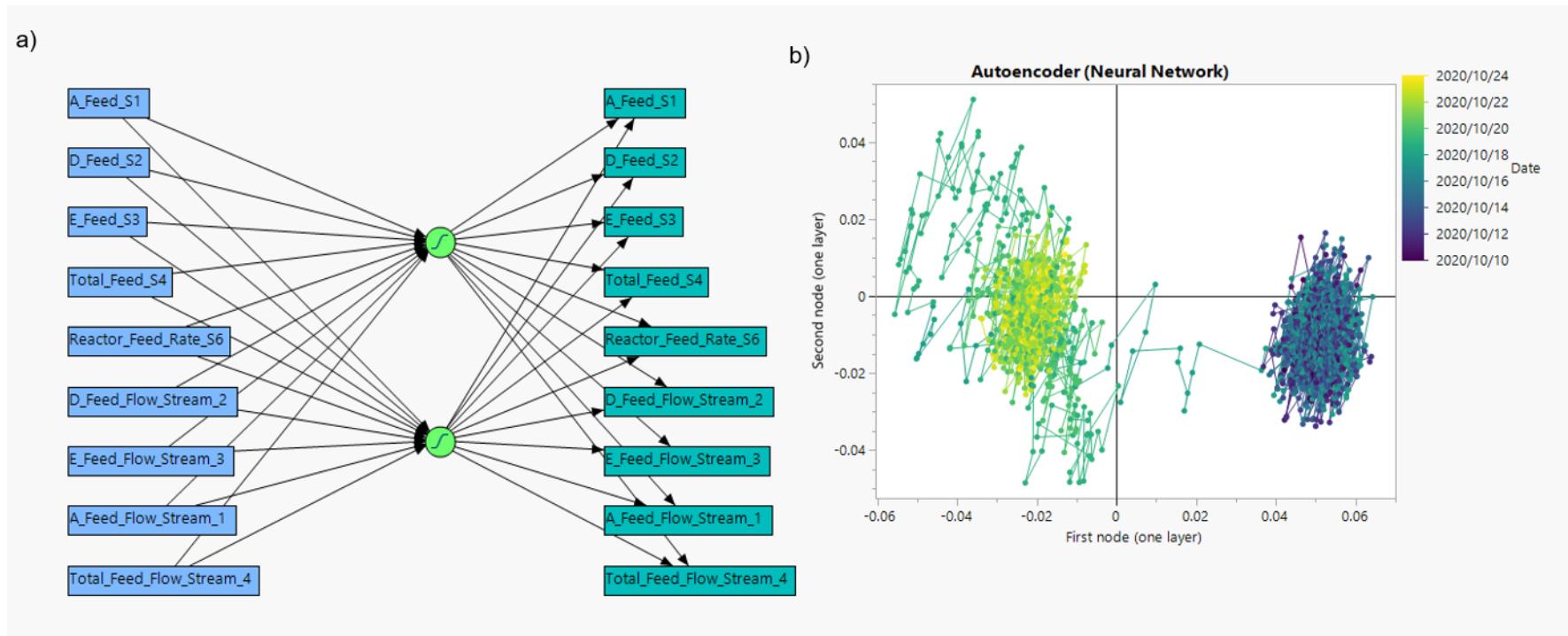
Anomaly detection (K Nearest Neighbors)



K Nearest Neighbor can identify an outlier based on distance. For each value of k , a Euclidean distance from each point to its k th nearest neighbor is calculated (color). A small value of k (left) can miss identifying points as outliers and a large value of k (right) can falsely classify points as outliers.

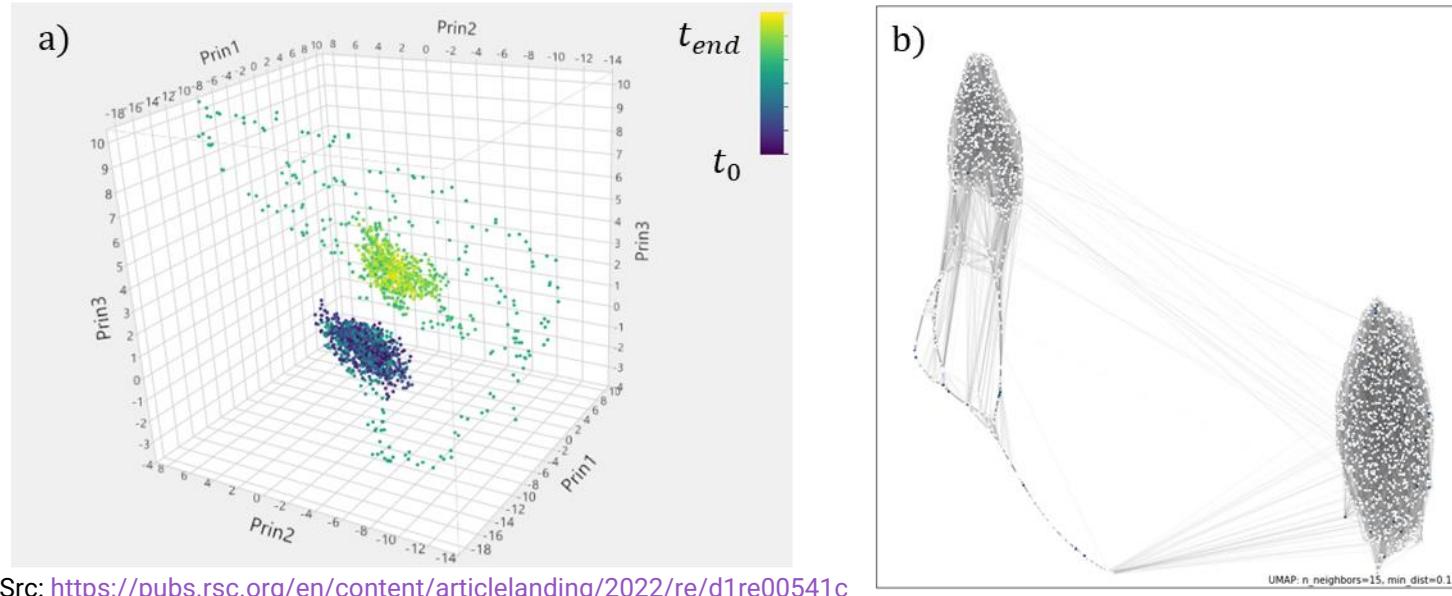


Anomaly detection (autoencoders)



Similar to PCA, autoencoders are neural networks (linear functions with saturation) that can reduce the dimensionality of the data. They achieved this by restricting the number of combinations in the middle layers. This way, correlated data (e.g., redundant measurements) can be reconstructed (a).

Anomaly detection (UMAP)



The transition between two steady-state regimes (Tennessee Eastman Plant) visualized with (a) PCA and (b) UMAP.

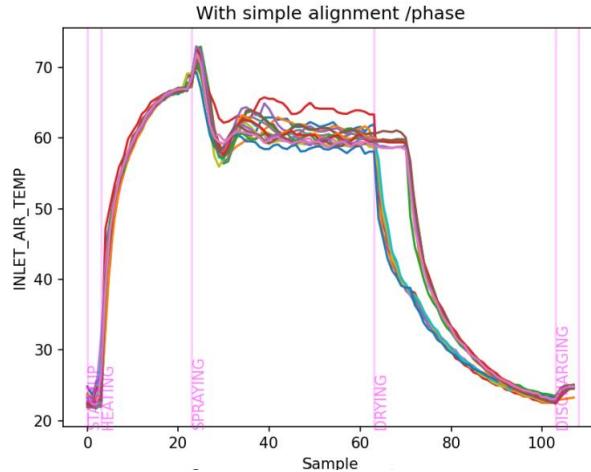
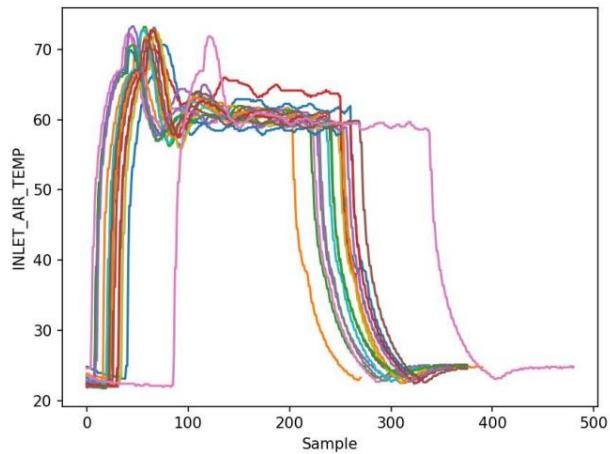
UMAP (Uniform Manifold Approximation and Projection) can outperform and bring new insights when compared with traditional methods. In this example, UMAP demonstrates better separation using only two dimensions, best local and global structure.

Detailed comparison: Dimensionality reduction for visualizing industrial chemical process data (Joswiak et al. 2019) - Dow

<https://www.sciencedirect.com/science/article/abs/pii/S0967066119301728>

Python package with applications to industrial batch data

PyPhiBatch – Multivariate Analysis of Batch Processes



Phi toolbox for multivariate analysis by Sal Garcia
<https://github.com/salvadorgarciamunoz/pyphi>

Years of industrial research made open-source:
Process Analytics Course at Sargent Centre

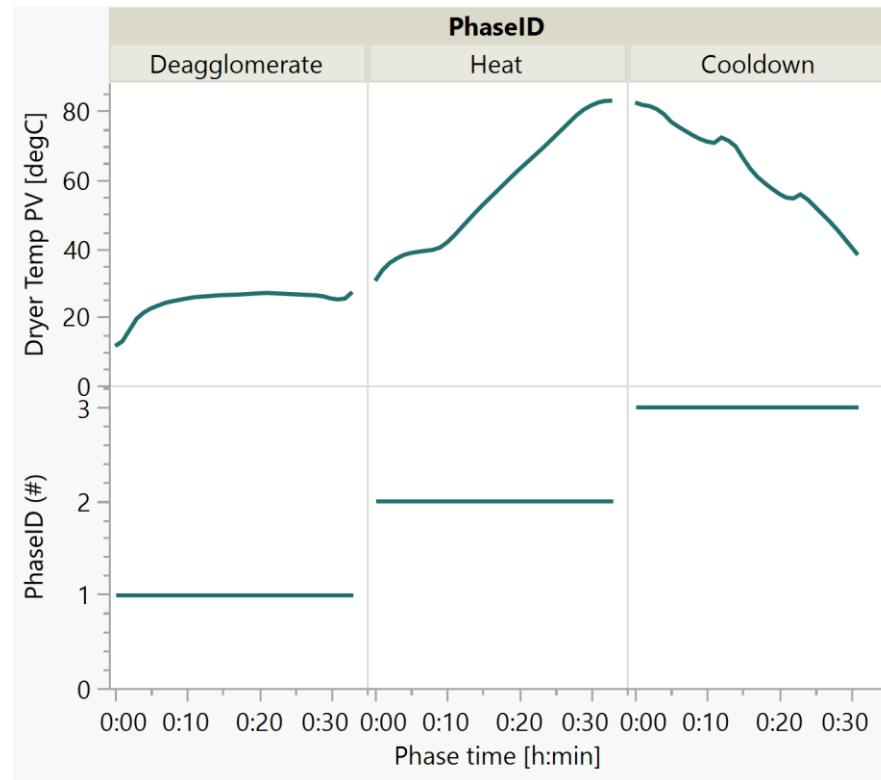
<https://www.imperial.ac.uk/process-systems-engineering/courses-and-seminars/workshops-and-courses/process-analytics-course/>

Monitoring and screening process data

Batch process monitoring



Batch data example – Industrial dryer



The dryer is charged with a variable amount of wet cake that evaporates and recollects a solvent material in an external tank. There are three distinct phases in the batch dataset:

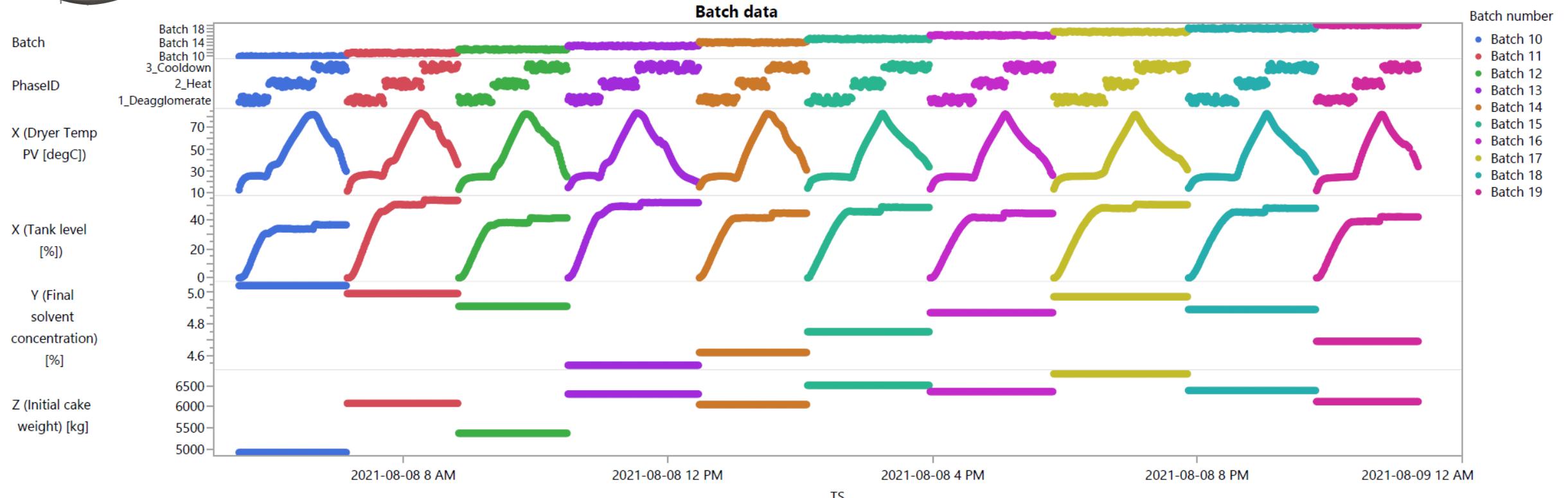
- **Deagglomeration** phase, where the cake reacts at low agitation speed while the solvent is collected
- **Heating** phase, where the temperature of the cake is increased until its set-point
- **Cooling** phase, where the batch temperature is reduced before unloading

Dataset and original paper:
Troubleshooting of an Industrial Batch Process Using
Multivariate Methods, Salvador Garcia et al., 2003
[Ind. Eng. Chem. Res. 2003, 42, 15, 3592–3601](https://doi.org/10.1002/anie.2003015392)





Batch data – time series view

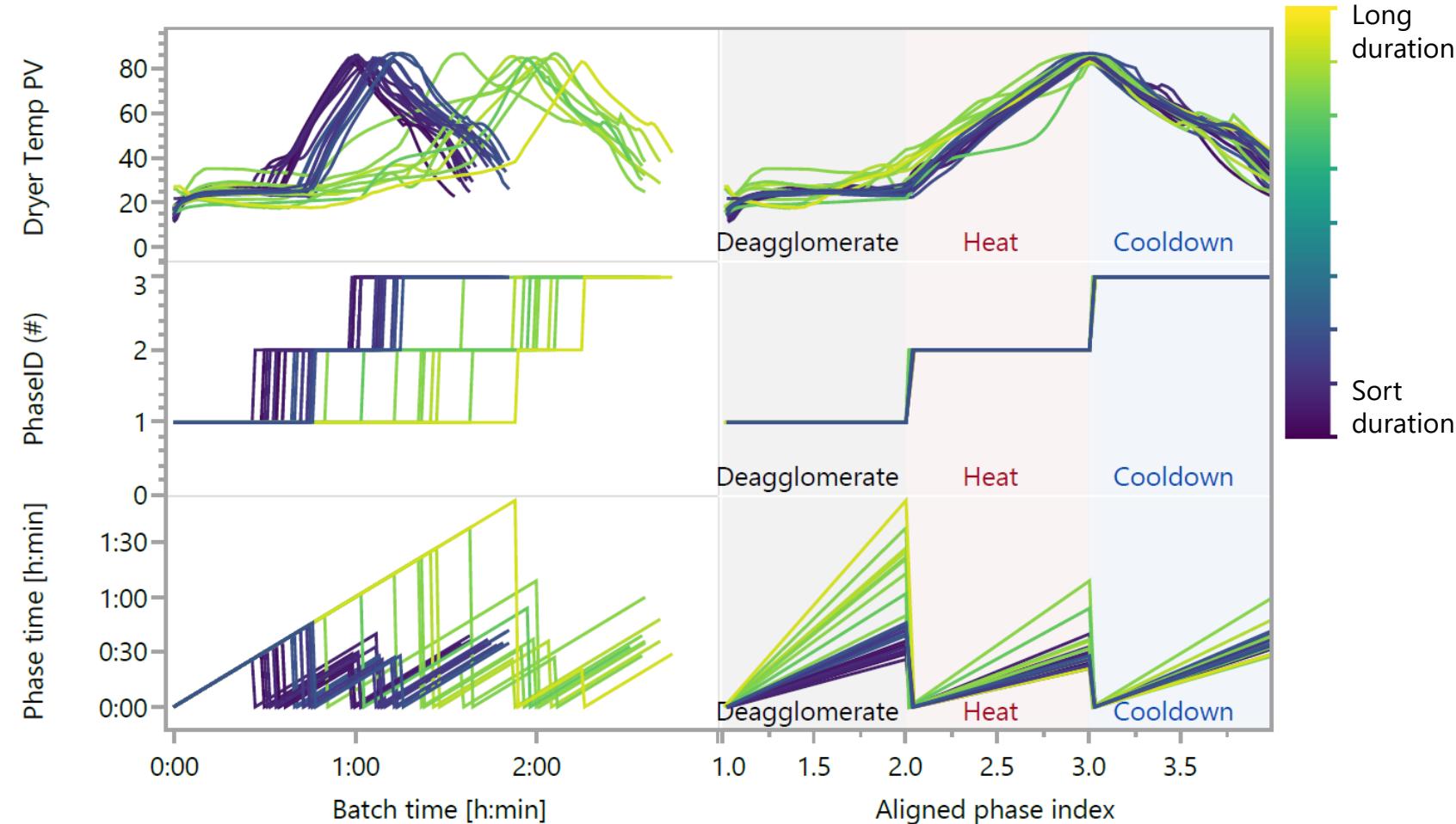


Batch processes have a combination of different data types. Event data which determines in which phase and in what batch the manufacturing process was. Time-varying data coming from sensors such as reactor temperature (c) which evolves through the process showing trends or trajectories (defined as inputs or X variables). Single measured values such as concentration, quantity, or quality for end products (named usually as target or Y variables) and initial conditions (defined as Z variables) that include raw materials properties.

Overlaying and aligning batch data

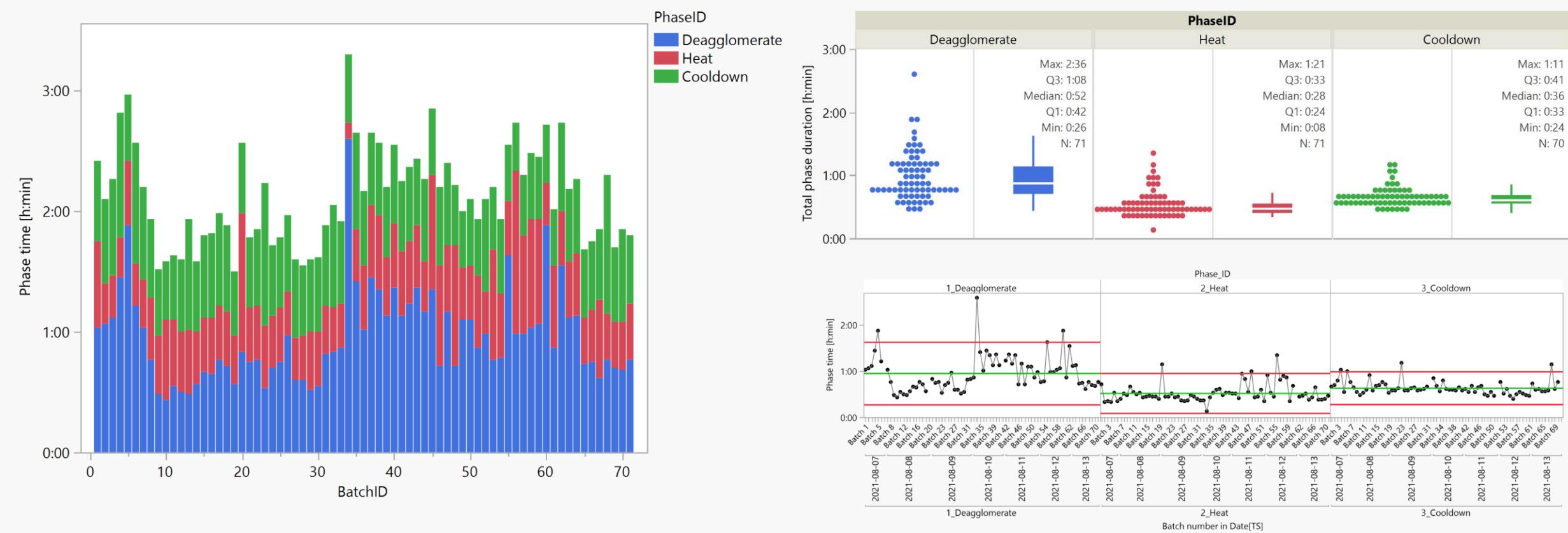
Overlay of the drying temperature profile and for a selection of historized batches colored by total batch duration.

Batches can be aligned using the automation data (phases or stages).



Source:
Industrial Data Science for Batch Manufacturing Processes
[arXiv:2209.09660](https://arxiv.org/abs/2209.09660)

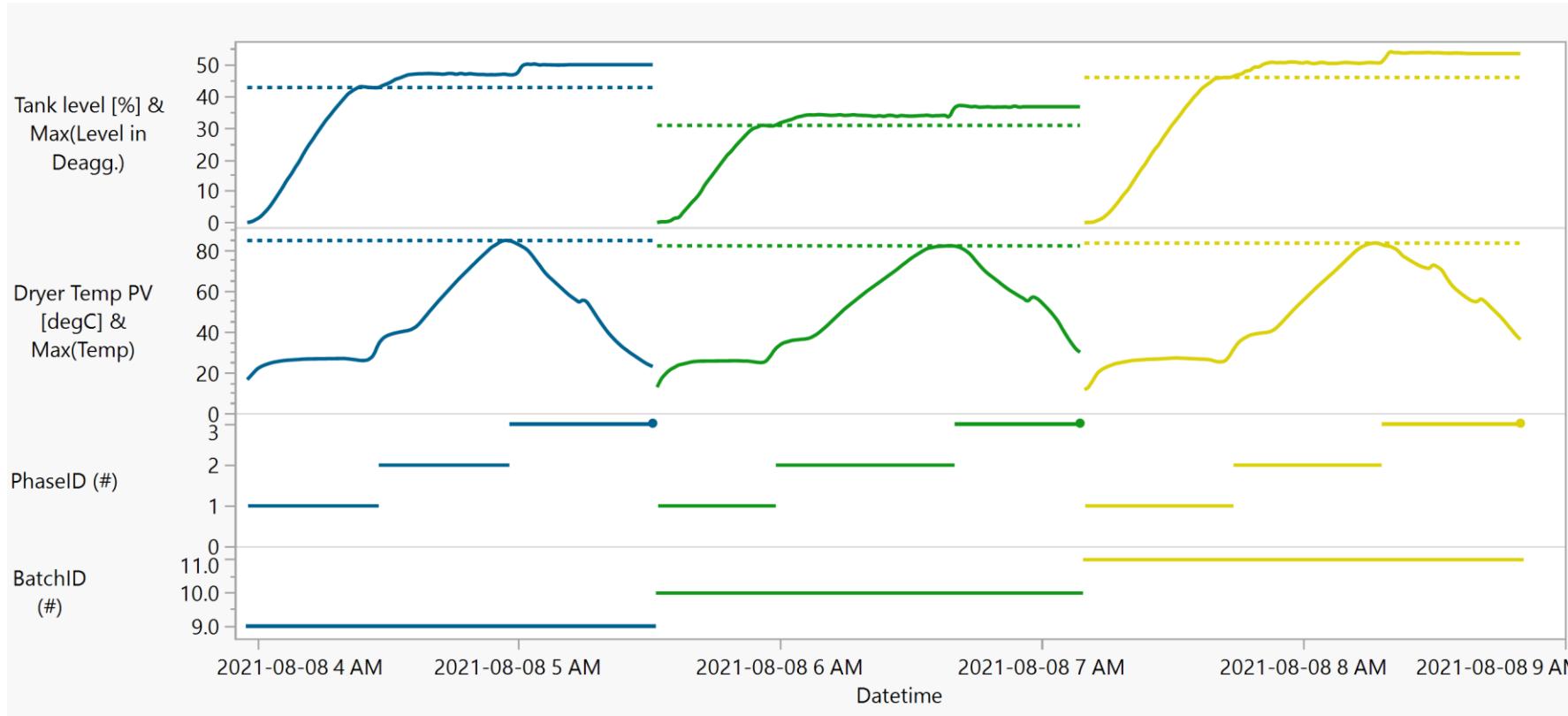
Time information is not lost after alignment



- Batch duration is one of the most important KPIs as it is directly related to productivity. Bar plots, distributions and control charts are often used to detect anomalies.

Source:
 Industrial Data Science for Batch Manufacturing Processes
[arXiv:2209.09660](https://arxiv.org/abs/2209.09660)

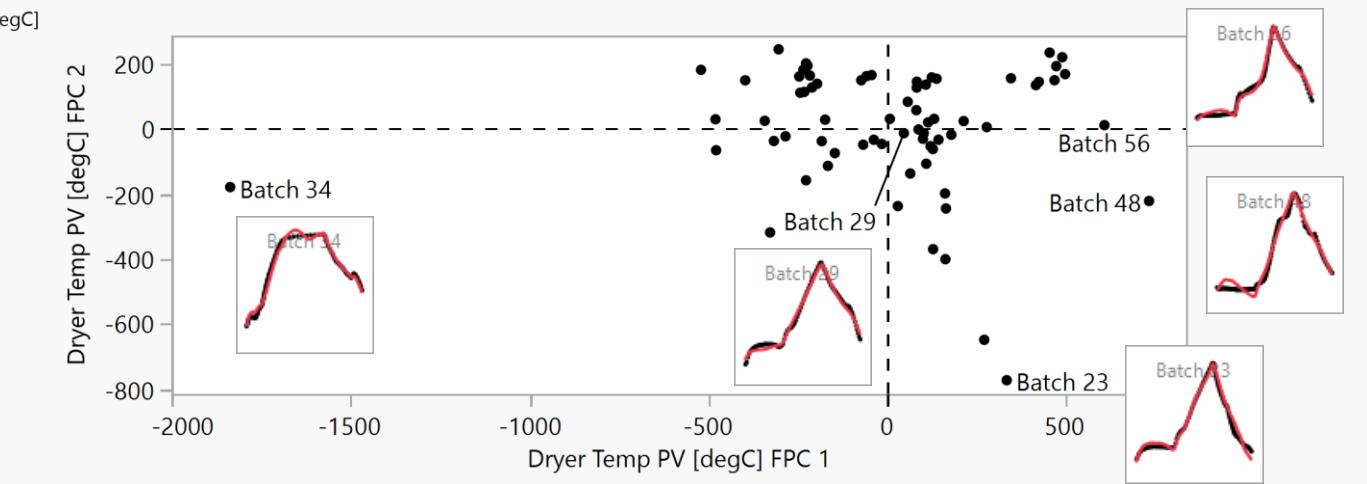
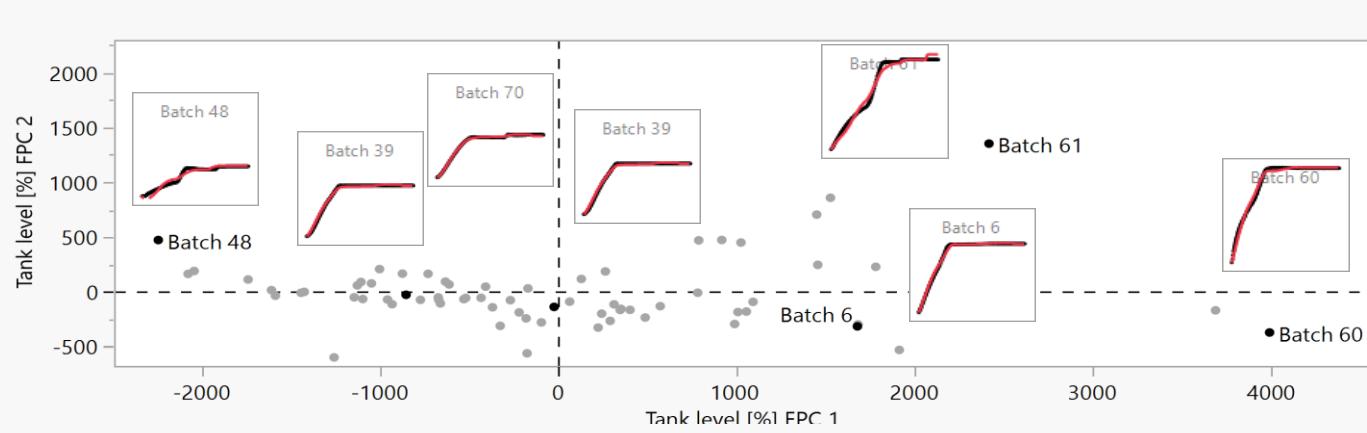
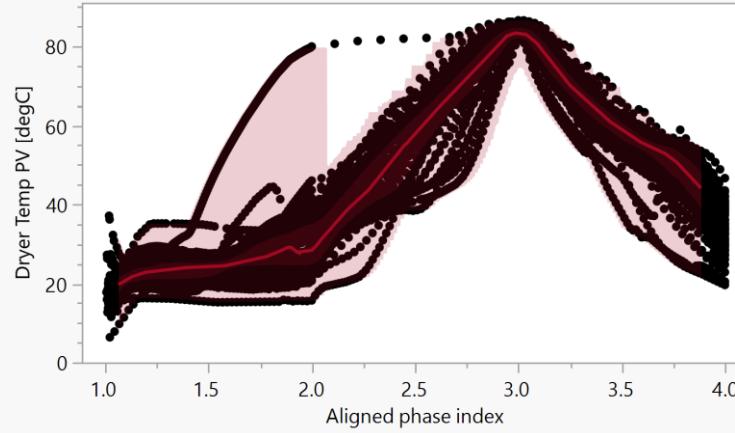
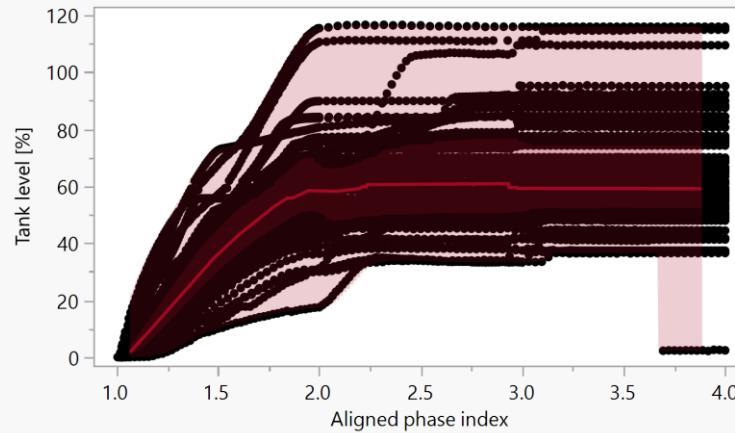
Batch monitoring can be done summarizing information into KPIs



Process and control knowledge is essential to summarize relevant batch information into useful batch statistics. These are often called KPIs, landmarks or fingerprints.

Source:
Industrial Data Science for Batch Manufacturing Processes
[arXiv:2209.09660](https://arxiv.org/abs/2209.09660)

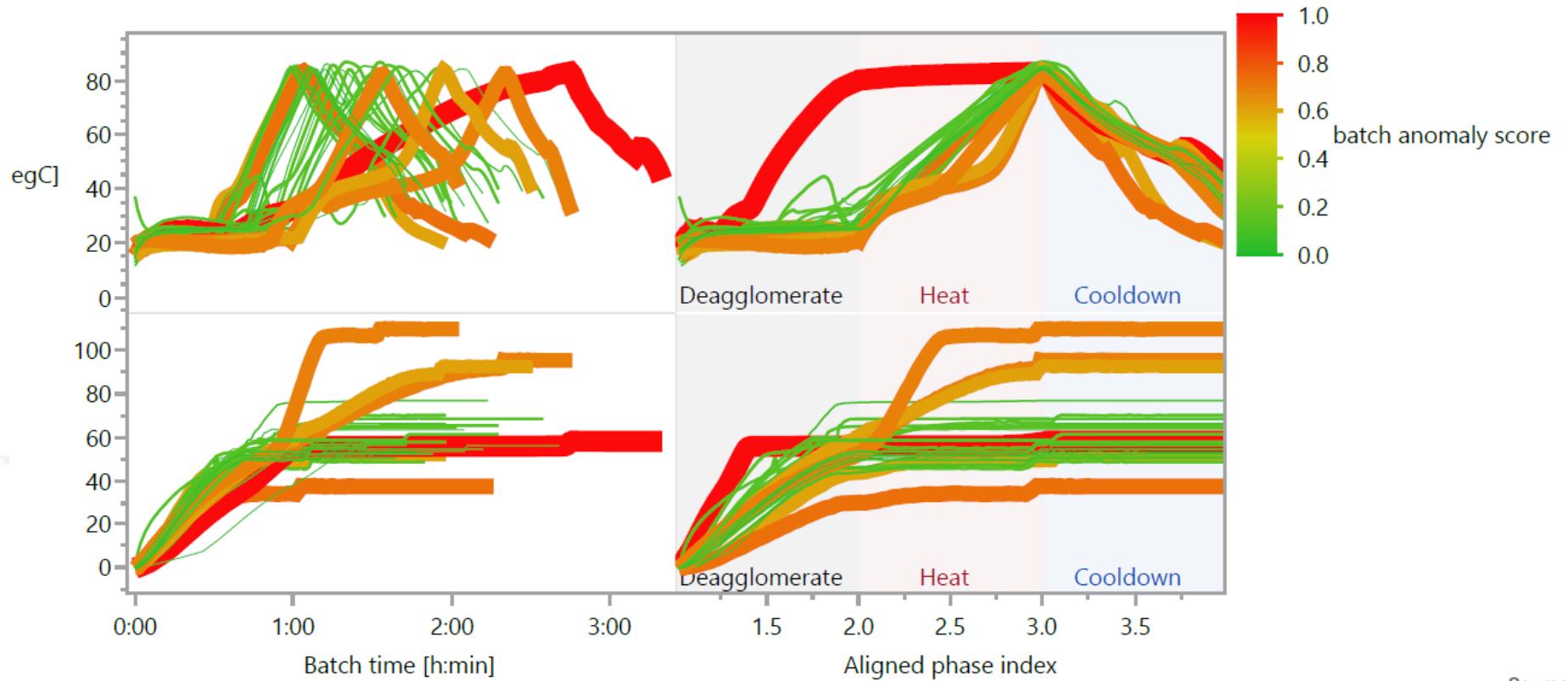
Functional PCA finds batch variability without defining KPIs



Without the need of asking for specific features, Functional PCA can capture and decompose the whole trajectories seen in batch processes.

Source:
Industrial Data Science for Batch Manufacturing Processes
[arXiv:2209.09660](https://arxiv.org/abs/2209.09660)

Application: Batch monitoring (unsupervised learning)



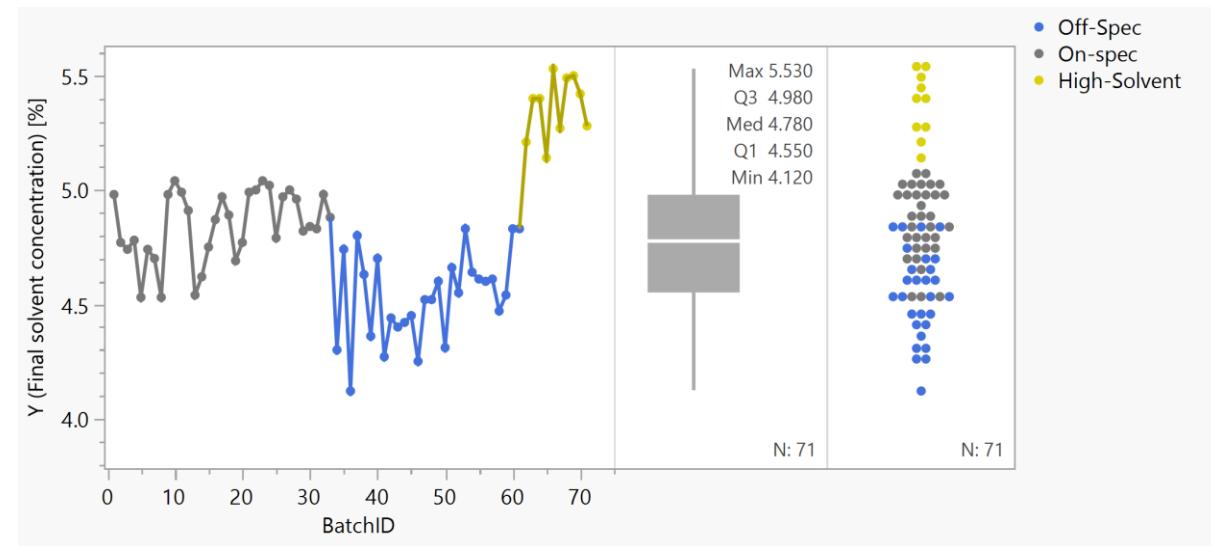
Source:
Industrial Data Science for
Batch Manufacturing
Processes
[arXiv:2209.09660](https://arxiv.org/abs/2209.09660)

Monitoring and screening process data

Batch process screening



Screening of batch processes using process engineering



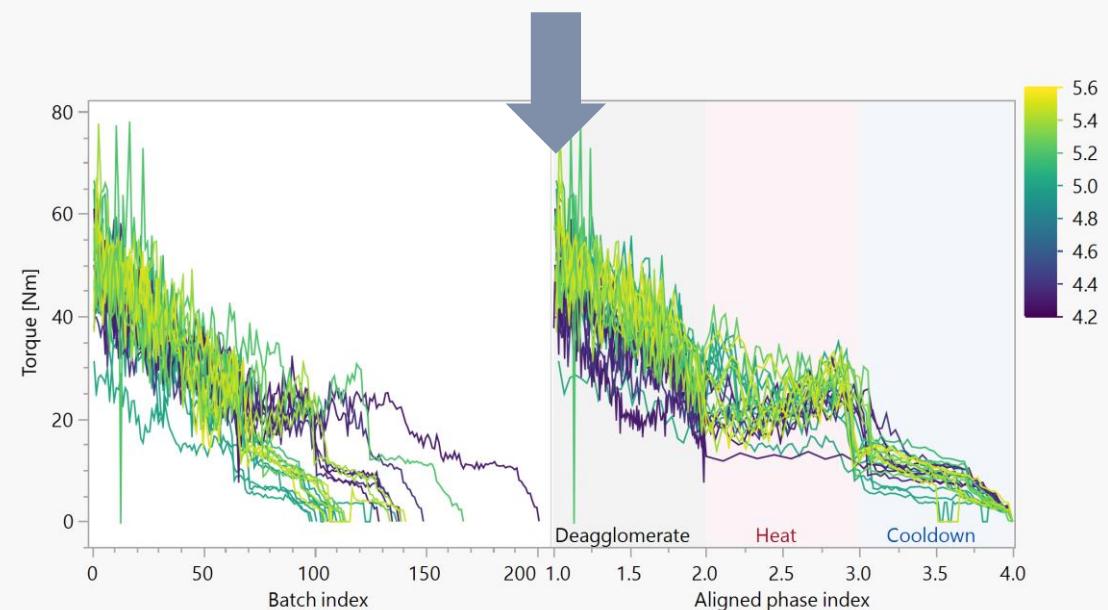
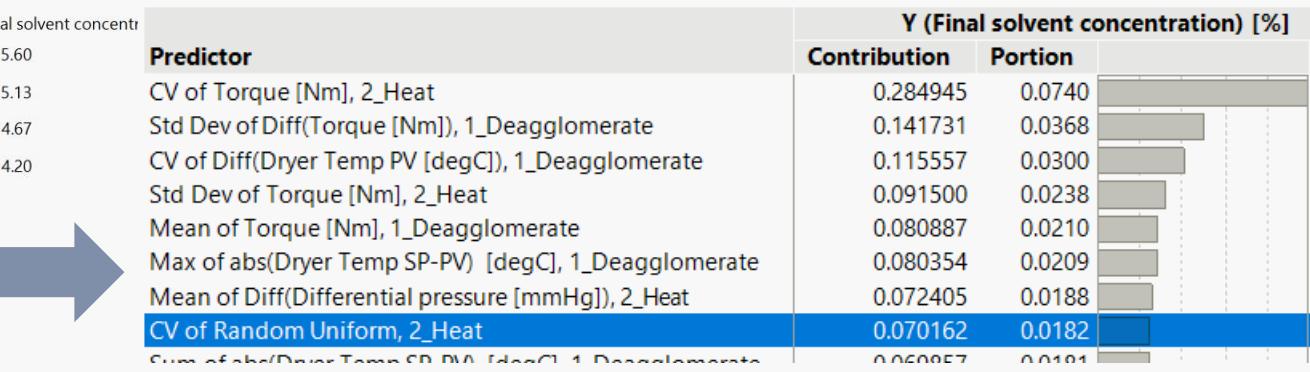
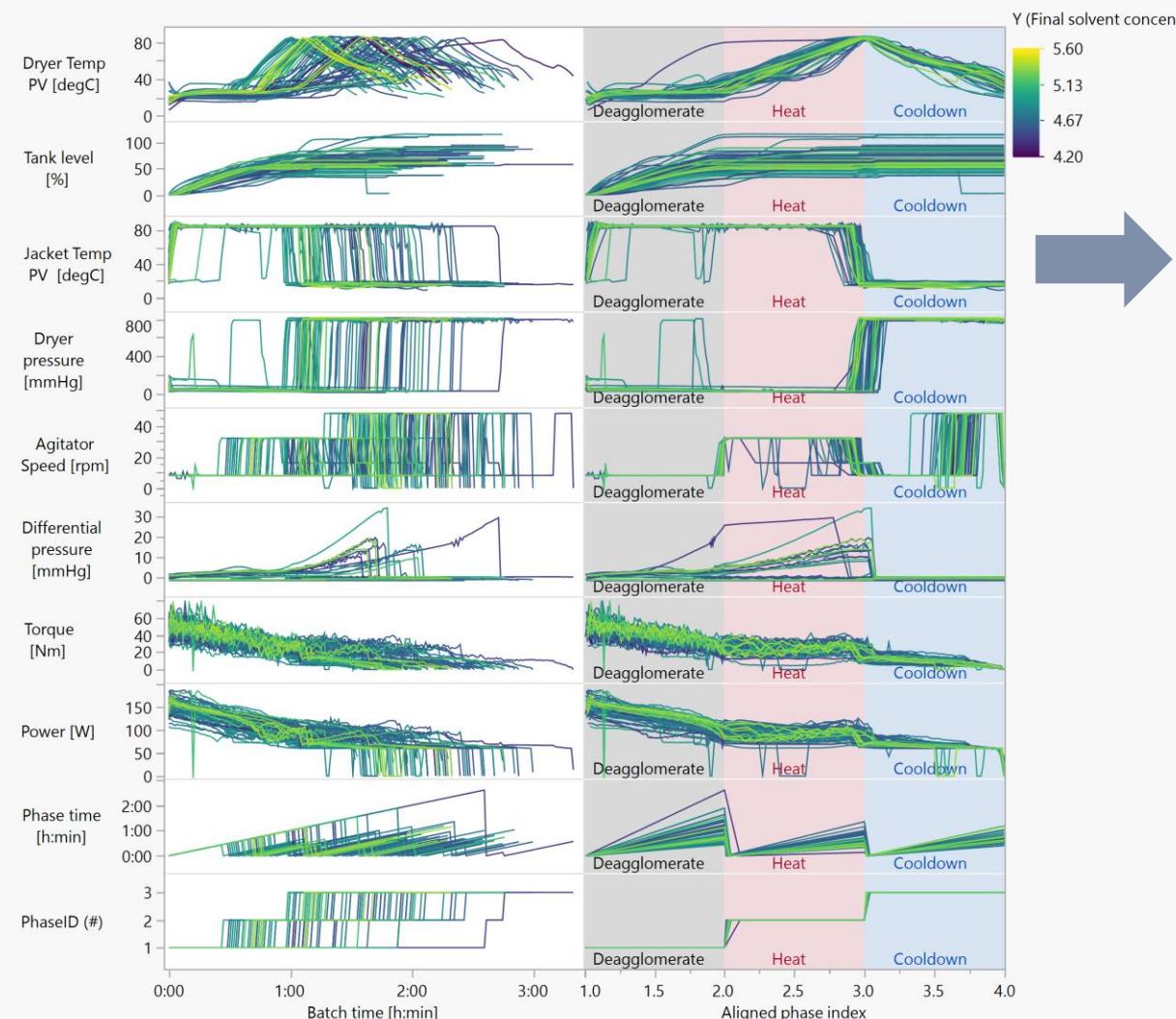
It is tempting to ask engineers for important KPIs that can explain solvent concentration, so a correlation analysis (supervised learning) can be done. Max pressures, temperatures rates, agitation power, jacket temperature, etc.

Do we need to?

Source:
Industrial Data Science for
Batch Manufacturing
Processes
[arXiv:2209.09660](https://arxiv.org/abs/2209.09660)

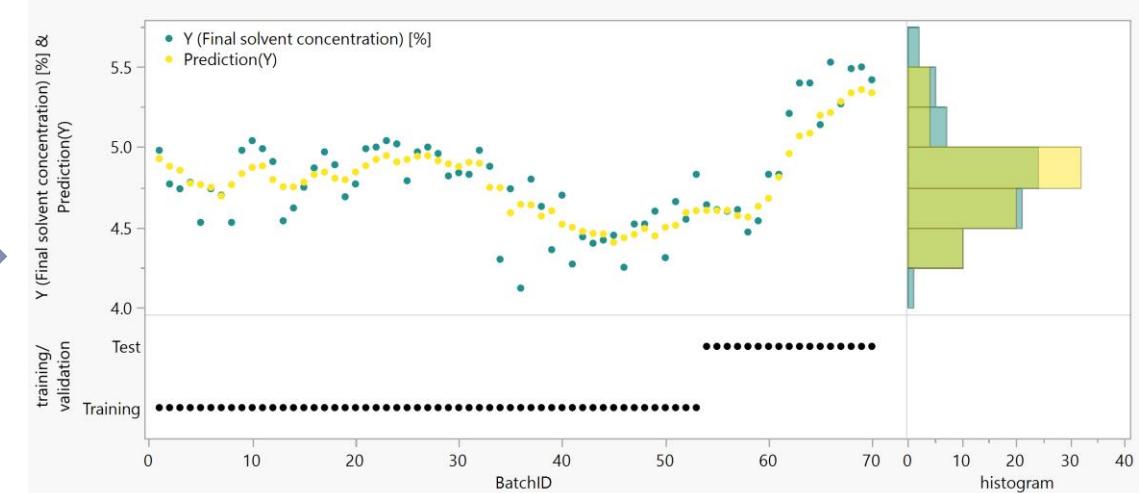
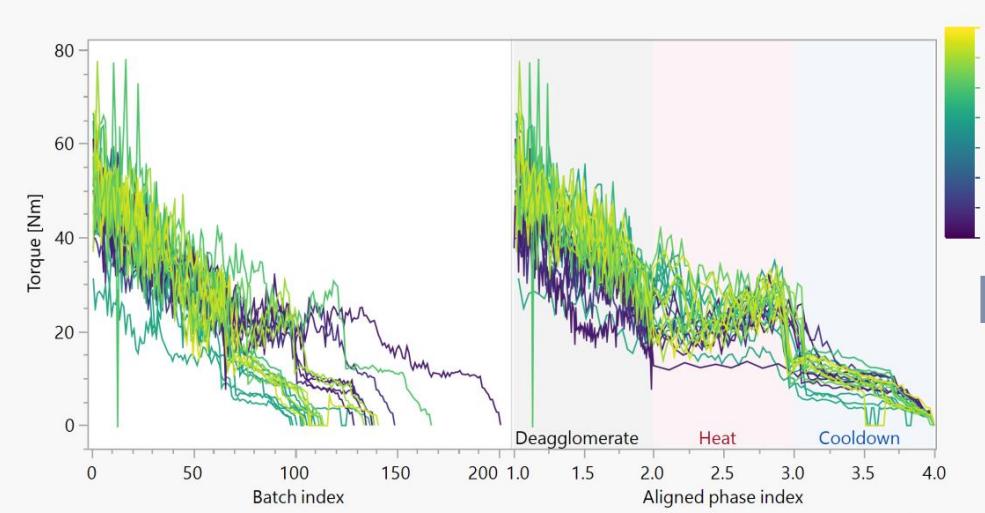
Screening of batch processes using feature engineering and variable selection

Source:
 Industrial Data Science for Batch Manufacturing
 Processes
[arXiv:2209.09660](https://arxiv.org/abs/2209.09660)



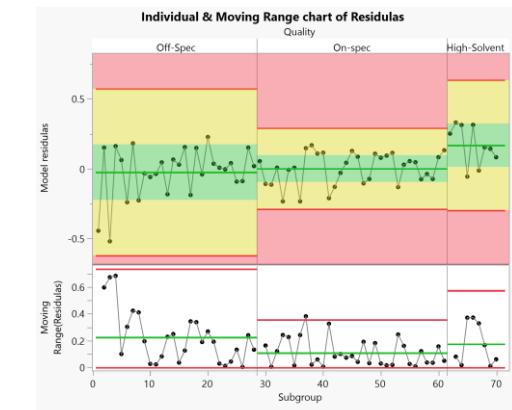
Where((Y (Final solvent concentration) [%]"n < 4.4389 | "Y (Final solvent concentration) [%]"n > 4.9845) and ("Y (Final solvent concentration) [%]"2"n >= 4.2627))

Application: Viscosity/concentration prediction



Agitation power correlates to solvent concentration due to changes of viscosity. In the literature, power consumption from re-circulating or rotary pumps and agitators are commonly used to predict quality measurements before they arrive from the laboratory.

Monitoring of residuals and adaptive models are often required to keep a model accuracy after deployment.

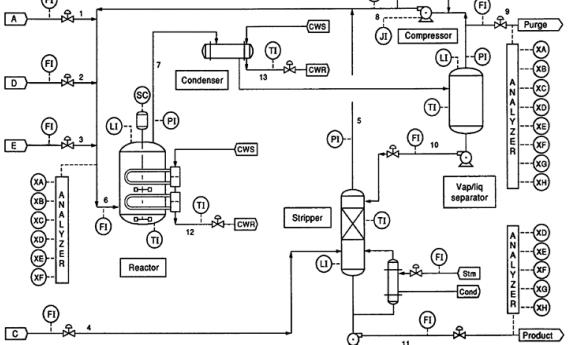


Monitoring and screening process data

Hands-on

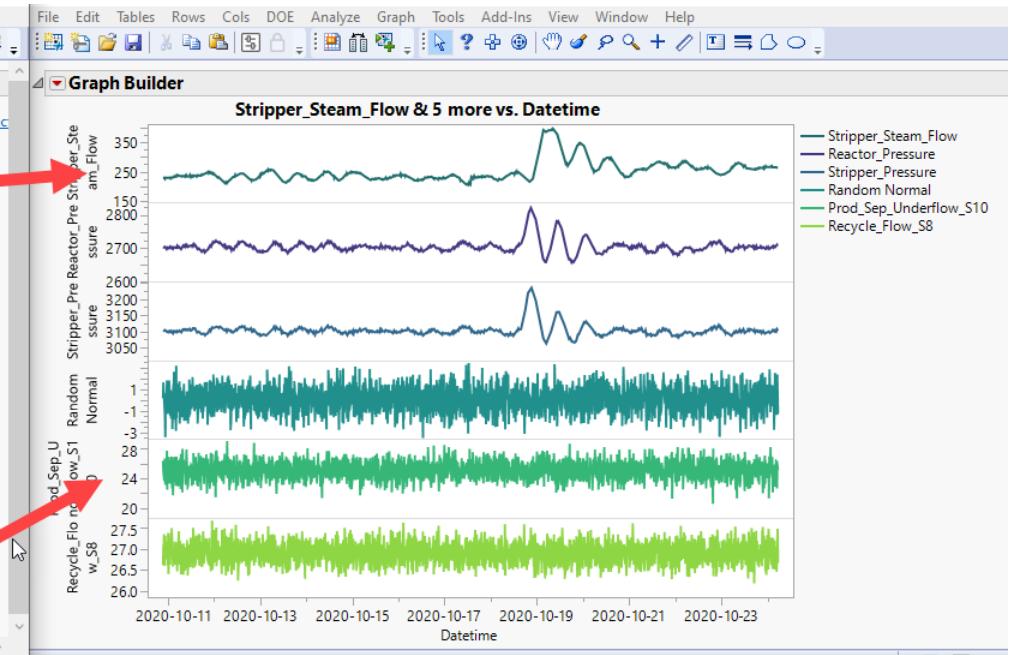


Predictor Explainer: Unsupervised monitoring



Predictor Screening

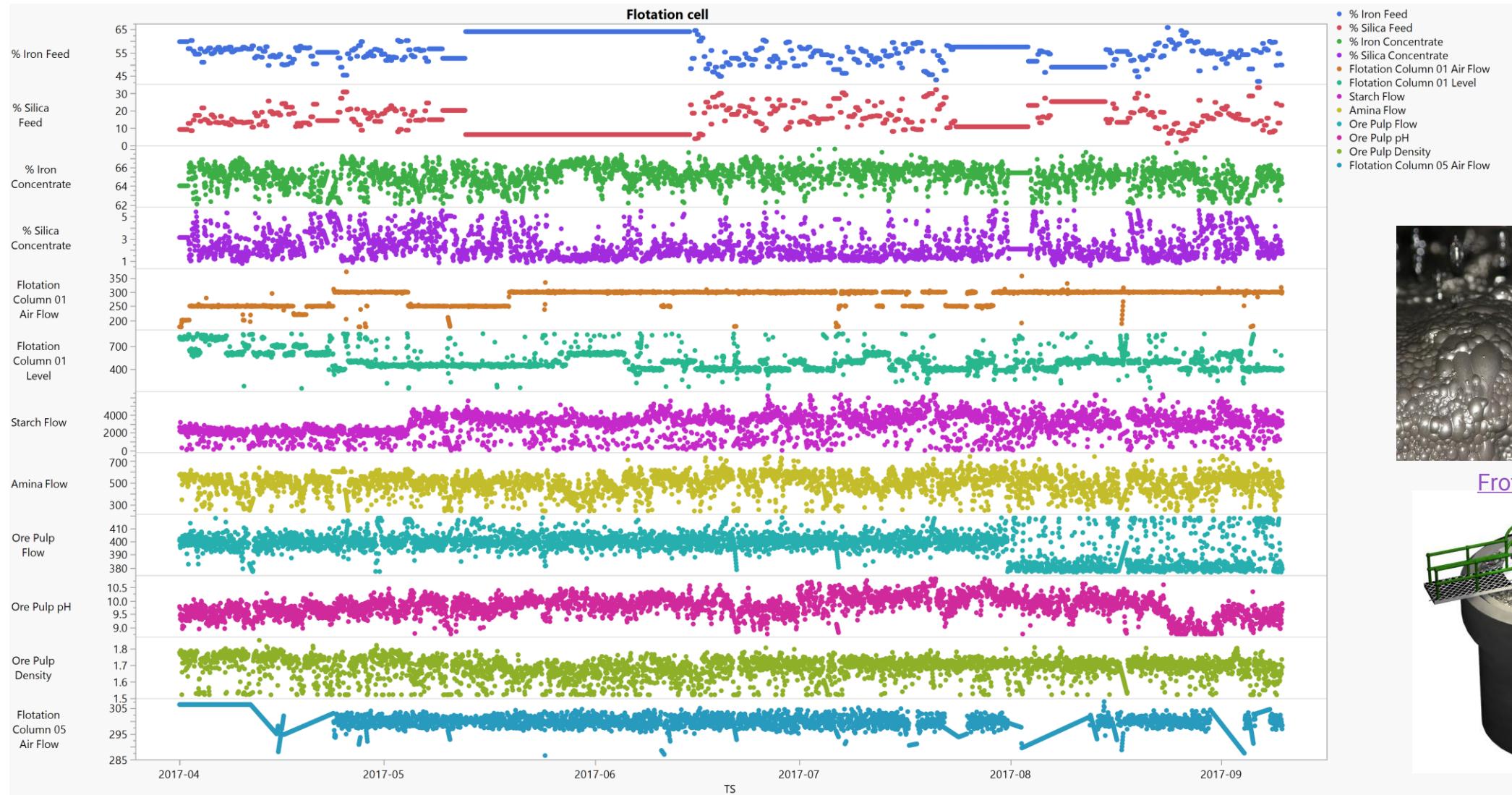
Predictor	Contribution	Portion	Rank
Stripper_Pressure	31526508	0.4082	1
Reactor_Pressure	9933557.3	0.1286	2
Stripper_Steam_Flow	8123483.2	0.1052	3
Product_Sep_Temp	7236590.9	0.0937	4
Stripper_Temp	4938697.6	0.0639	5
Prod_Sep_Pressure	4791168.3	0.0620	6
Compressor_Work	4717225.1	0.0611	7
Purge_Rate_S9	2131639.1	0.0276	8
Radiator_Cooling_Water_Outlet_Temp	1027985.6	0.0133	9
A_Feed_S1	954638.96	0.0124	10
Total_Feed_S4	664358.38	0.0086	11
Separator_Cooling_Water_Outlet_Temp	363415.55	0.0047	12
Radiator_Level	132093.83	0.0017	13
Stripper_Level	126919.18	0.0016	14
Stripper_Underflow_S11	114667.12	0.0015	15
Product_Sep_Level	85844.932	0.0011	16
E_Feed_S3	83116.678	0.0011	17
Radiator_Temperature	71746.005	0.0009	18
D_Feed_S2	60417.116	0.0008	19
Radiator_Feed_Rate_S6	50217.317	0.0007	20
Random Normal	352133.75	0.0005	21
Prod_Sep_Underflow_S10	30449.878	0.0004	22
Recycle_Flow_S8	30284.243	0.0004	23



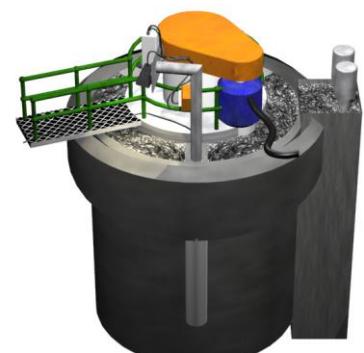
- Frontend: JMP
- Backend (optional): Python

<https://github.com/industrial-data/predictor-explainer>

Process data from mining industry

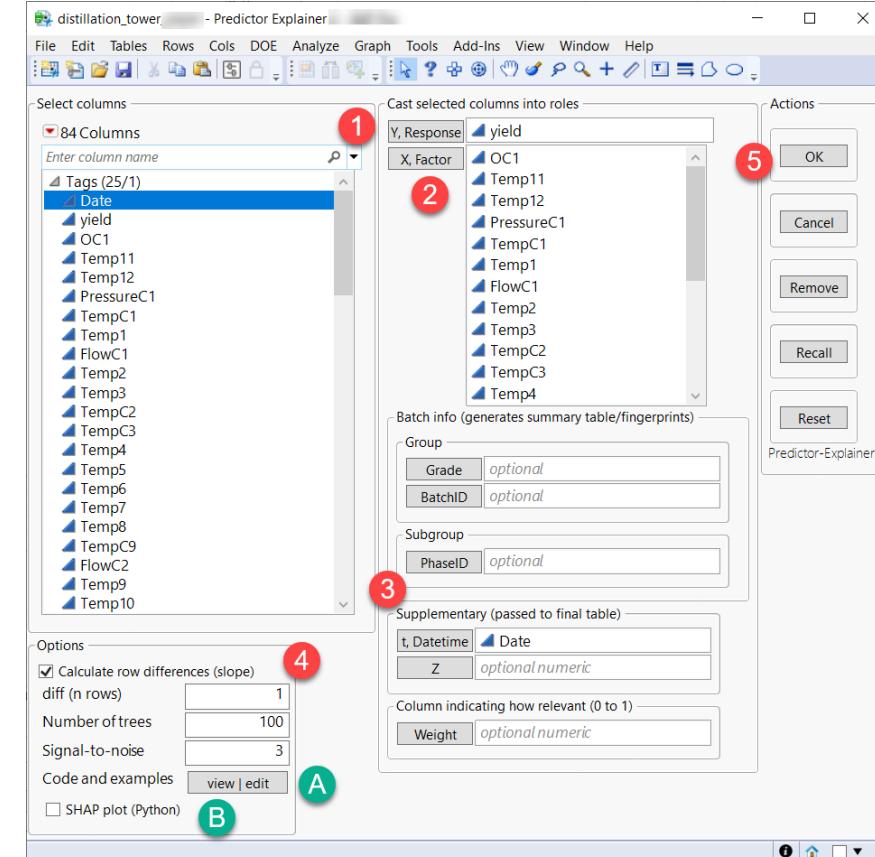
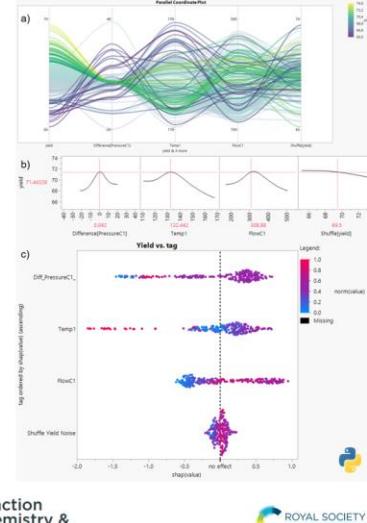
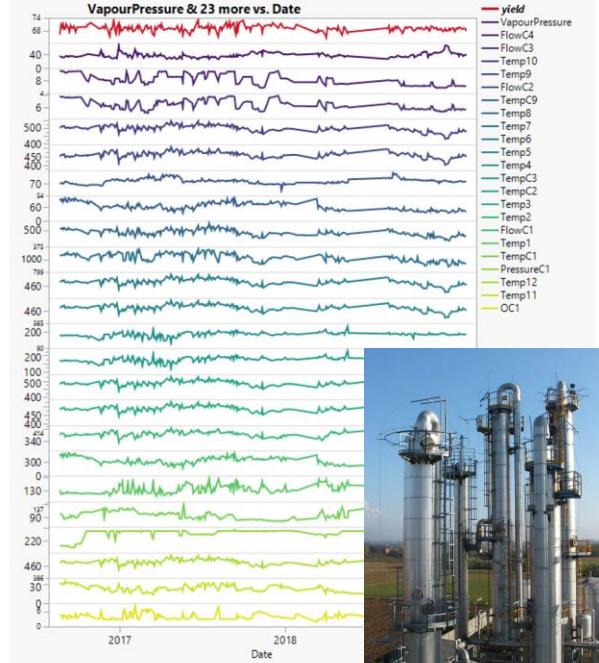


Froth flotation



<https://www.kaggle.com/datasets/edumagalhaes/quality-prediction-in-a-mining-process>

AutoML and ExplainableAI applied to process industries



- Frontend: JMP
- Backend (optional): Python

<https://github.com/industrial-data/predictor-explainer>



Similar project in pure Python:
<https://explainerdashboard.readthedocs.io/>

Summary

- In machine learning, something that doesn't change is irrelevant, both for anomaly and screening models.
Example: a controlled reactor temperature won't appear in models as its variation will be minimal while measured disturbances and manipulated variables will.
- System dynamics are always present in process data, several methods exist to remove unstable periods (see BASF example in Day 3)
- Anomaly detection can be used to identify major process challenges
 - Univariate control charts
 - Multivariate control charts using Hotelling's T^2 , SPE and KNN distances
- Batch alignment using automation triggers allows simpler batch anomaly detection
- Screening of batch processes can be automated using summary statistics and screening models with synthetic noise for variable selection
- Online estimation of quality parameters are a common application of batch predictive control (e.g., estimating viscosity as a function of agitation power). These are often called inferential sensors or software sensors and require continuous adaptation (re-calibration).



*All models are wrong, but some
are useful*

George Box

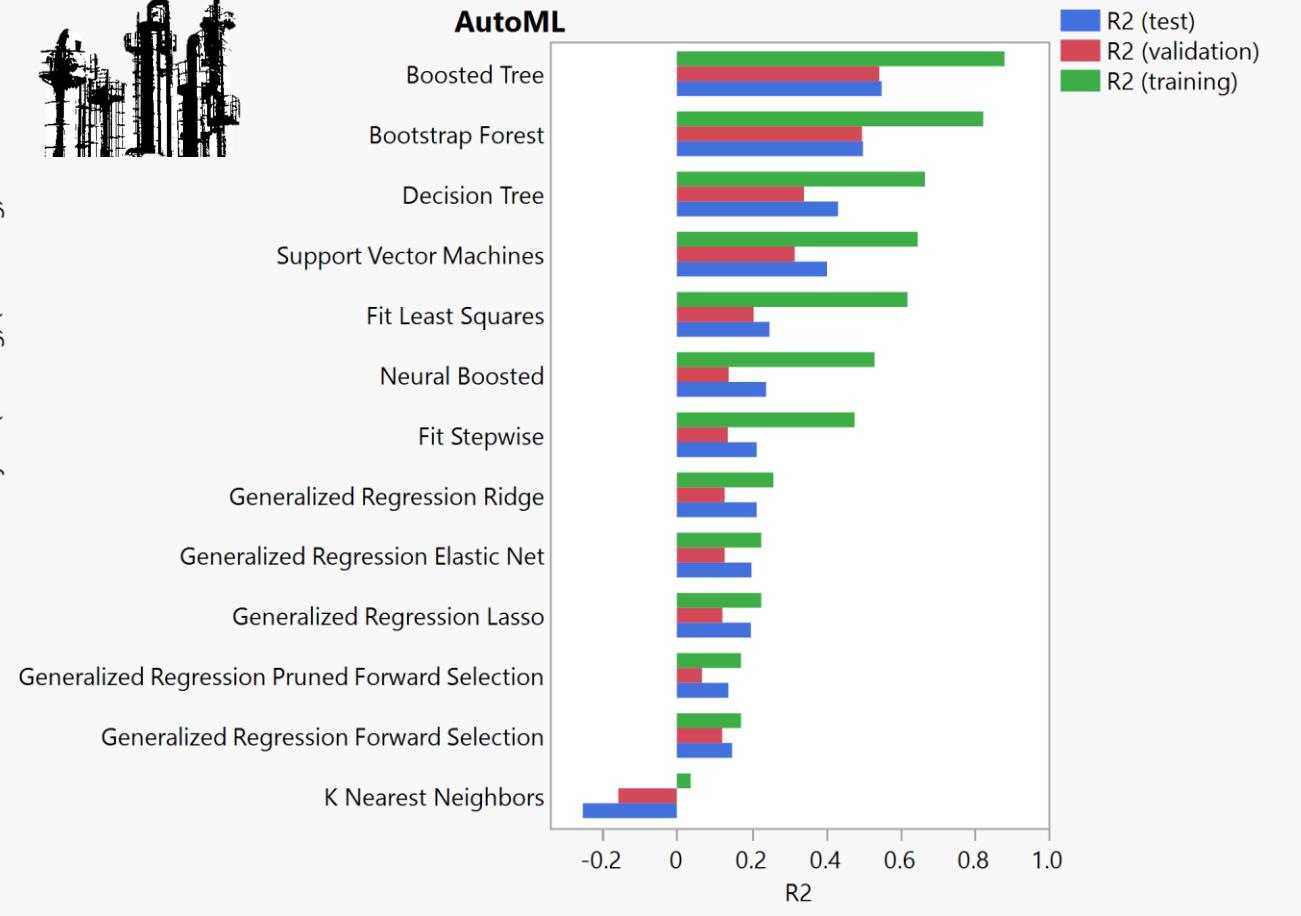
Modeling and understanding



AutoML, trying all models and tuning parameters at once



Method ordered by R2 (training) (ascending)



JMP Pro Model Screening

<https://pycaret.org/>



Data Preparation



Model Training



Hyperparameter Tuning



Analysis & Interpretability



Model Selection



Experiment Logging

```
● ● ●
# Time Series Forecasting Functional API Example
# loading sample dataset
from pycaret.datasets import get_data
data = get_data('airline')

# init setup
from pycaret.time_series import *
s = setup(data, fh = 3, session_id = 123)

# model training and selection
best = compare_models()

# plot trained model
plot_model(best)

# predict on hold-out/test set
pred_holdout = predict_model(best)

# predict in unseen future
predictions = predict_model(best, fh = 36)

# save model
save_model(best, 'best_pipeline')
```

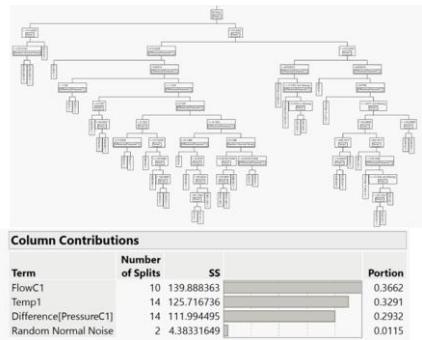
Careful optimizing the wrong problem

*Possibly the most common error of a smart engineer
is to optimize a thing that should not exist*

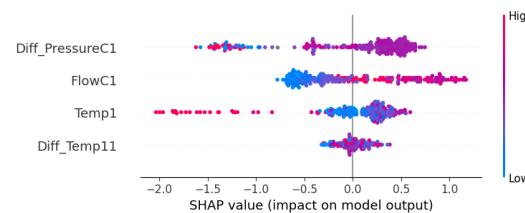
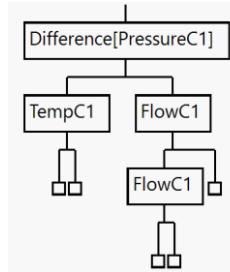
Elon Musk

Modeling steps in practice

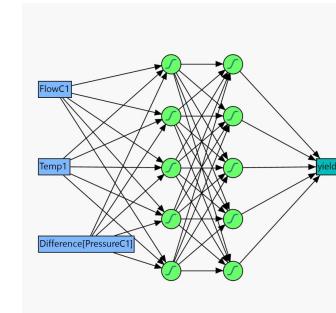
Variable selection



Understanding

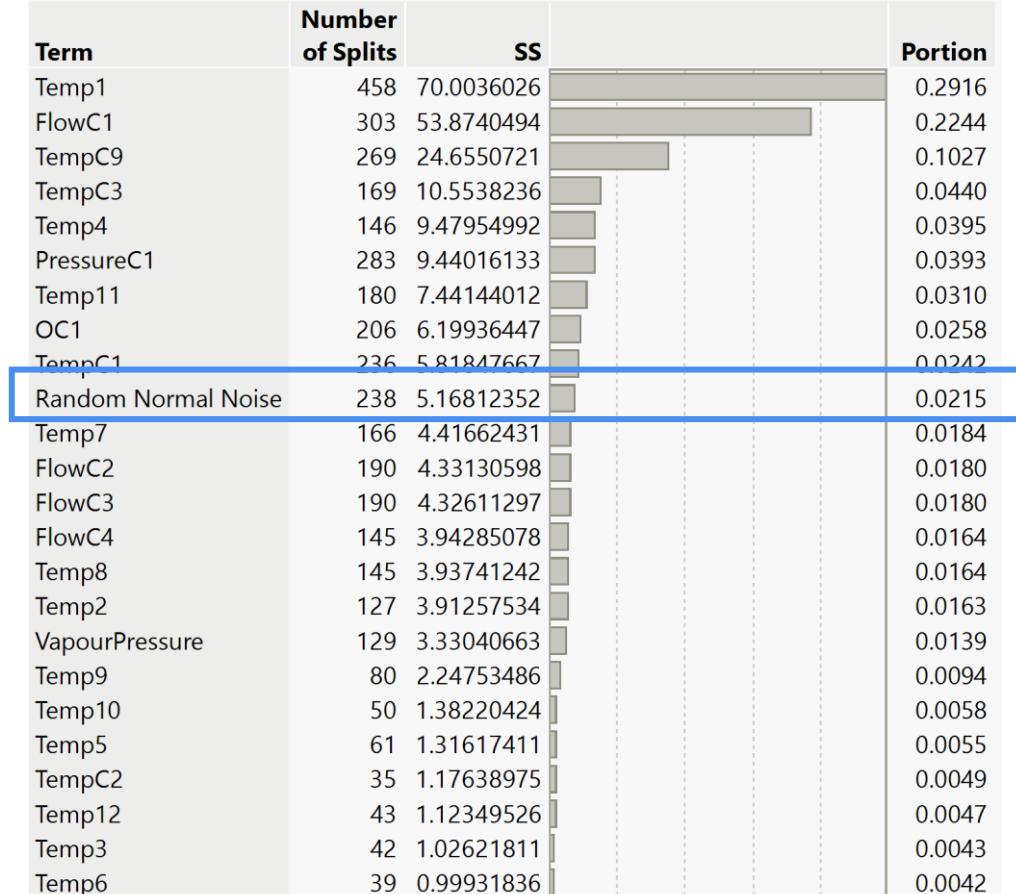
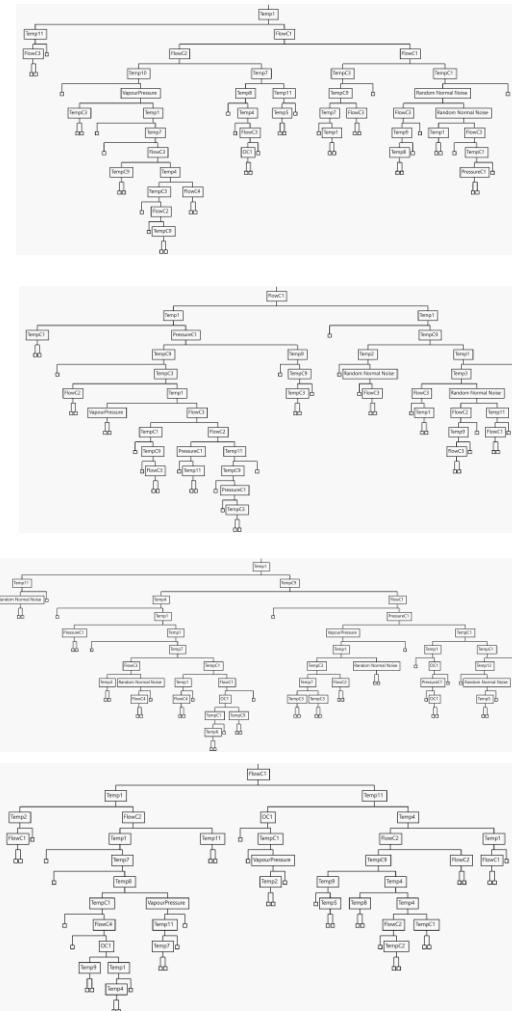


Prediction (if needed)

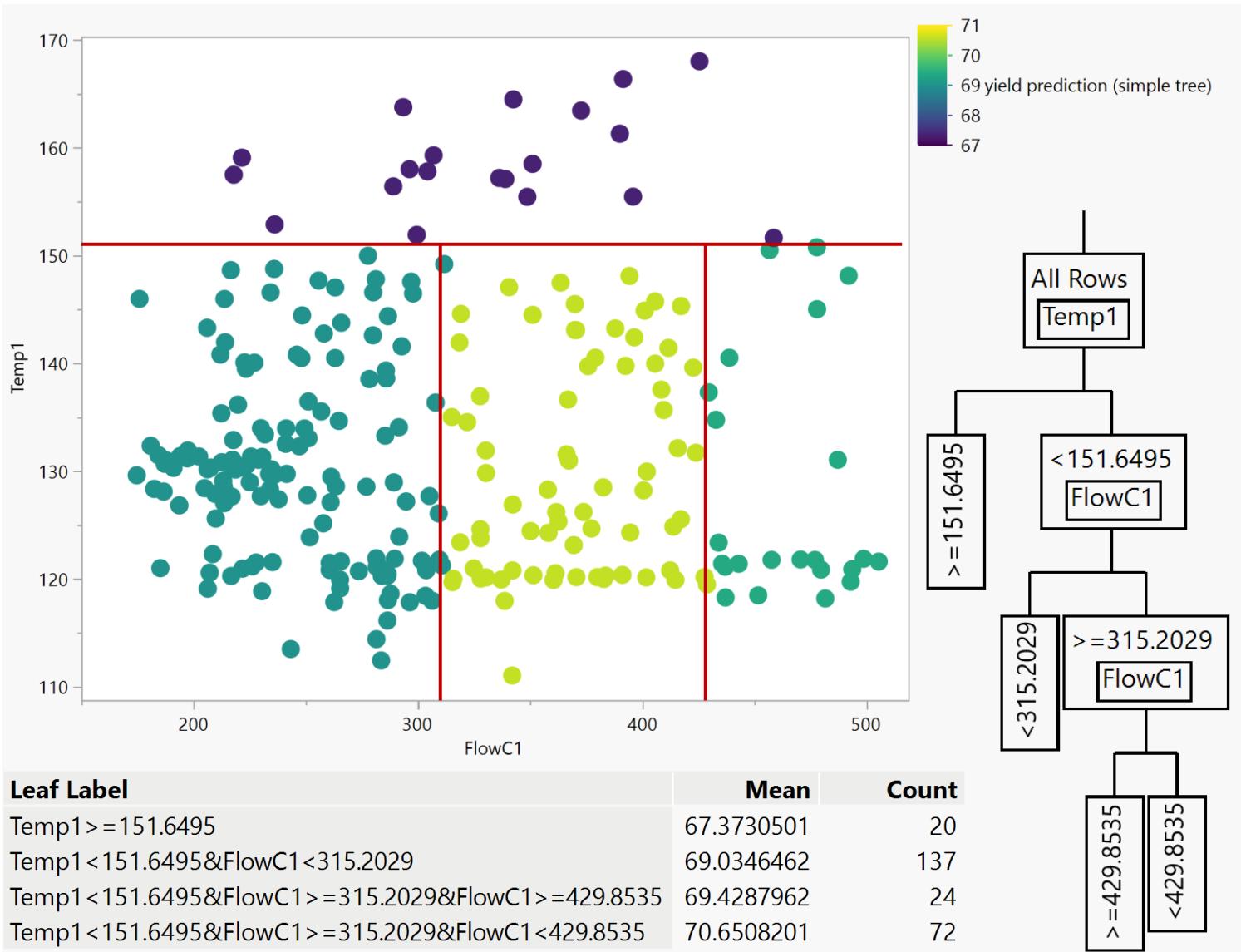


Variable selection

First, a **screening of variables** and selection of tags (sensors) using **tree-based models**. Many tags will end up being weakly correlated to the target, perhaps trying to explain noise in the yield measurement. By adding known noise additional tags, the selection of variables to keep in the model with a certain contribution is facilitated.

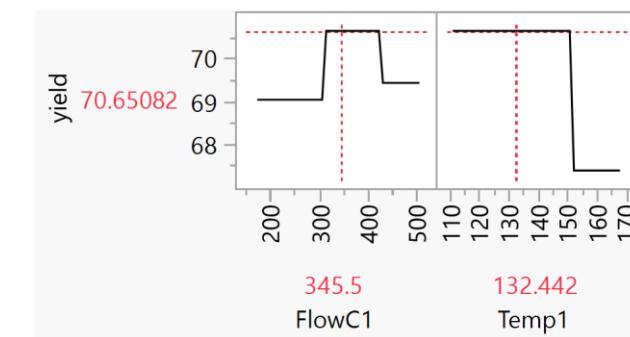


Partition trees, creating expert systems from data



Tree-based models create simple if-then logics via data partitions that can better explain the target.

They are common for screening and model interpretation, as they can handle tags with different units, the presence of missing values, and outliers while uncovering non-linear relationships.

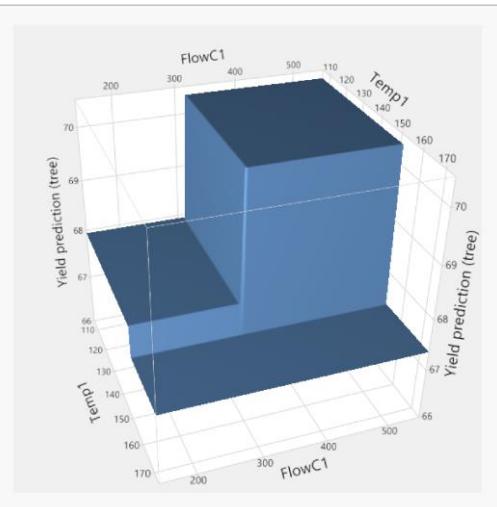
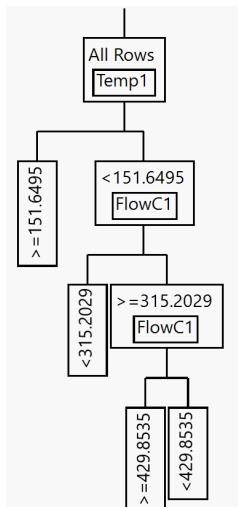


[Industrial data science – a review of machine learning applications for chemical and process industries](#)

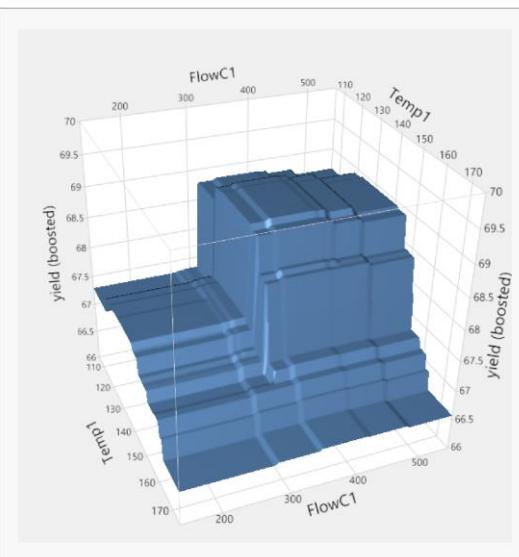
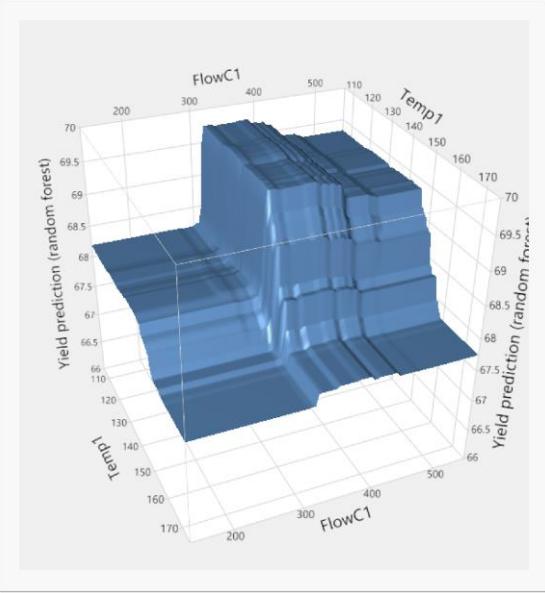
Tree-based models

Bootstrap tree
(Random forest)

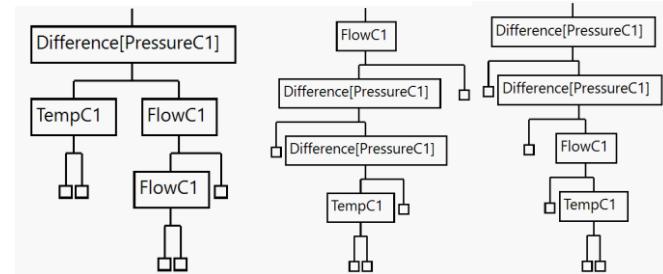
Partition tree



Boosted tree
(LightGBM, CatBoost,
XGBoost...)

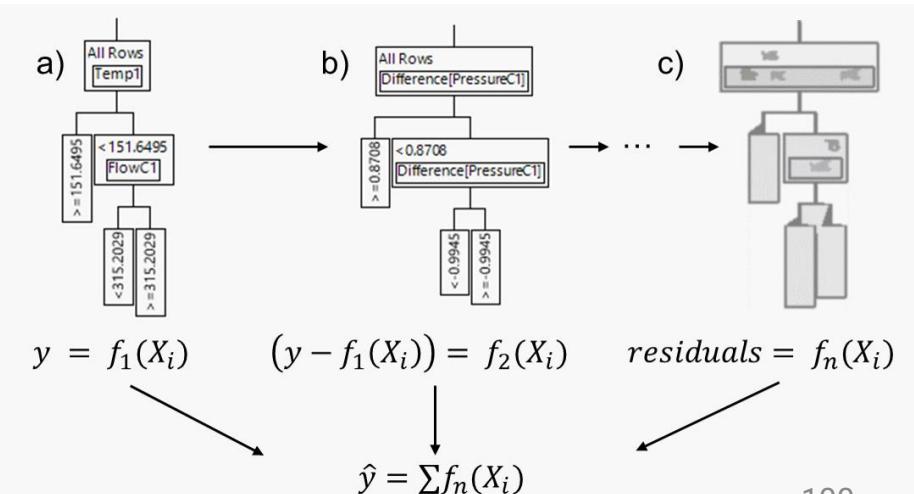


Average prediction of multiple trees after sampling



$$\hat{y} = \frac{\sum f_n(X_i)}{n}$$

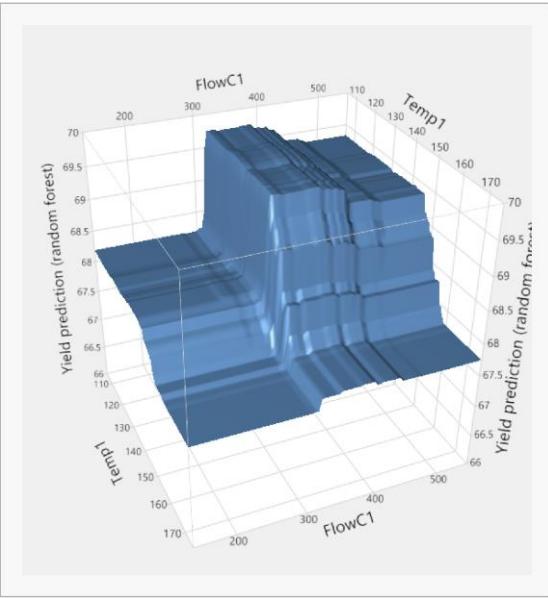
Sum of subsequent predictions



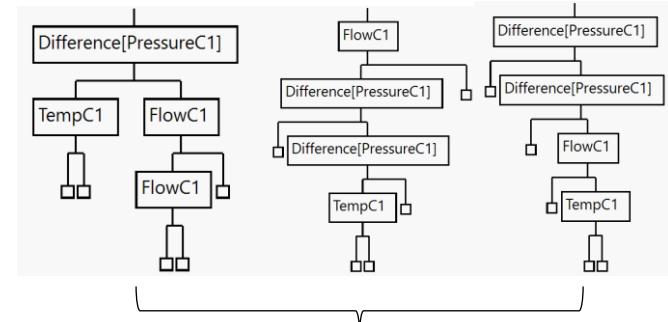
Tree-based models

Bootstrap tree (Random forest)

- Faster
- Robust, resampling and averages
- **Overfits**



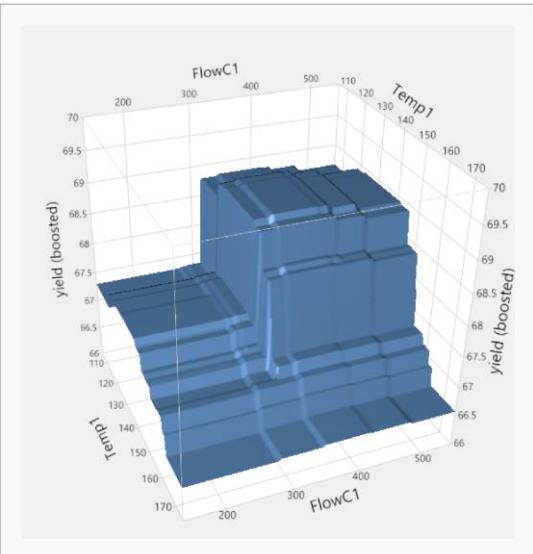
Average prediction of multiple trees after sampling (bagging)



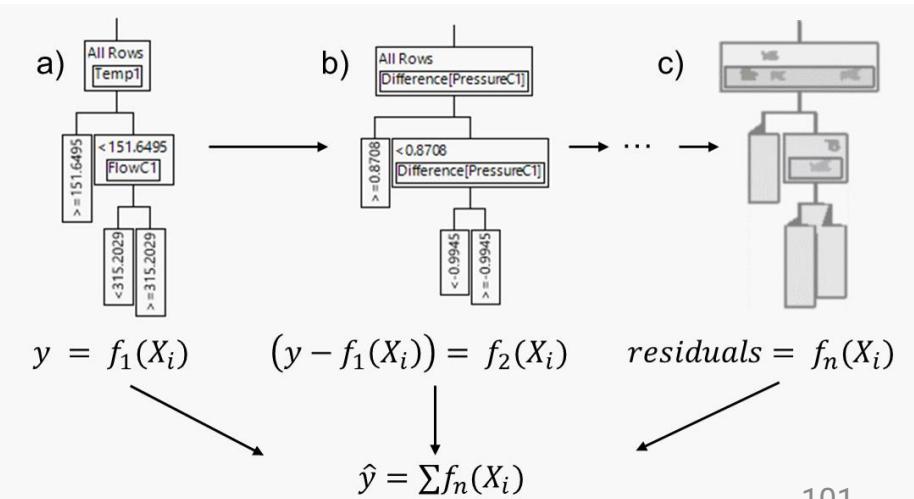
$$\hat{y} = \frac{\sum f_n(X_i)}{n}$$

Boosted tree (LightGBM, CatBoost, XGboost,)

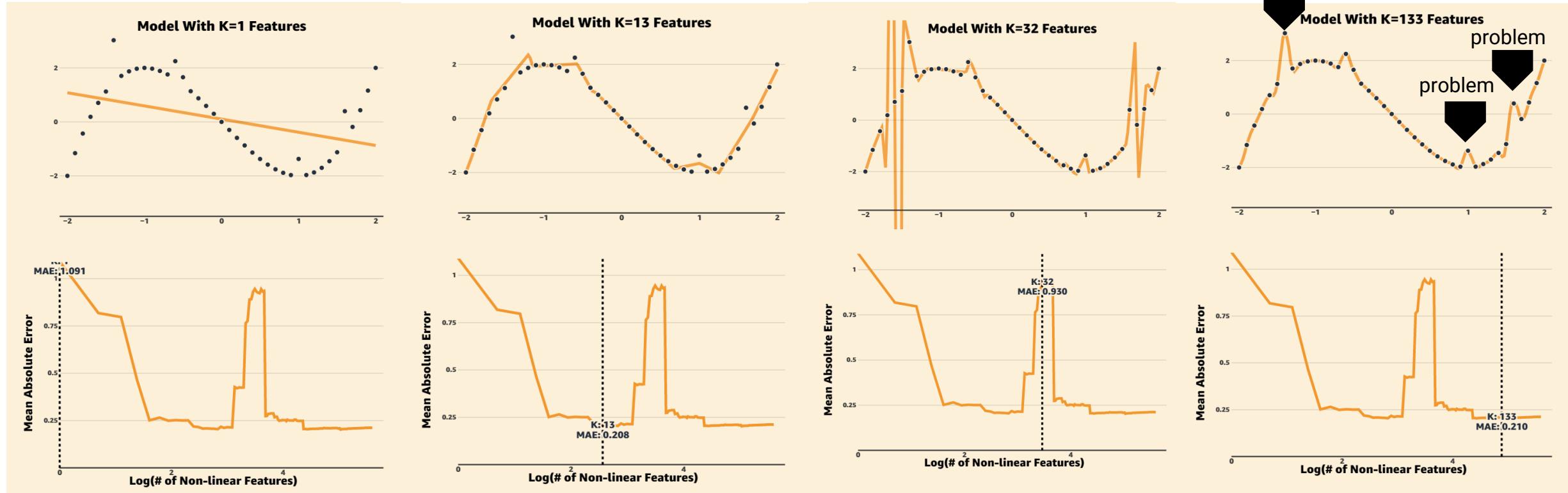
- Fast, specially LightGBM
- Handles known factors (used during the first trees)
- **Overfits**



Sum of subsequent predictions



Overfitting and the double descent



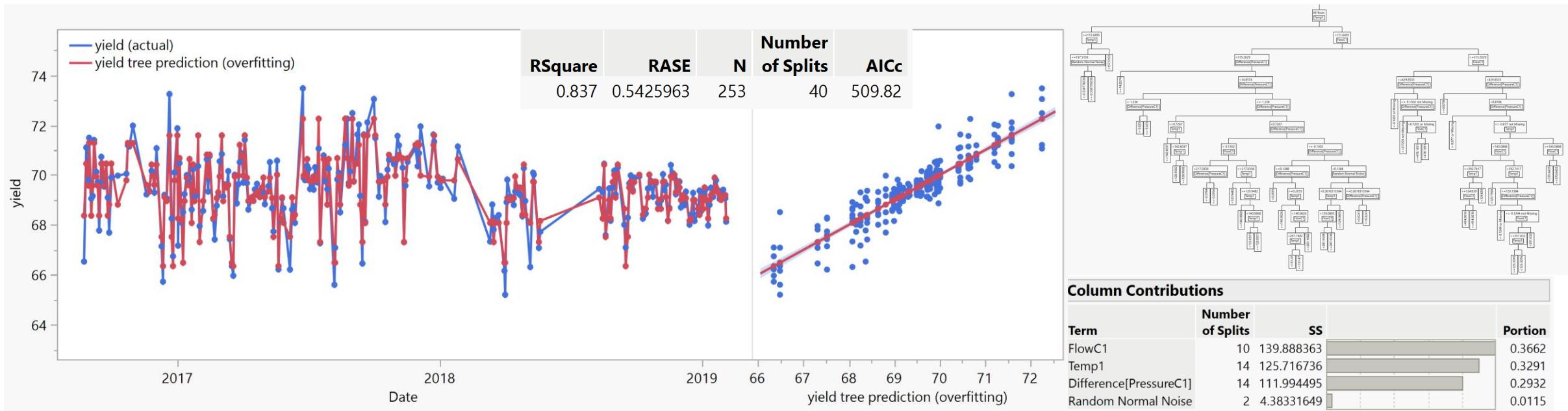
<https://mlu-explain.github.io/double-descent/>



MLU-EXPLAIN

Visual explanations of core machine learning
concepts

Overfitting is a problem if model captures noise

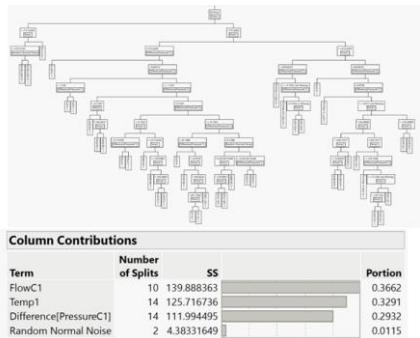


- Data-driven models can adjust perfectly to data given enough parameters, even to noise

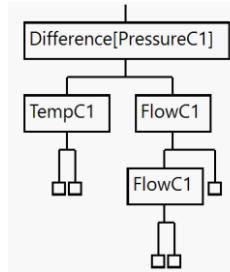
$$data = model + noise$$

Modeling steps in practice

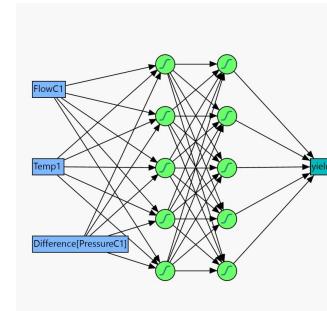
Variable selection



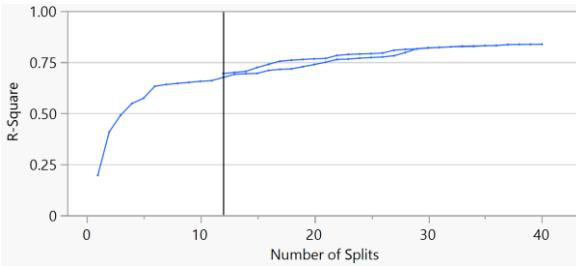
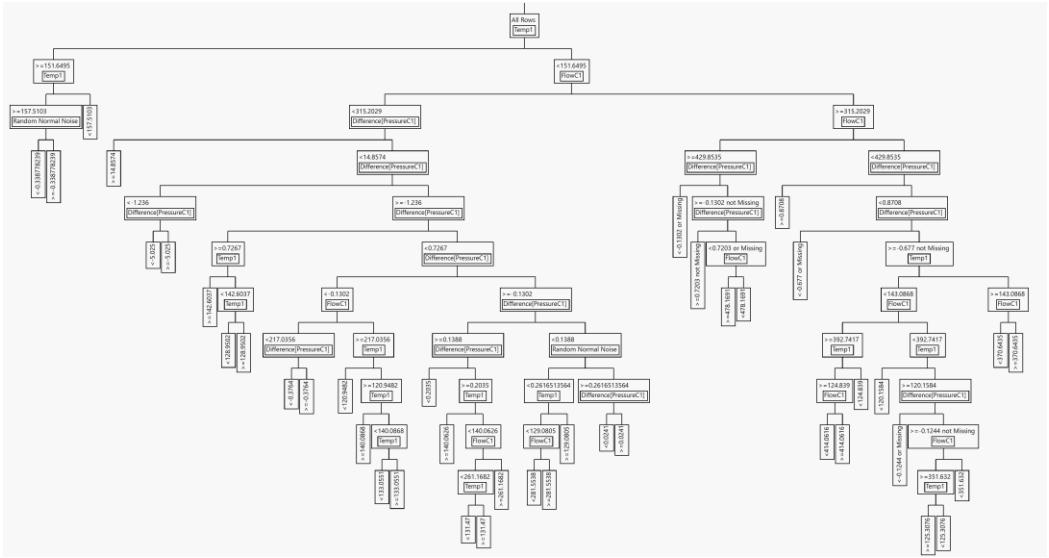
Understanding



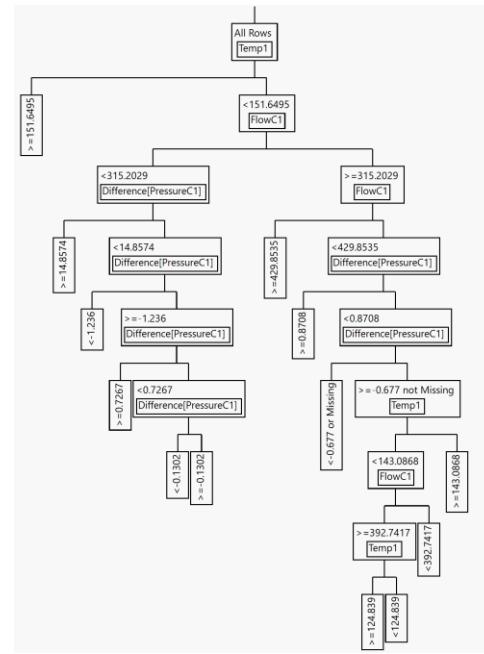
Prediction (if needed)



Tree will be simpler to interpret and implement



Simpler model
(less number of splits)

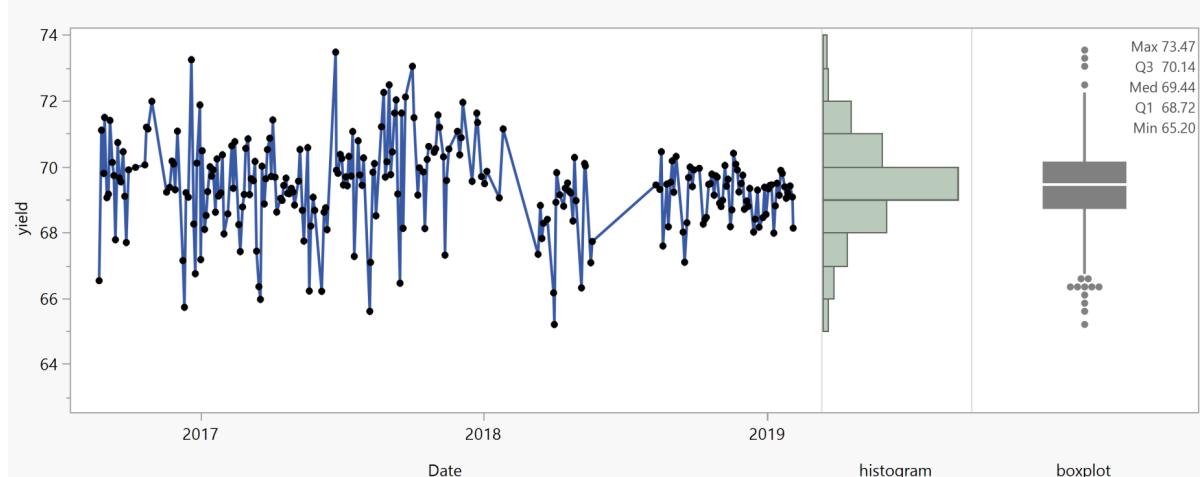


Term	Number of Splits	SS	Portion
FlowC1	3	126.590026	0.3997
Temp1	3	101.429523	0.3202
Difference[PressureC1]	6	88.7241233	0.2801
Random Normal Noise	0	0	0.0000

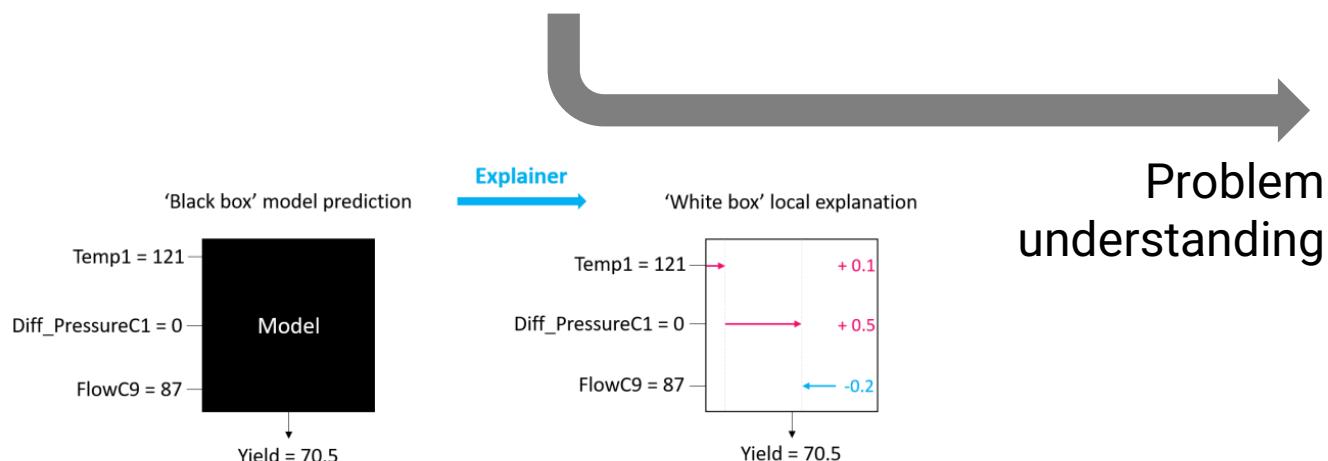
- Adding synthetic noise, a model can be used to describe the main effect of variables aside from common variation.

$$\text{data} + \text{known noise} = \text{model} + \text{noise}$$

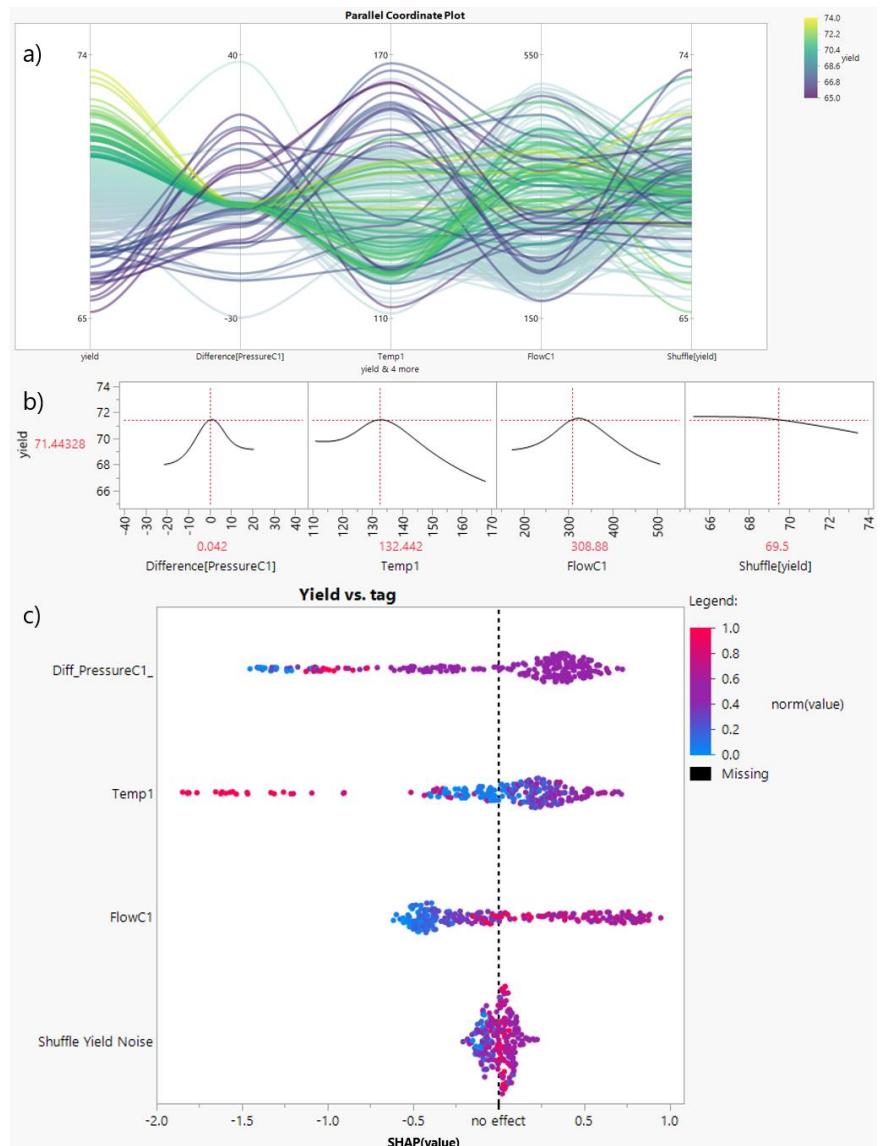
Understanding black boxes with SHAP (ExplainableAI)



Problem definition

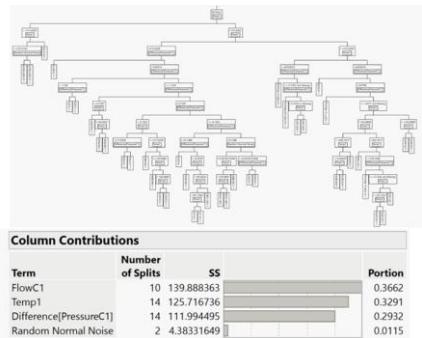


Problem understanding

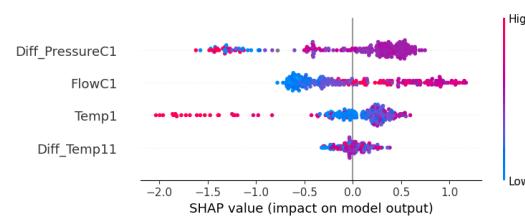
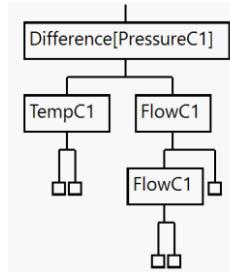


Modeling steps in practice

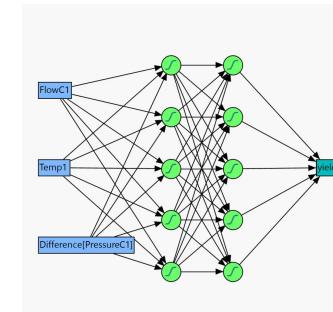
Variable selection



Understanding



Prediction (if needed)



Prediction

*It is difficult to make predictions,
especially about the future*

Anonymous*

Better modeling by subsampling the data

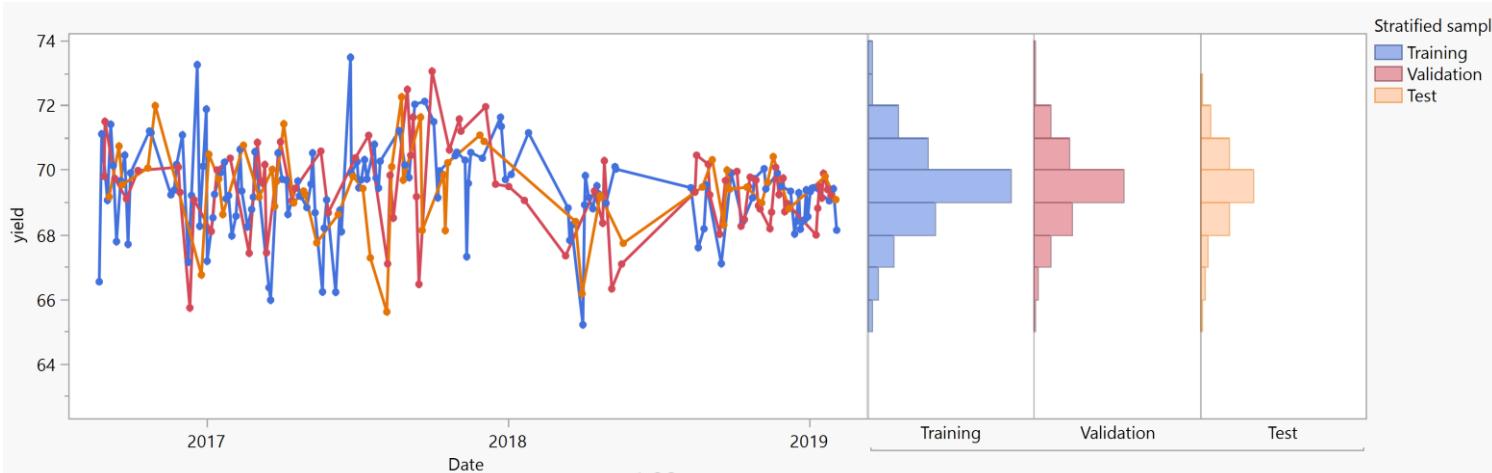


Training and validation splits can also be done multiple times (cross validation)

A rookie mistake? Train/validation/test

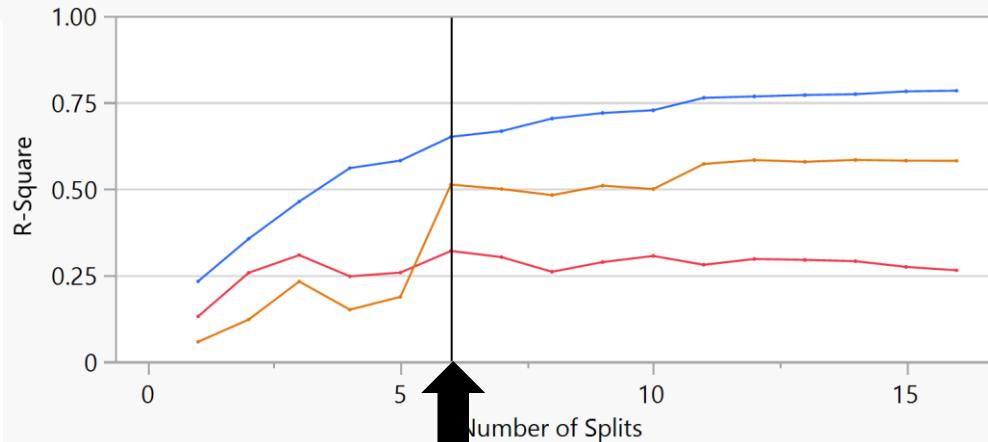
Training = fitting
Validation = overfitting
Test = unseen data

Stratified split
(random split)

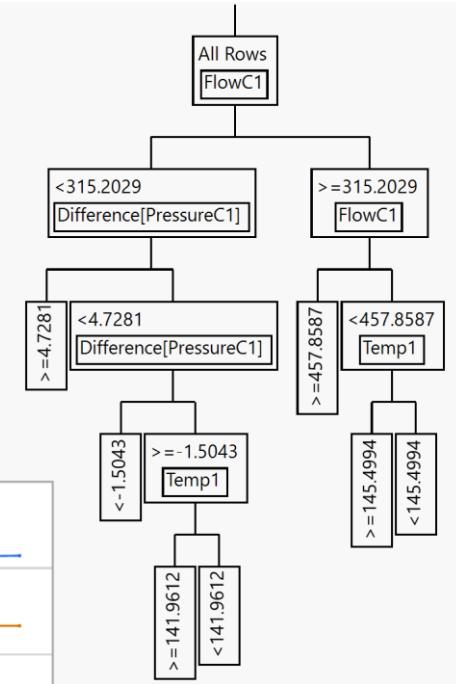


If sampling higher than process dynamics, data points will be very similar (autocorrelated).

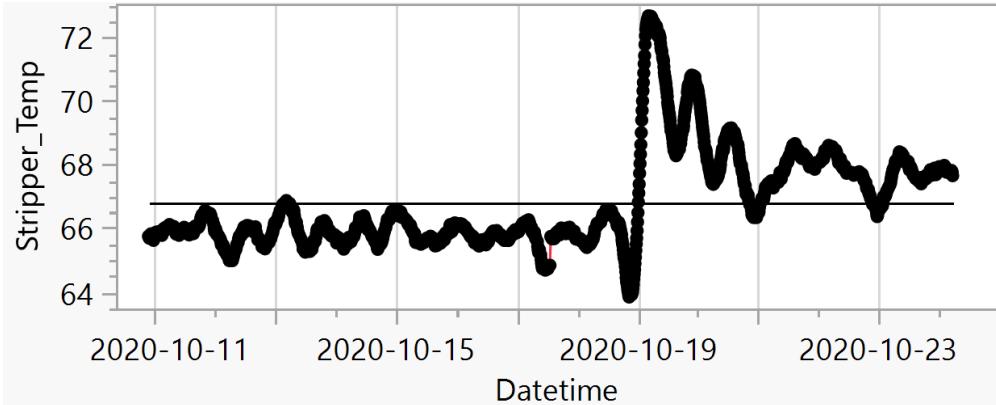
Model will be overfitting without you knowing



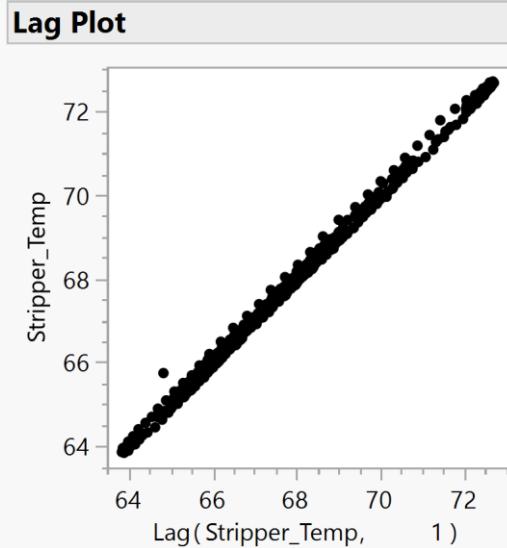
Training stops when validation performance decays



High autocorrelation



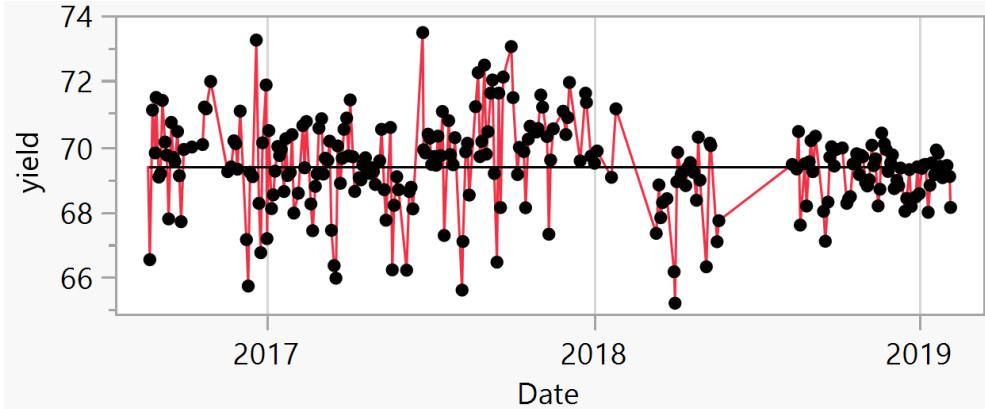
As time series, sensor data often autocorrelation (higher sampling rates). Historians can also interpolate data if high compression was configured (data quality)



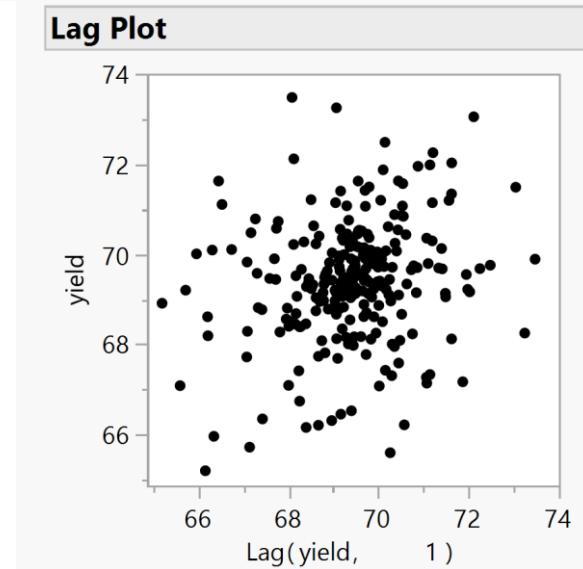
Time Series Basic Diagnostics

Lag	AutoCorr	- .8 .-6 .-4 .-2 .0 .2 .4 .6 .8	Ljung-Box Q	p-Value	Lag	Partial	- .8 .-6 .-4 .-2 .0 .2 .4 .6 .8
0	1.0000	.	.	.	0	1.0000	.
1	0.9988	.	.	.	1	0.9988	.
2	0.9962	.	.	.	2	-0.5363	.
3	0.9929	.	.	.	3	-0.0549	.
4	0.9890	.	.	.	4	-0.0597	.
5	0.9845	.	.	.	5	-0.1442	.
6	0.9791	.	.	.	6	-0.1914	.
7	0.9728	.	.	.	7	-0.1709	.
8	0.9656	.	.	.	8	-0.1307	.
9	0.9575	.	.	.	9	-0.0913	.
10	0.9486	.	.	.	10	-0.0828	.
11	0.9390	.	.	.	11	-0.0731	.
12	0.9287	.	.	.	12	-0.0416	.
13	0.9178	.	.	.	13	-0.0284	.
14	0.9062	.	.	.	14	0.0220	.
15	0.8942	.	.	.	15	0.0309	.
16	0.8816	.	.	.	16	0.0264	.
17	0.8687	.	.	.	17	0.0460	.
18	0.8553	.	.	.	18	0.0473	.
19	0.8417	.	.	.	19	0.0459	.
20	0.8277	.	.	.	20	0.0509	.
21	0.8136	.	.	.	21	0.0482	.
22	0.7993	.	.	.	22	0.0541	.
23	0.7849	.	.	.	23	0.0418	.
24	0.7704	.	.	.	24	0.0503	.
25	0.7560	.	.	.	25	0.0503	.

No autocorrelation



Distillation column example sampled only when process reached steady state, hence yield rows are not correlated to each other

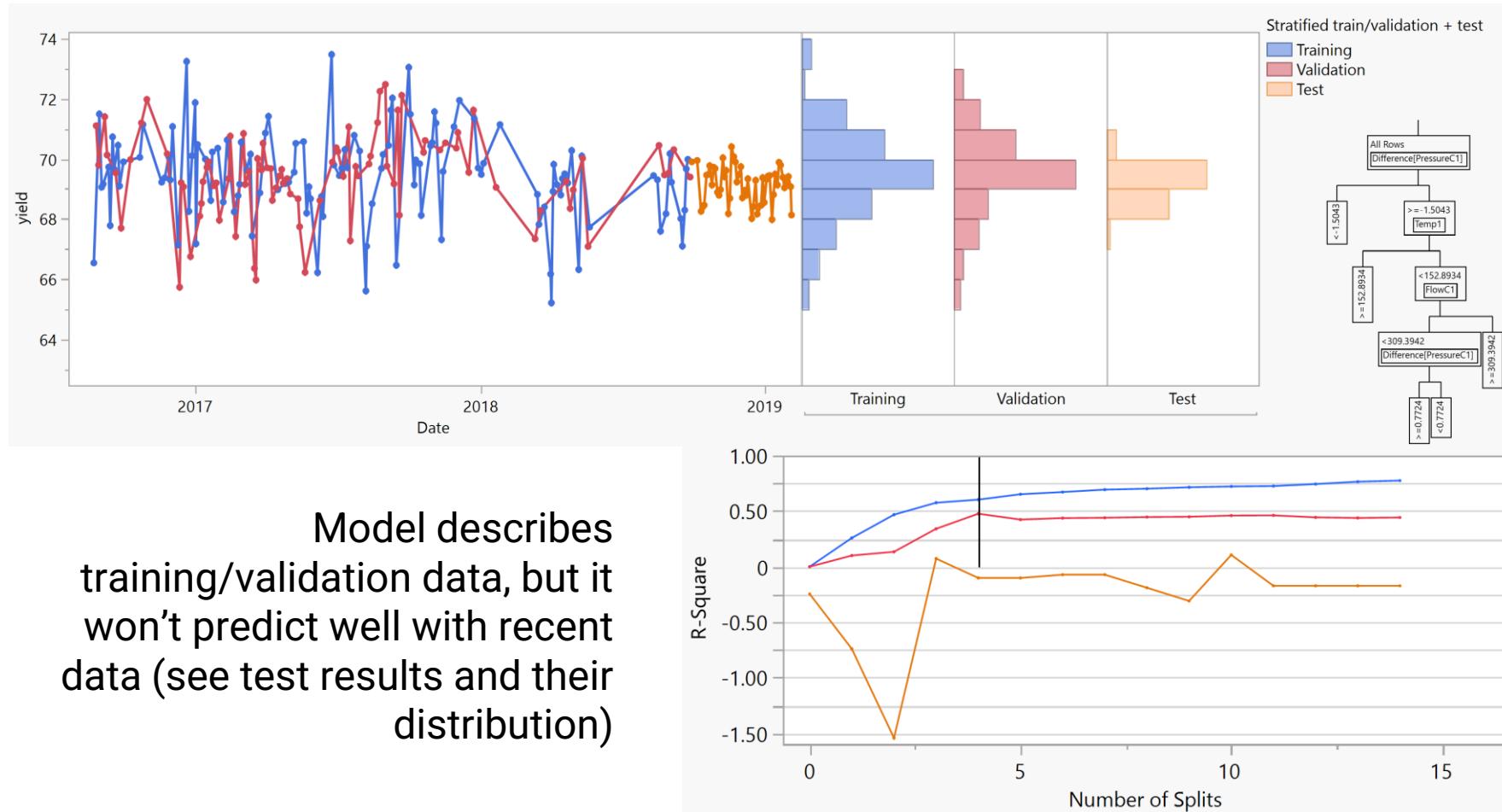


Time Series Basic Diagnostics

Lag	AutoCorr	-0.8	-0.6	-0.4	-0.2	0	0.2	0.4	0.6	0.8	Ljung-Box Q	p-Value	Lag	Partial	-0.8	-0.6	-0.4	-0.2	0	0.2	0.4	0.6	0.8
0	1.0000												0	1.0000									
1	0.1971										9.9490	0.0016*	1	0.1971									
2	0.0593										10.8526	0.0044*	2	0.0213									
3	0.1152										14.2758	0.0026*	3	0.1036									
4	0.0558										15.0835	0.0045*	4	0.0141									
5	0.0305										15.3251	0.0091*	5	0.0116									
6	0.0732										16.7242	0.0104*	6	0.0556									
7	0.0302										16.9629	0.0176*	7	-0.0014									
8	-0.0093										16.9857	0.0303*	8	-0.0229									
9	0.0355										17.3189	0.0439*	9	0.0292									
10	0.0832										19.1558	0.0383*	10	0.0690									
11	0.0746										20.6398	0.0373*	11	0.0496									
12	0.0061										20.6498	0.0558	12	-0.0299									
13	0.0183										20.7399	0.0782	13	0.0033									
14	-0.0003										20.7400	0.1085	14	-0.0186									
15	0.0934										23.1041	0.0820	15	0.0976									
16	0.0319										23.3804	0.1039	16	-0.0147									
17	0.0823										25.2318	0.0896	17	0.0750									
18	-0.0233										25.3802	0.1148	18	-0.0719									
19	0.0157										25.4485	0.1463	19	0.0262									
20	0.0506										26.1586	0.1606	20	0.0230									
21	0.0607										27.1841	0.1648	21	0.0386									
22	0.1495										33.4282	0.0561	22	0.1341									
23	-0.0890										35.6479	0.0448*	23	-0.1691									
24	0.0410										36.1214	0.0534	24	0.0894									
25	0.0436										36.6594	0.0622	25	-0.0205									

Keeping most recent data as test

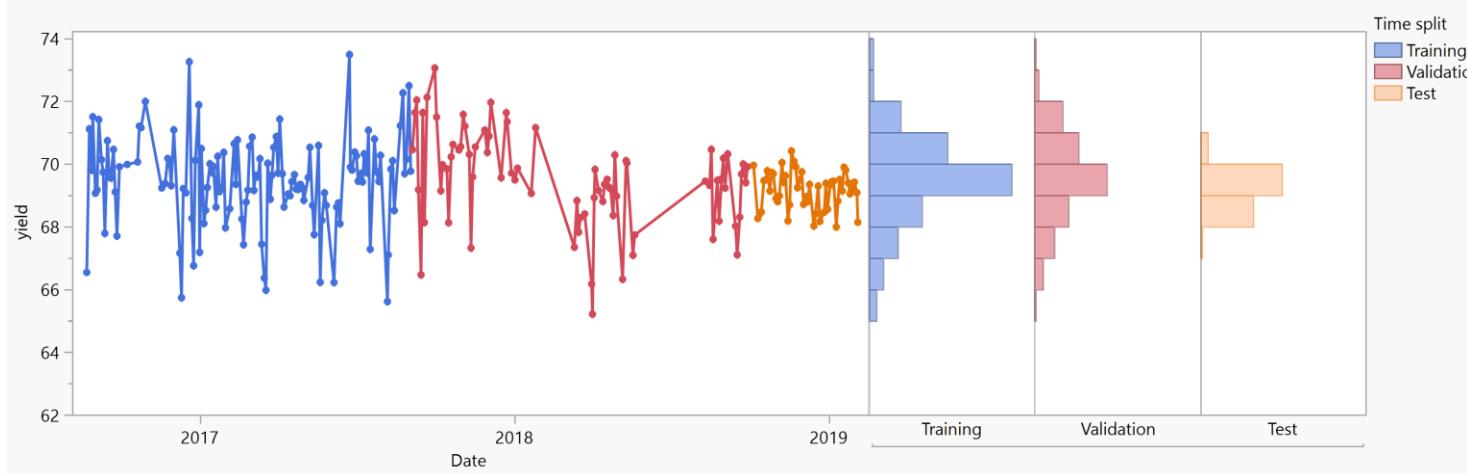
Stratified split
+ test



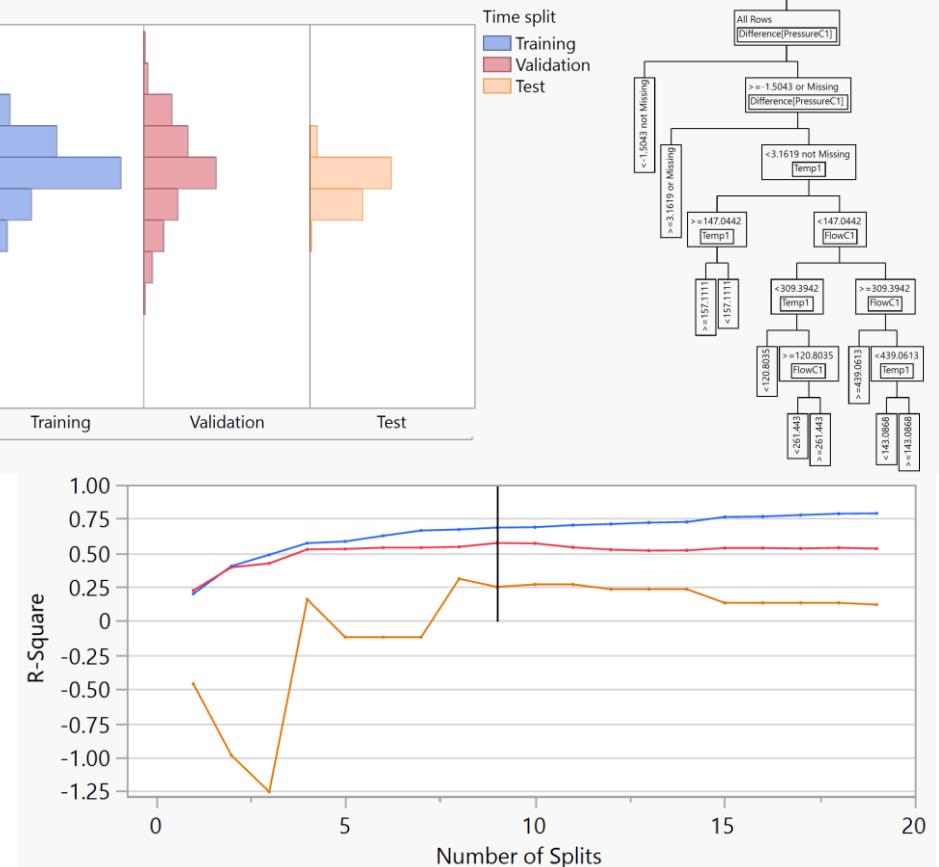
Time split, usually a better approach

Training = fitting
Validation = overfitting
Test = unseen data

Time-split
(or cut point)

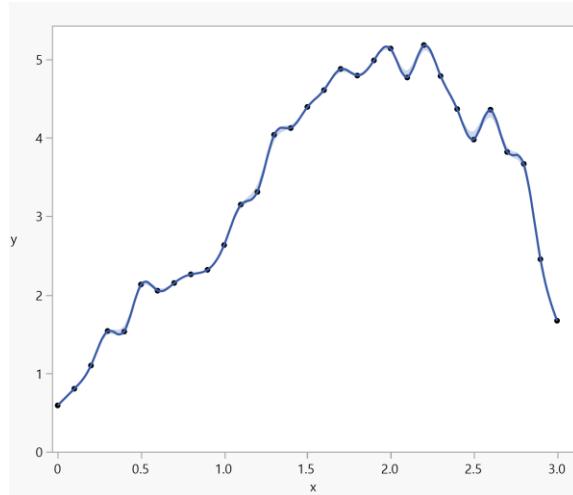


Model describes training/validation data, but it won't predict well with recent data if process has changed over time. **Use your test data wisely and try to obtain simple models.**

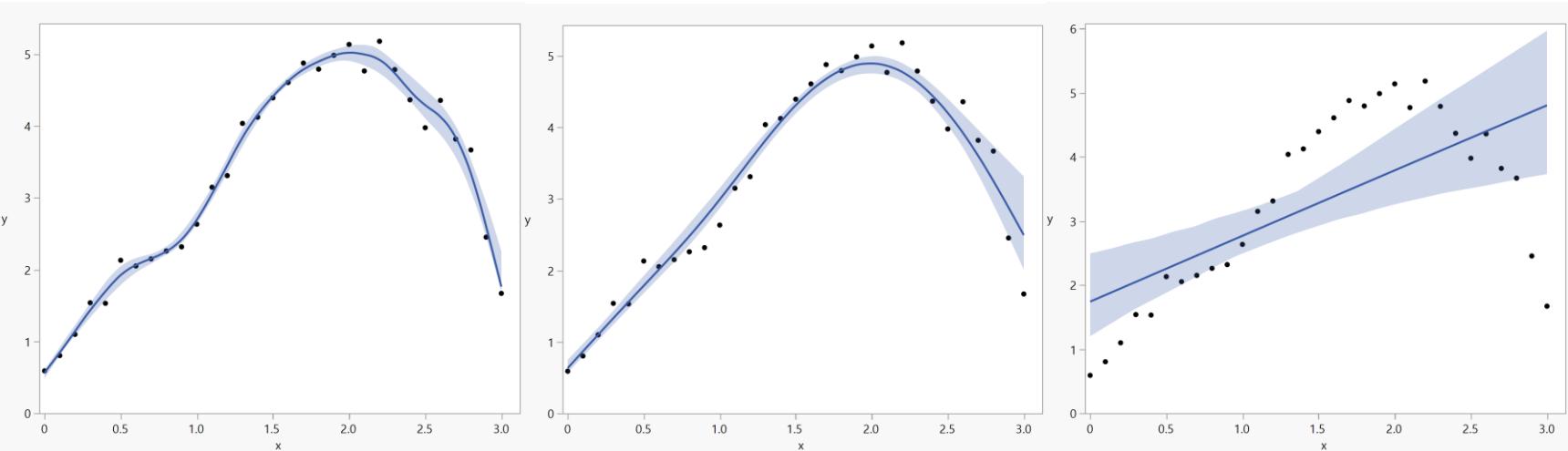


An accuracy-simplicity tradeoff (penalized regression)

overfitting



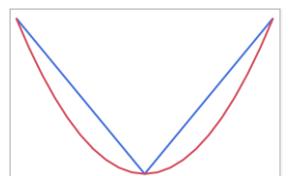
underfitting



$$Model: \hat{y} = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + \beta_4 x^4 + \beta_5 x^5 + \dots$$

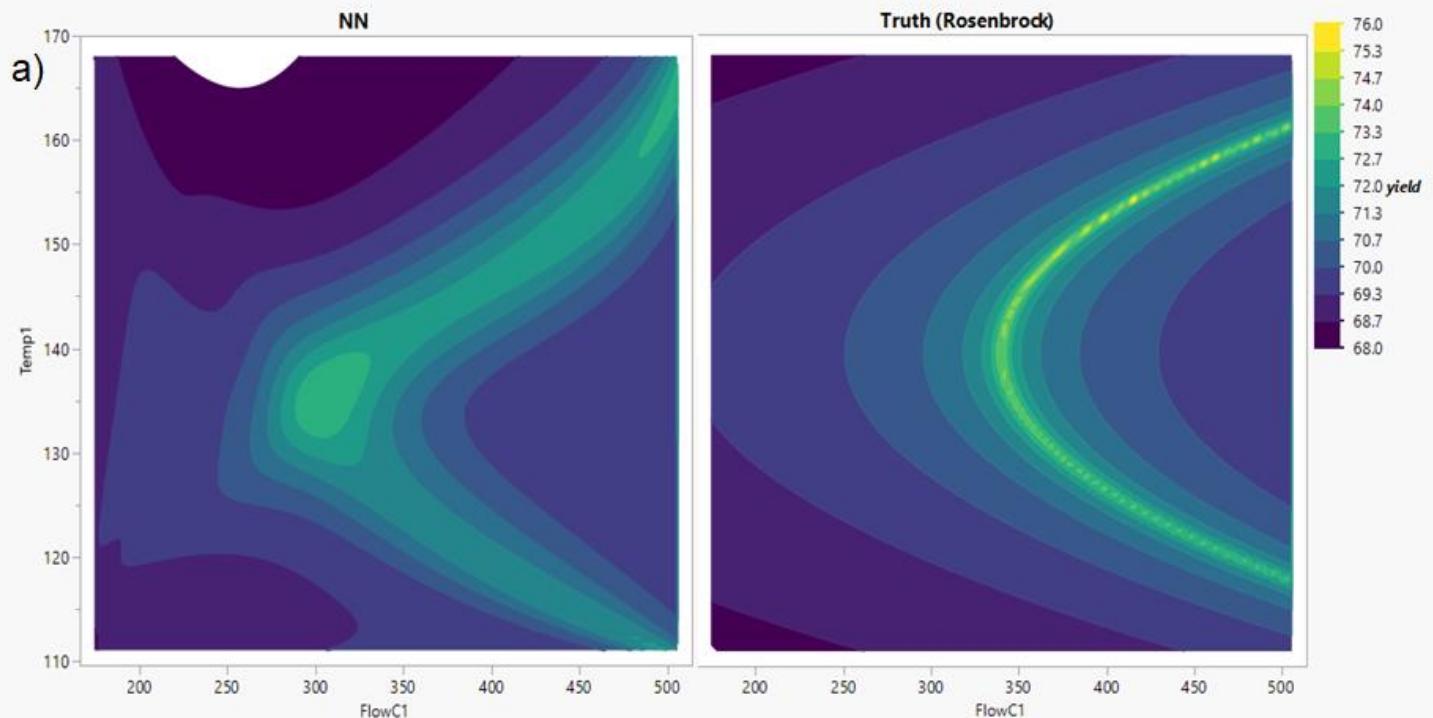
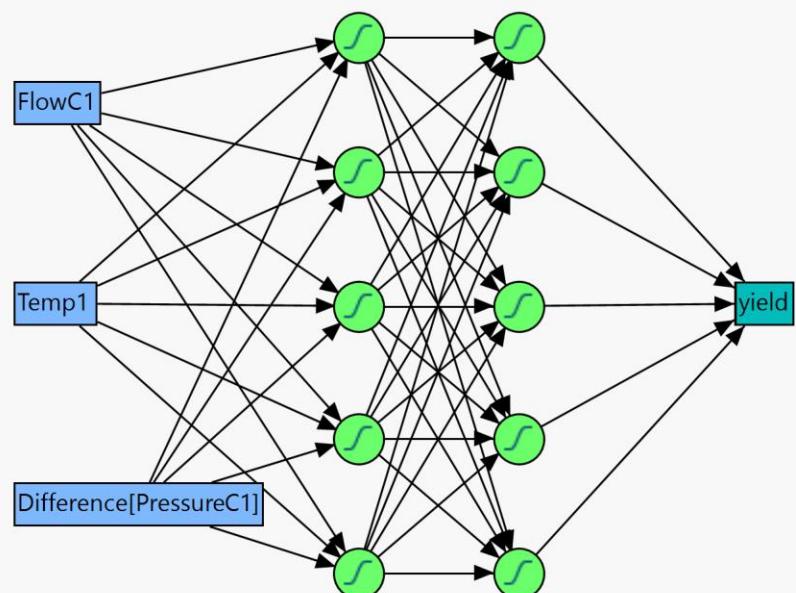
$$\text{Minimize: } \sum(y - \hat{y})^2 + \lambda \sum |\beta_j|$$

Lasso regression minimizes **prediction error** as well as the **model coefficients**, being able to eliminate them entirely. Other penalization techniques exist (e.g., ridge regression, BIC, AIC...) How to select the best one? Don't use all your data ☺



Complex models can approximate anything

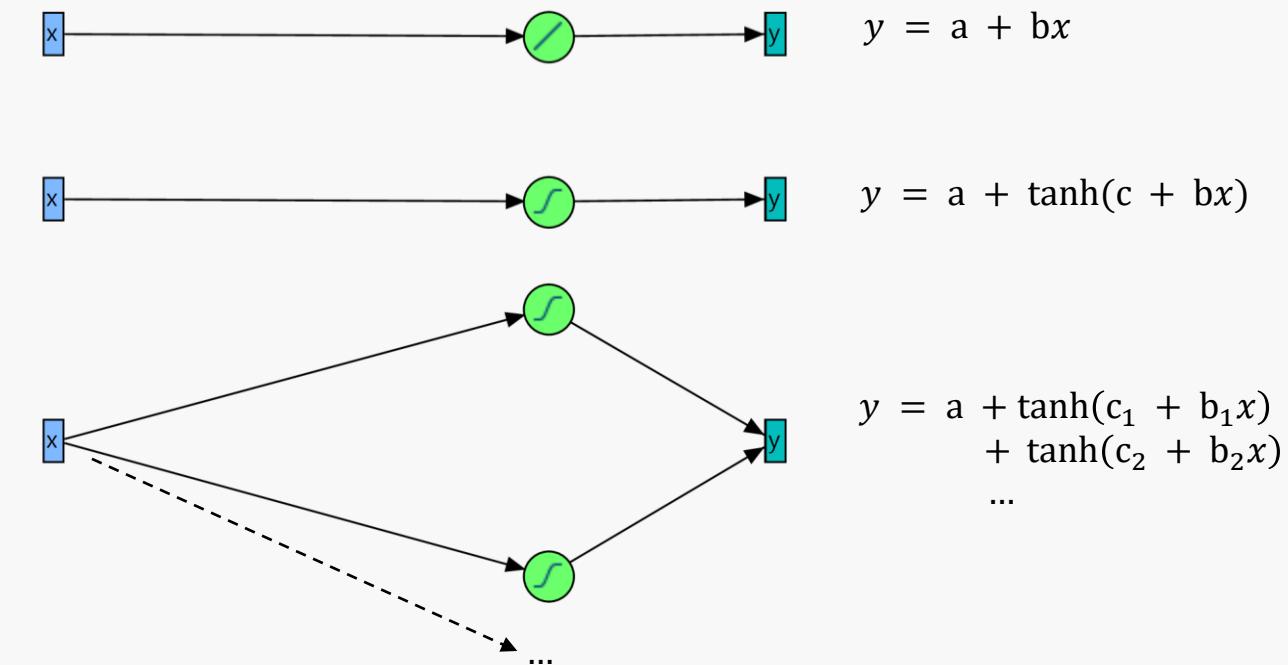
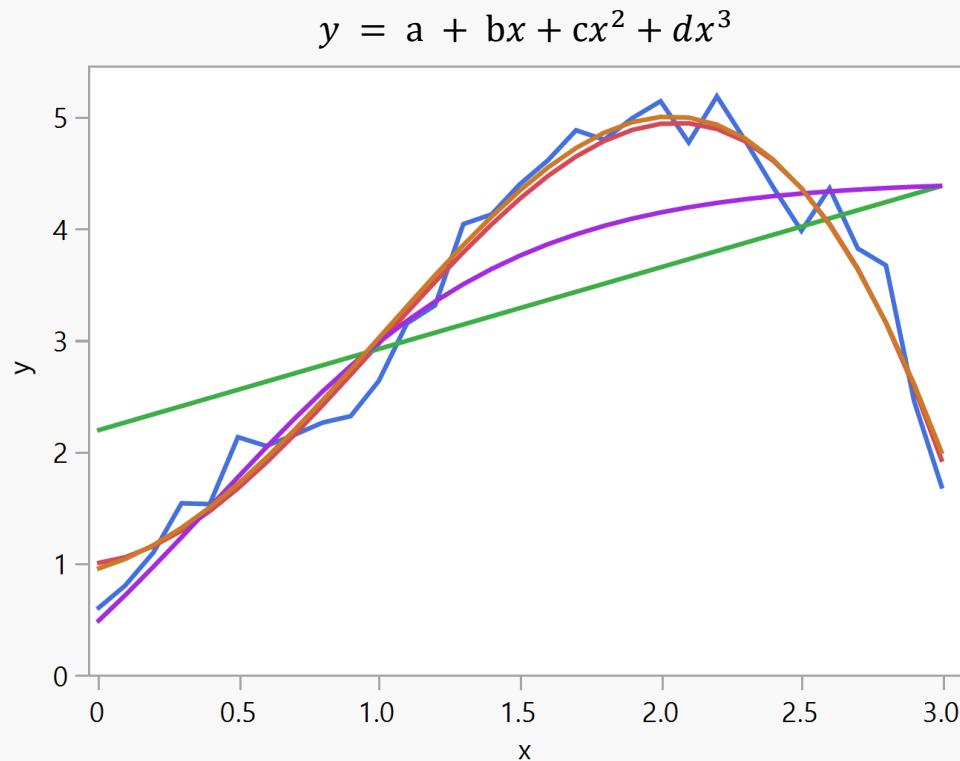
- By combining multiple nodes and layers, Neural Networks can be used to capture intricate non-linear correlations within datasets, so what are they?



Source:
[Industrial data science – a review of machine learning applications for chemical and process industries](#)
React. Chem. Eng., 2022, 7, 1471-1509

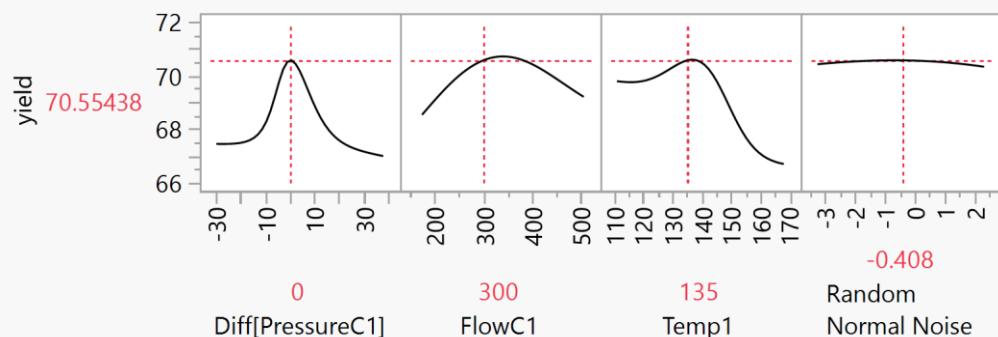
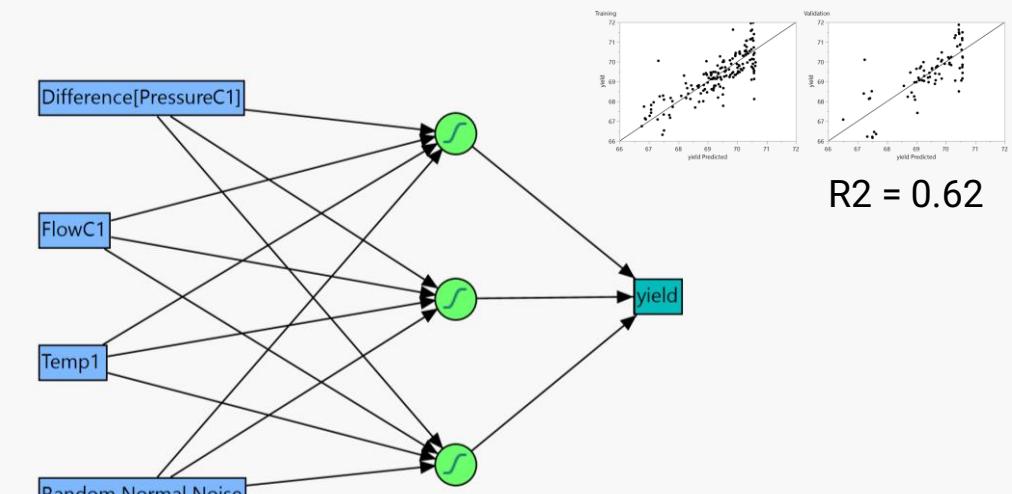
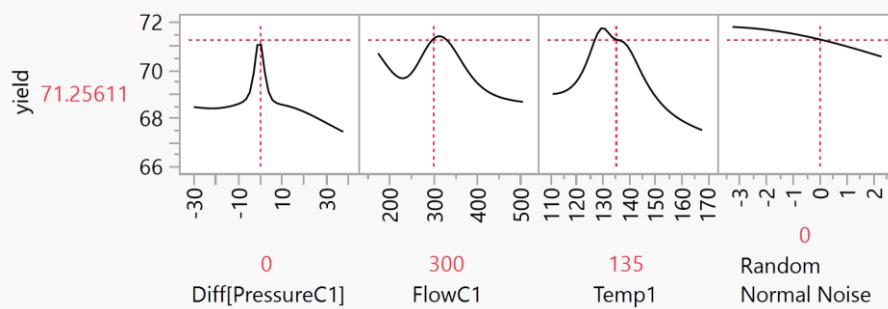
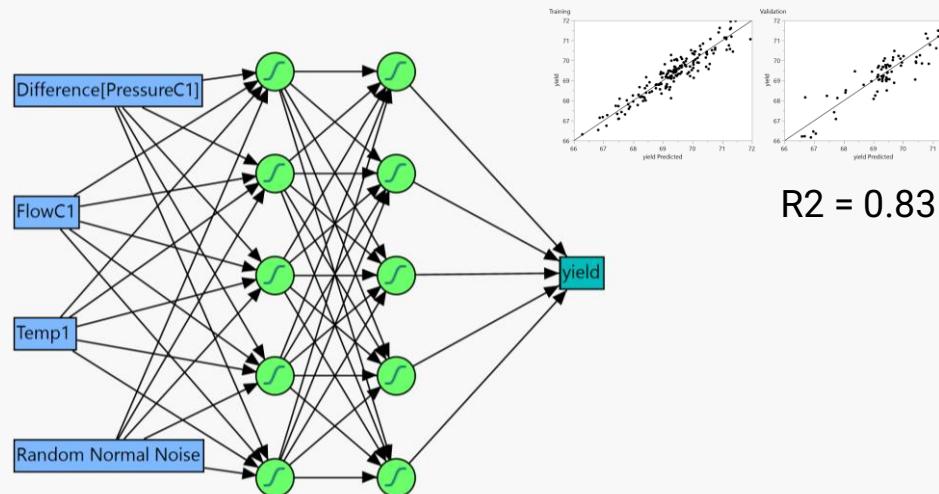
Neural Networks are function approximators

- Like polynomials, Neural Networks approximate functions by increasing the number of non-linear elements (sigmoid-like functions call nodes or activation functions)



You probably don't need deep learning

- Their can also fit the noise in the data, adding noise to see variable importance is recommend to detect overfitting and discard less impactful parameters.

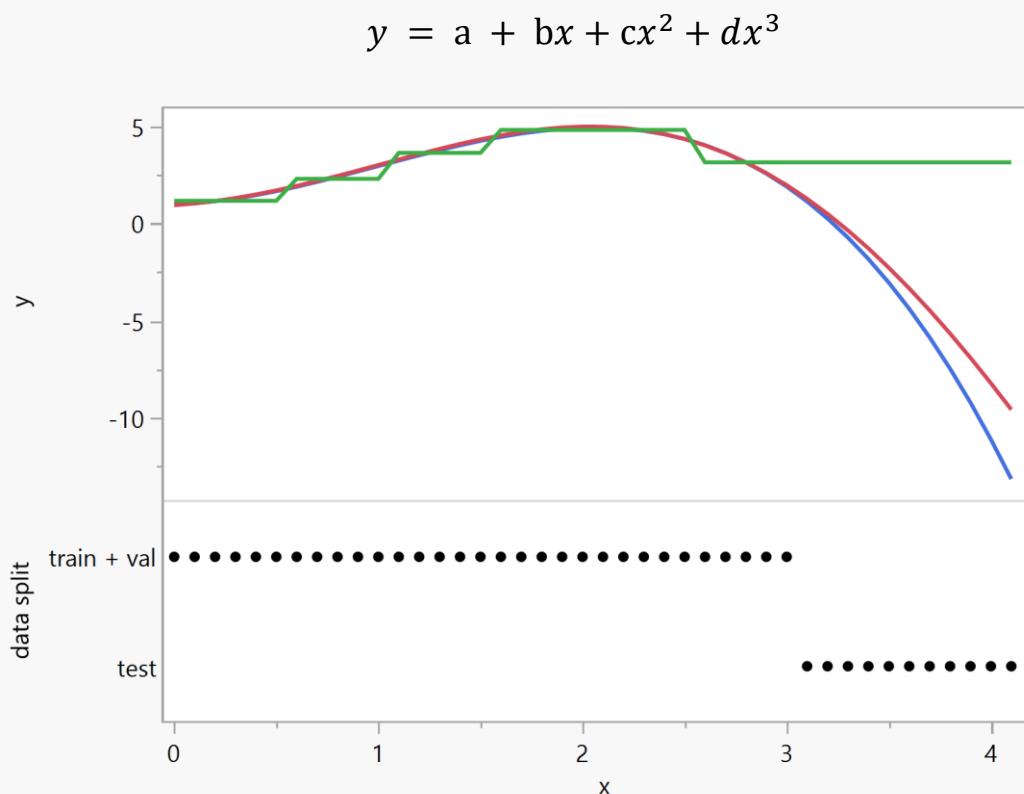


Source:

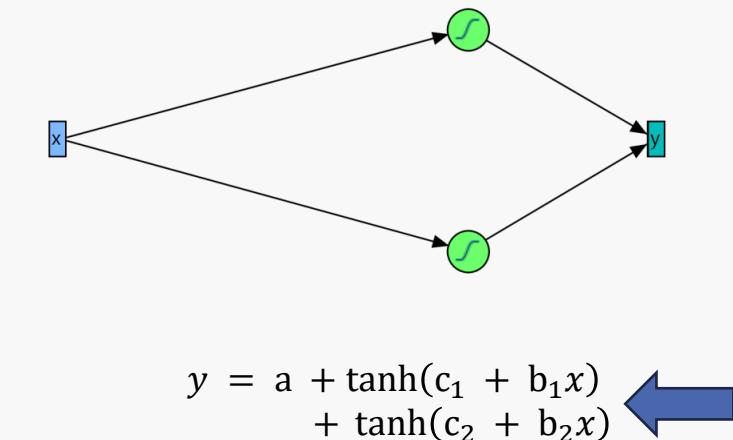
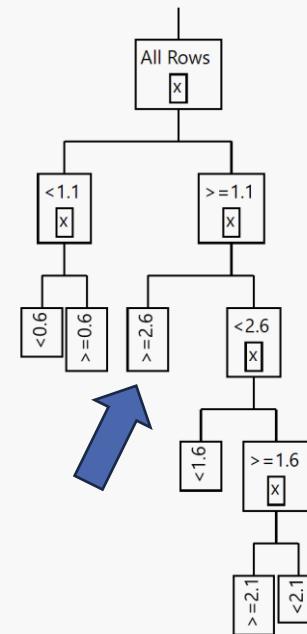
Industrial data science – a review of machine learning applications for chemical and process industries
React. Chem. Eng., 2022, 7, 1471-1509

Interpolation vs extrapolation

- Like polynomials, Neural Networks extrapolate outside training data
- In practice, Neural Networks require anomaly detection (unsupervised learning) to know how far from known data the model is
- Tree-based models flatten outside training data



— y (truth)
— \hat{y} (NN 2 nodes)
— \hat{y} (tree)
● x

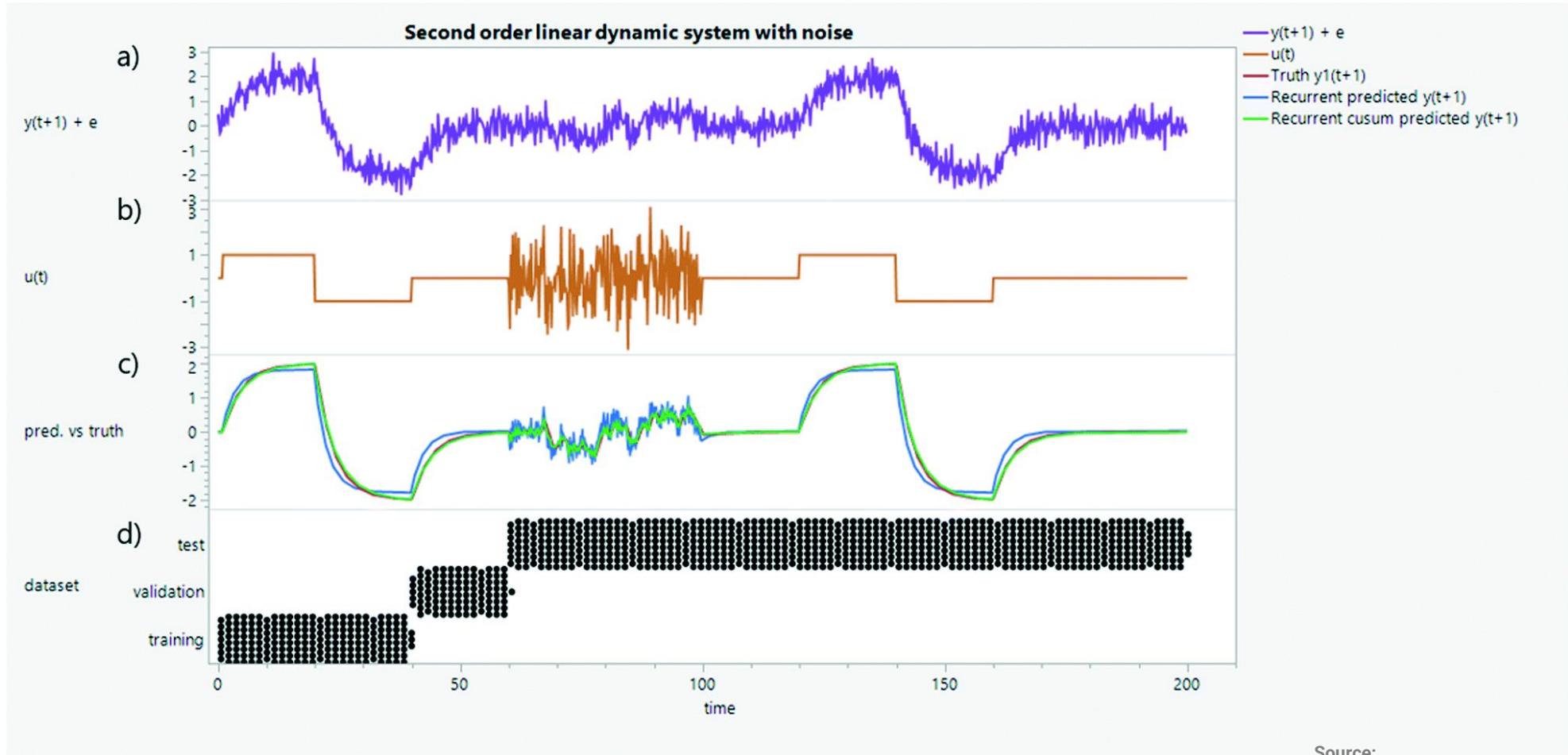


Source:

Industrial data science – a review of machine learning applications for chemical and process industries
React. Chem. Eng., 2022, 7, 1471-1509

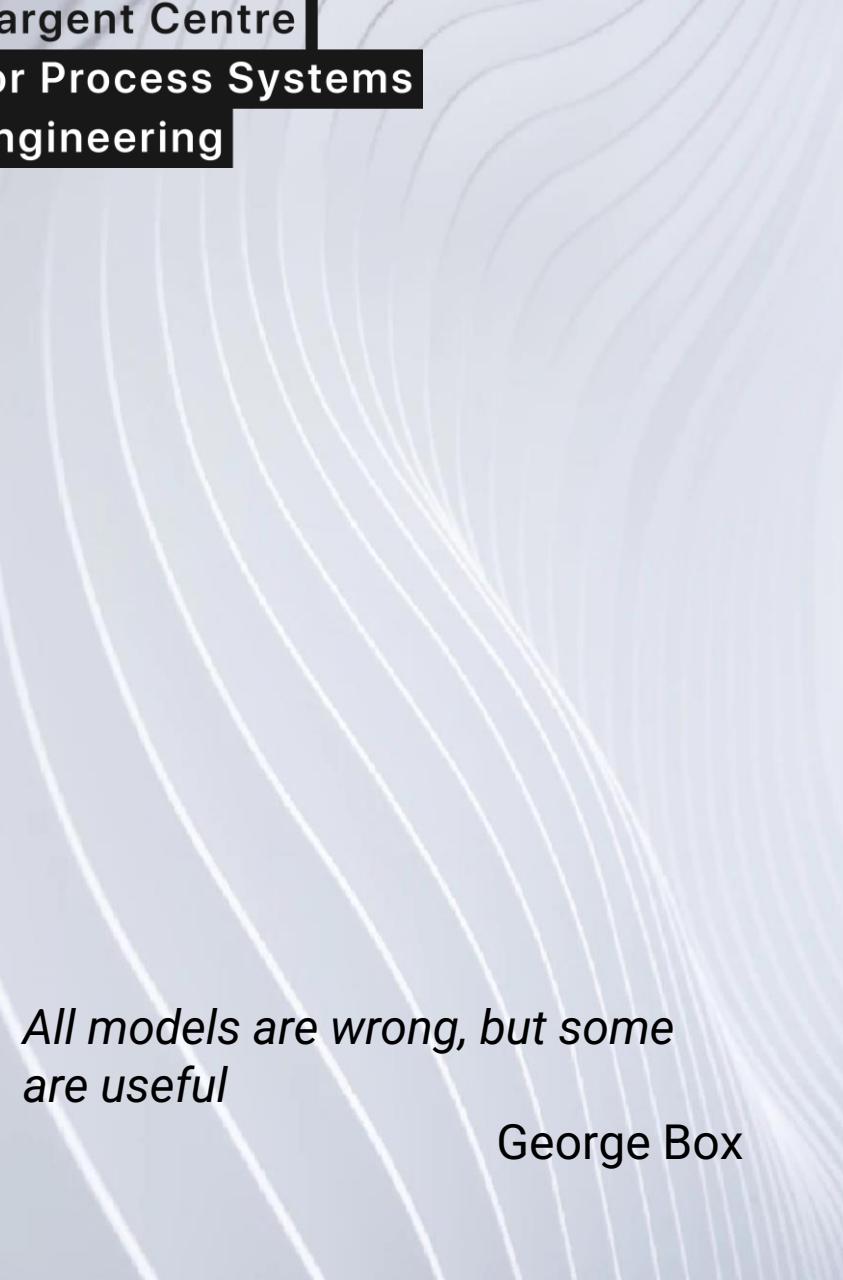
Application: System Identification

- Predicting one step ahead



Summary

- AutoML compares multiple models and optimizes their tuning parameters, so you don't have to
- Prediction is only important after problem is understood, focus on descriptive models first
- Tree-based models are best for screening and understanding. SHAP is a common approach to explain them
- Neural networks can produce more accurate results at higher risk of overfitting and extrapolation
- Use synthetic noise and data splits for better/simpler models
- We aware of process dynamics (autocorrelation) when splitting data, and the field of system identification for one-step ahead model predictions.

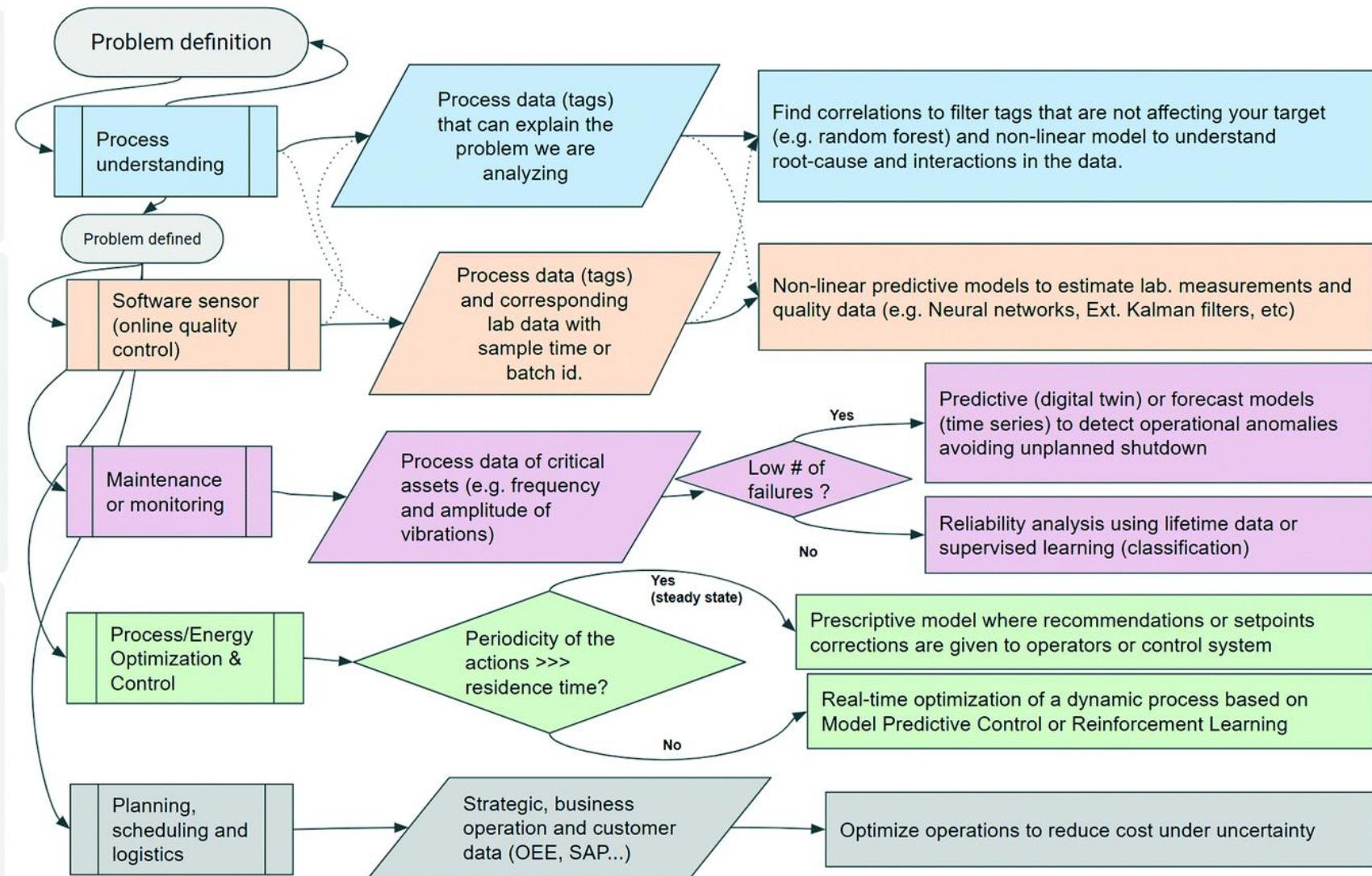


Industrial applications



Industrial data science applications

Descriptive analytics (diagnostics)



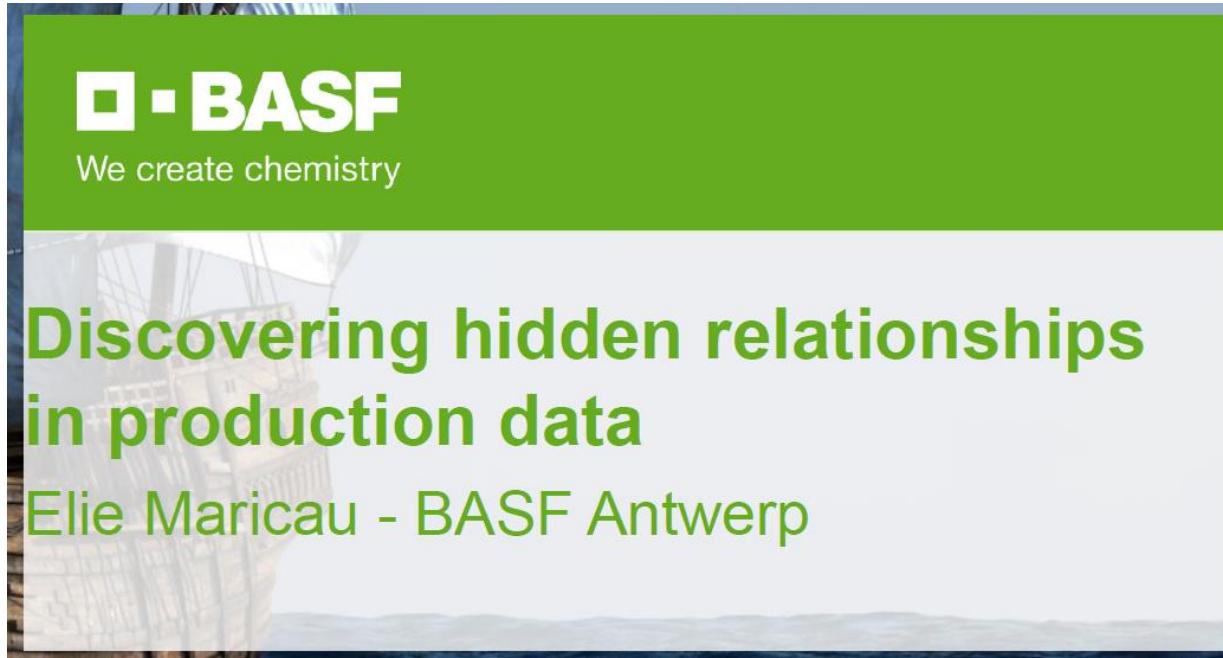
ONLINE
(implementation +
industrialization)

Predictive
analytics

Prescriptive
analytics

Source:

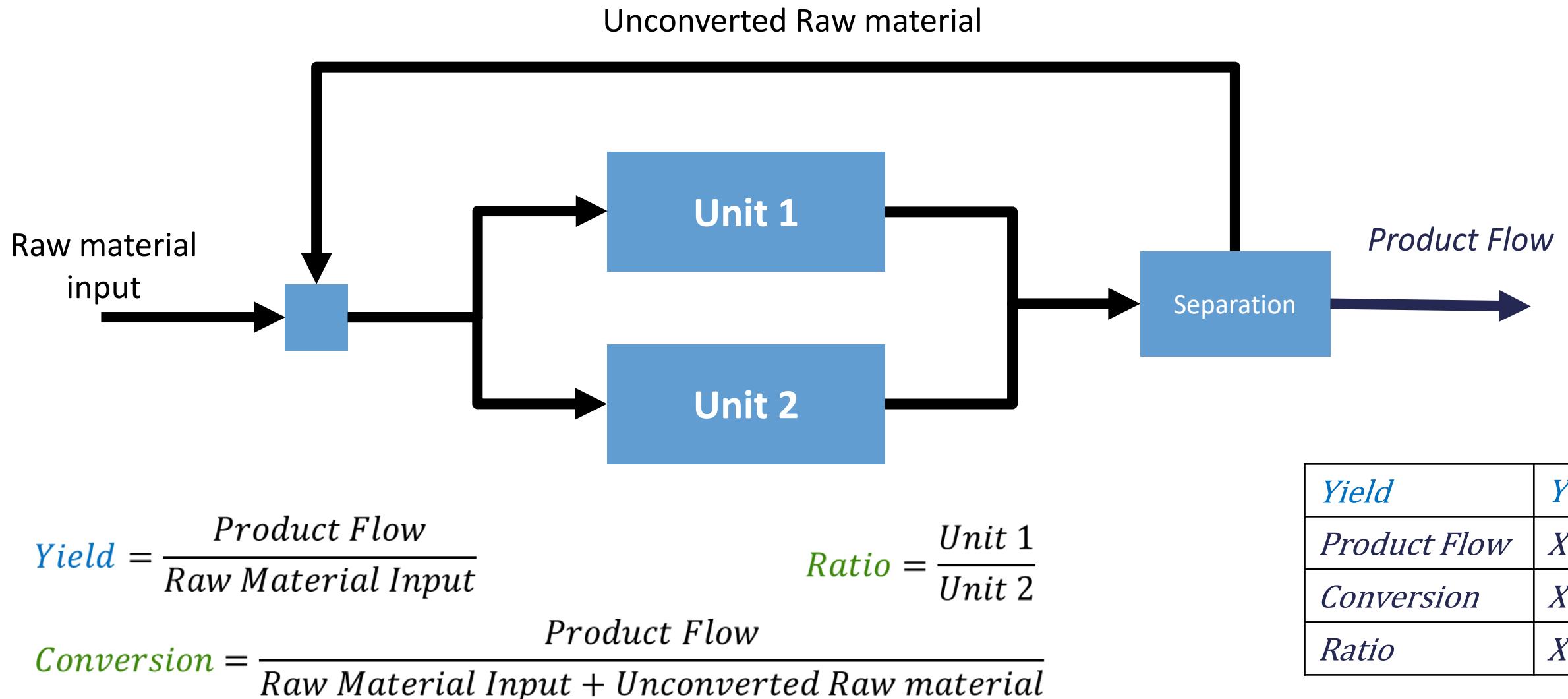
[Industrial data science – a review of machine learning applications for chemical and process industries](#)
React. Chem. Eng., 2022, 1471-1509



<https://community.jmp.com/t5/Discovery-Summit-Europe-2018/Discovering-hidden-relationships-in-production-data-EU2018-113/ta-p/51283>

<https://community.jmp.com/t5/Discovery-Summit-Europe-2019/DOE-for-World-Scale-Manufacturing-Processes-Can-We-Do-Better/ta-p/184245>

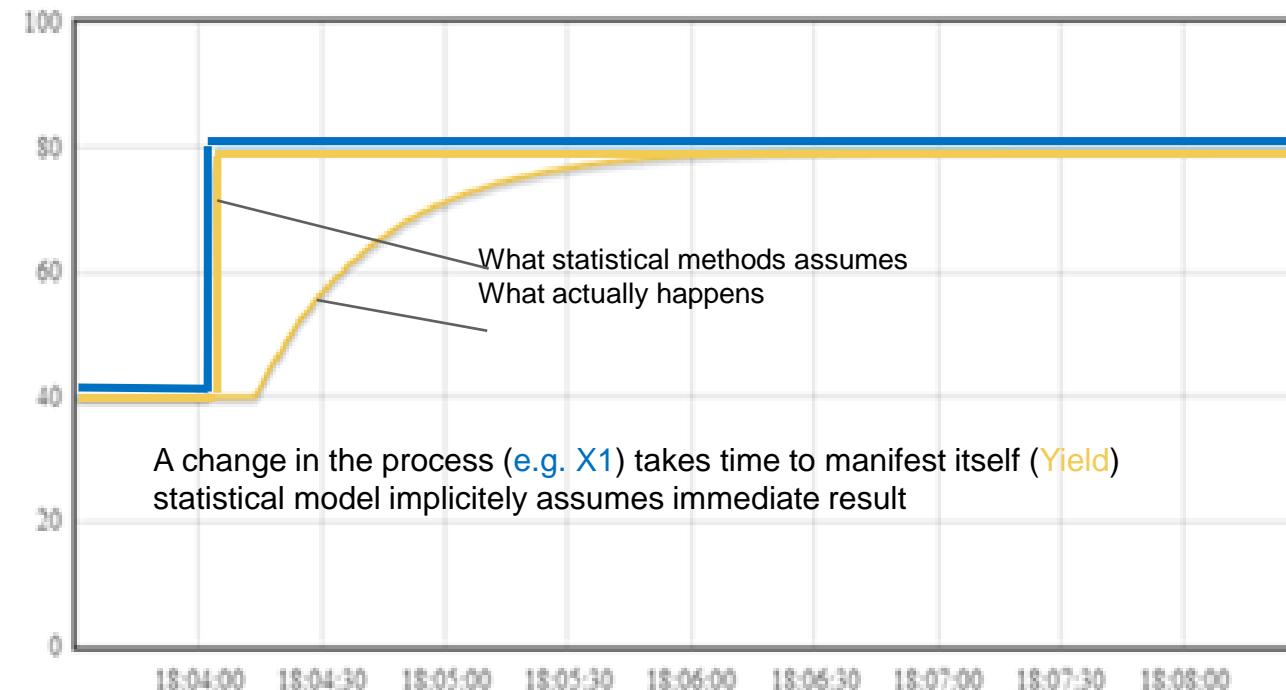
The problem



Data analysis in continuous plant

Do we need non steady-state data?

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

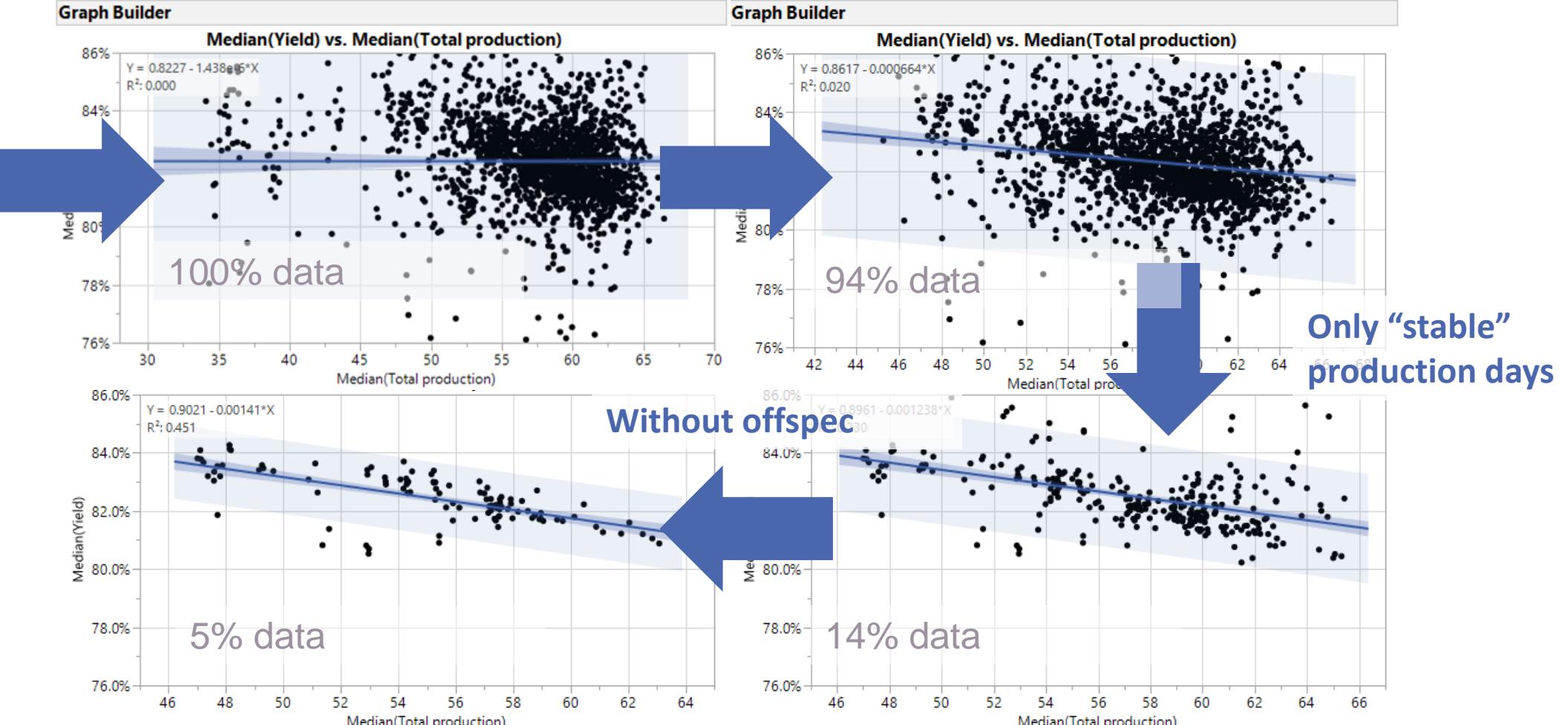


Stability raw material 1 last 48H ²
+ Stability production 1 last 48H ²
+ Stability raw material 2 last 48H ²
+ Stability production 2 last 48H ²
+ Stability offspec last 48H ²
+ Stability
+ Stability
Std Dev
Raw_material_1,
Lag (Raw_material_1, 1),
Lag (Raw_material_1, 2),
Lag (Raw_material_1, 3),
Lag (Raw_material_1, 4),
Lag (Raw_material_1, 5),
Lag (Raw_material_1, 6),
Lag (Raw_material_1, 7),

Data cleaning is key

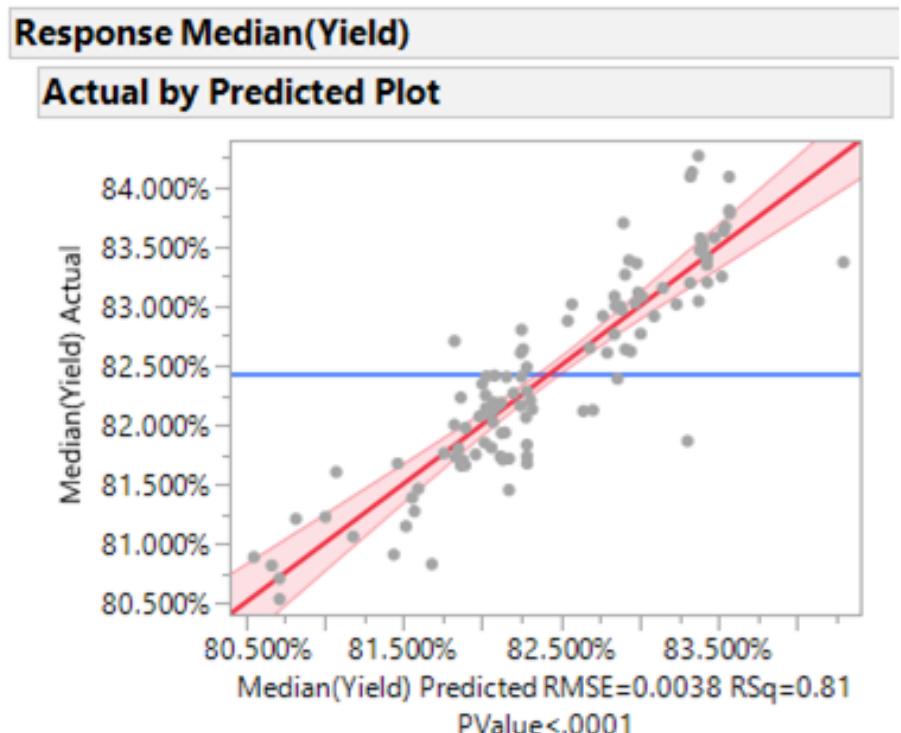
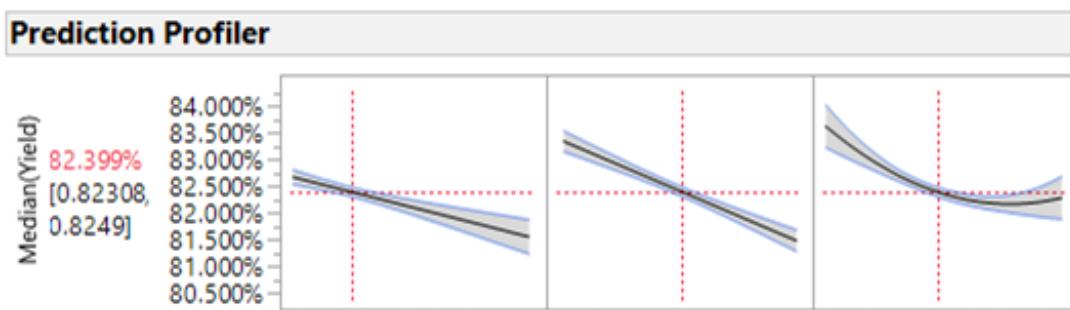
Only medium to high production volumes

Data: 24H average values

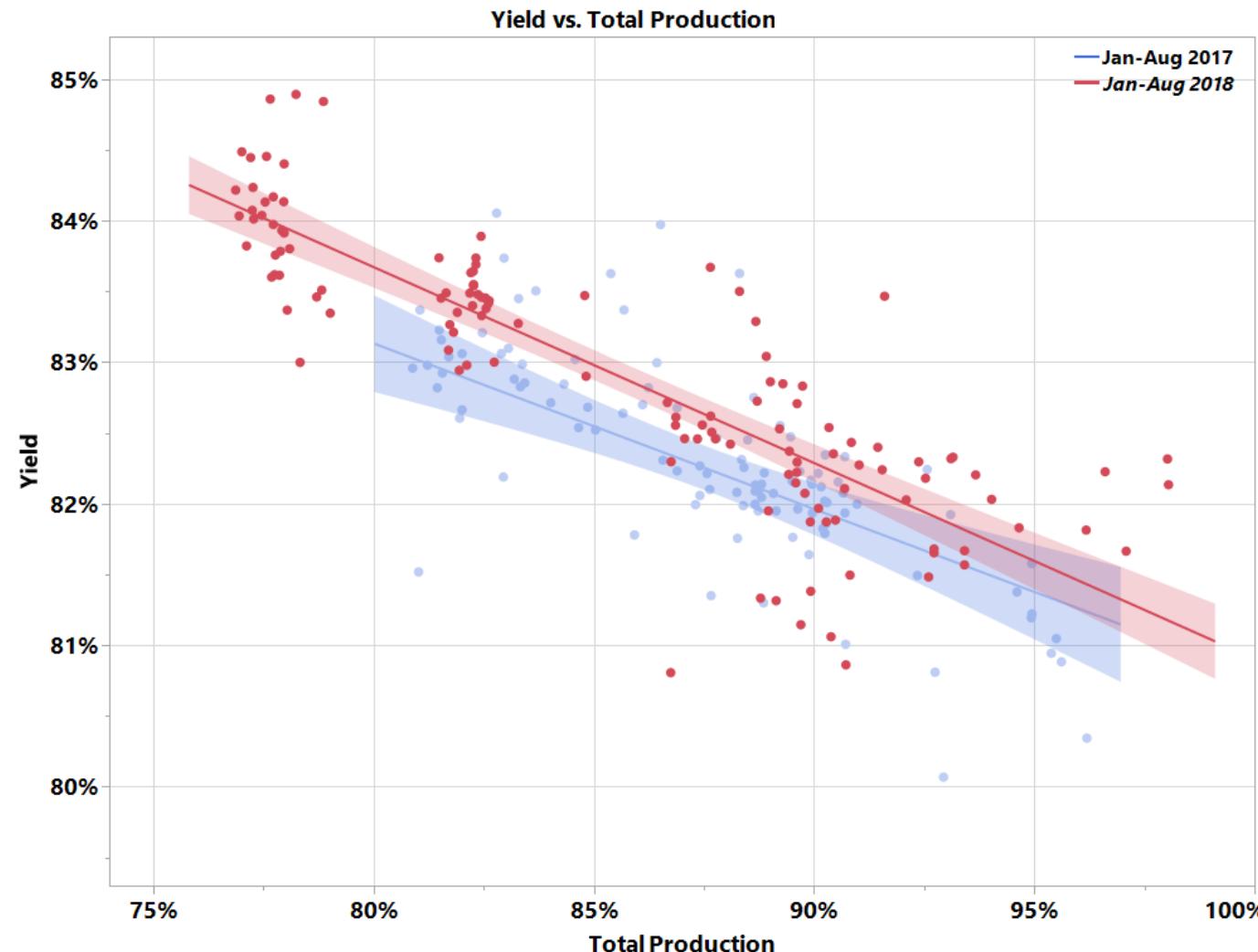


A regression model is enough

- 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 (and detection of deviations)



What benefit did we achieve



Ambition at the start of the project

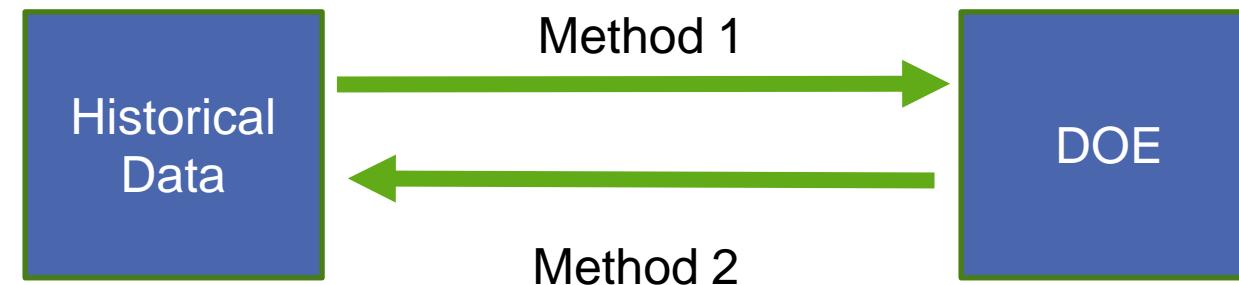
- 0.3% yield improvement (0.4 M€/y)

Result

- Total Yield improvement of 1.1% (1.4 M€/y)
(Jan-Aug 2018 vs Jan-Aug 2017)
- Similar benefit for 2019 wrt 2017

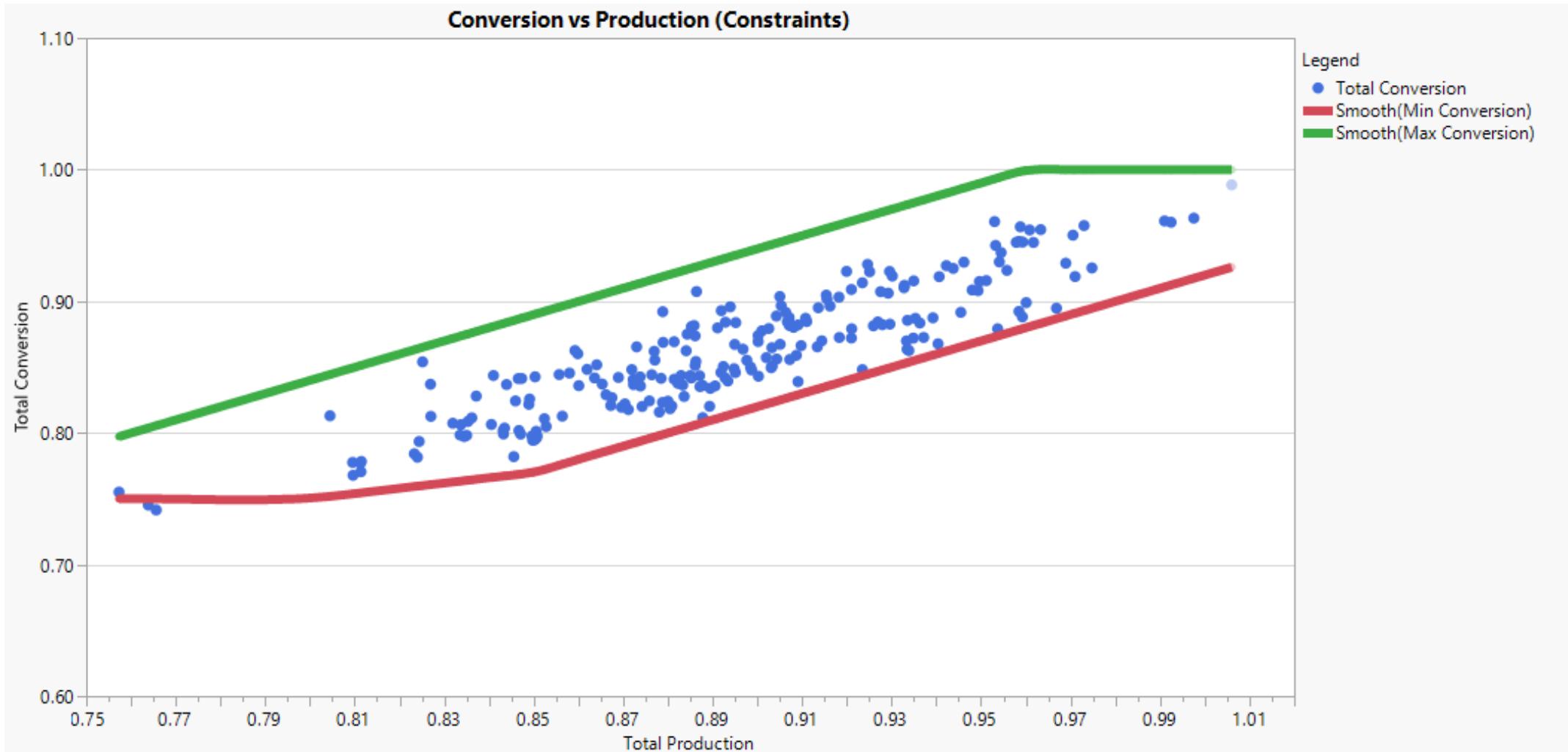
How can we further improve our insights? a cost-effective DOE on an industrial process?

- The idea ... leverage as much as possible historical data:
- Method 1: Augment "most reliable" historical datapoints with Space Filling Design
- Method 2: Reduce experiment cost of Advanced Design by leveraging historical data

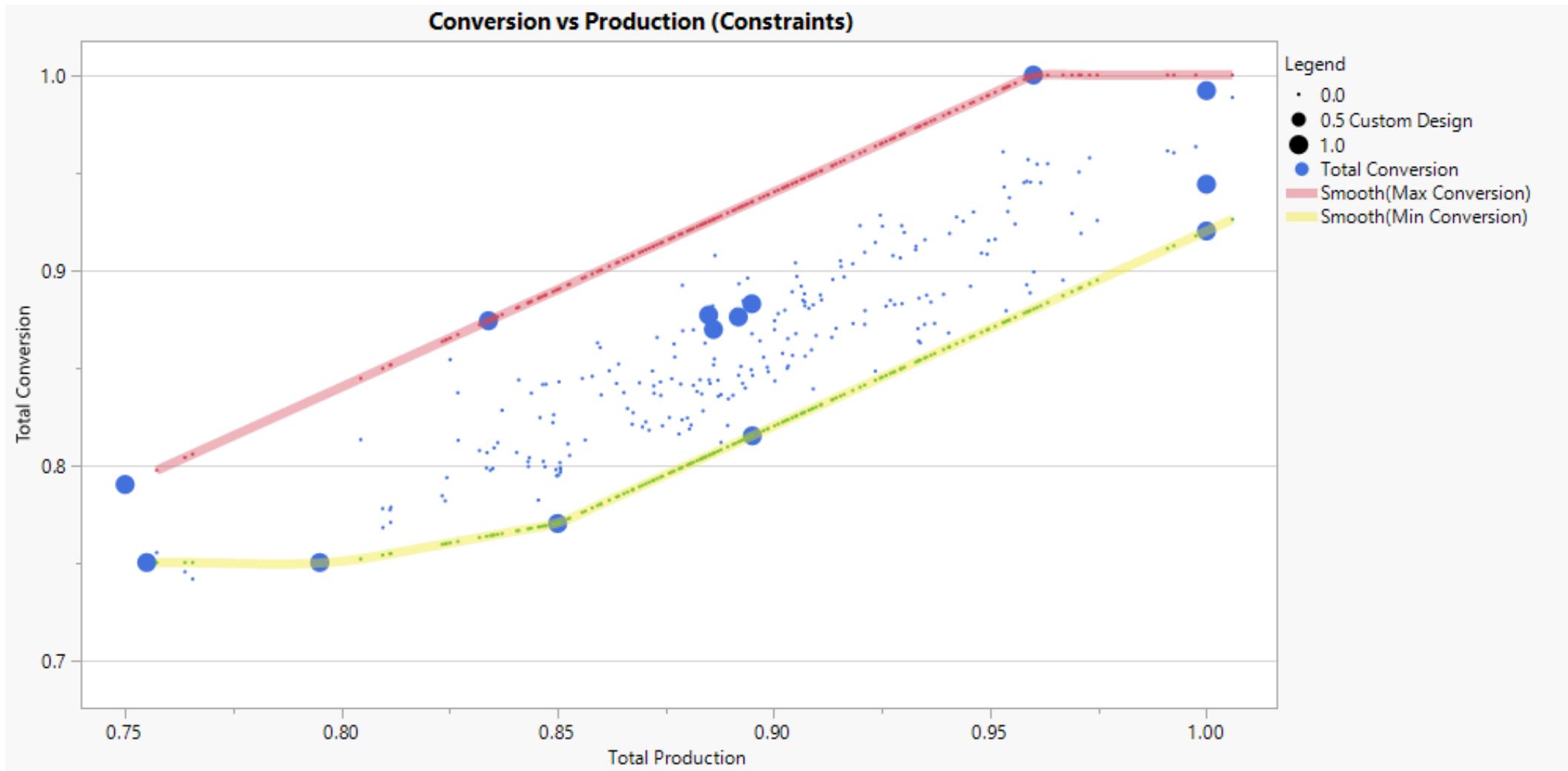


Are we missing information?

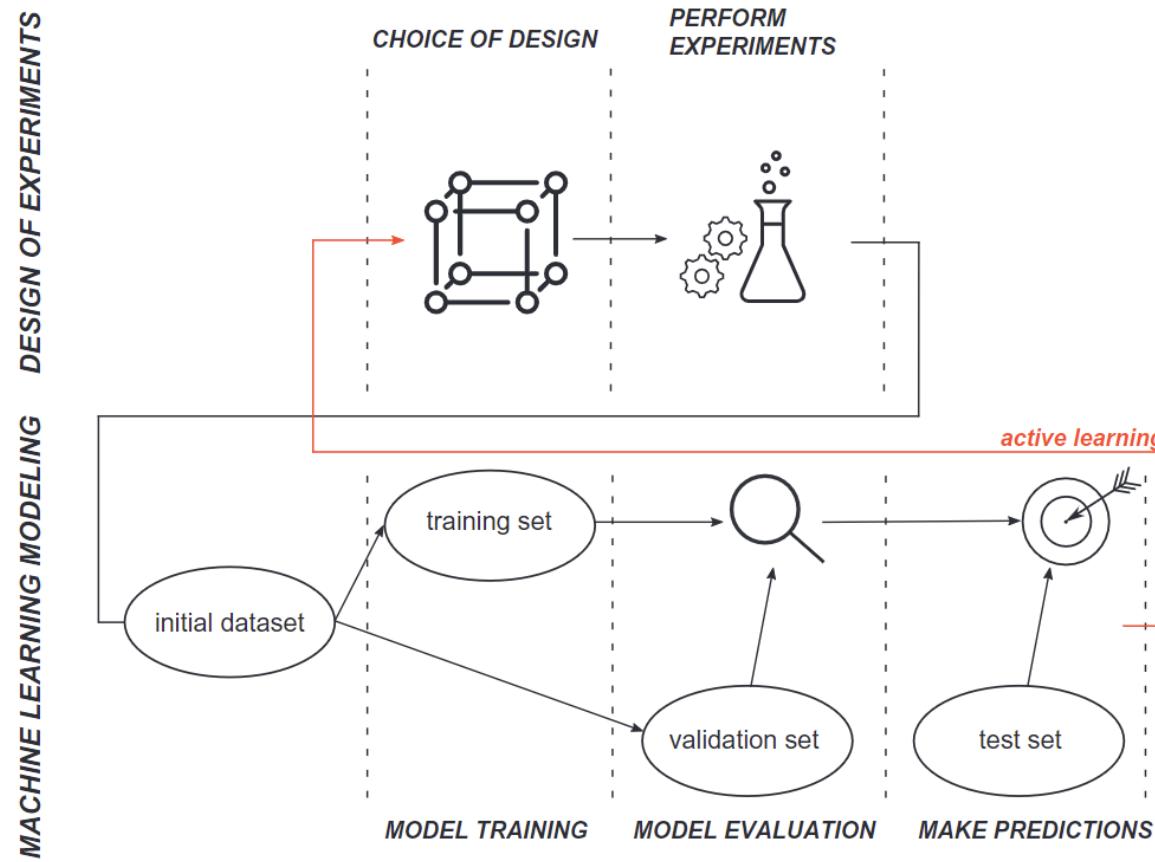
Historical data does not cover the feasible space evenly



Space Filling DOE with historical data as previous experiment



Where to learn more about DoE and Machine Learning



Source:

[Design of Experiments and machine learning for product innovation: A systematic literature review](#) Qual Reliab Engng Int.2022;38:1131–1156

[Design choice and machine learning model performances](#) Qual Reliab Engng Int.2022;38:3357–3378

<https://www.solvay.com/en/brands/soda-solvay>

Scaling up the Use of Machine Learning in Chemical Process Industries - Solvay

Discovery Summit JMP Europe

Sitges (Barcelona), 7th March 2023

David Peig - Digital Champion GBU Soda Ash & Derivatives,
Solvay

Carlos Pérez - Industrial Data Scientist, Solvay

page 1



<https://community.jmp.com/t5/Discovery-Summit-Europe-2023/Scaling-up-the-Use-of-Machine-Learning-in-Chemical-Process/ta-p/572644>

Data-driven approach improve plant performance

Problem

Historically low productivity across site

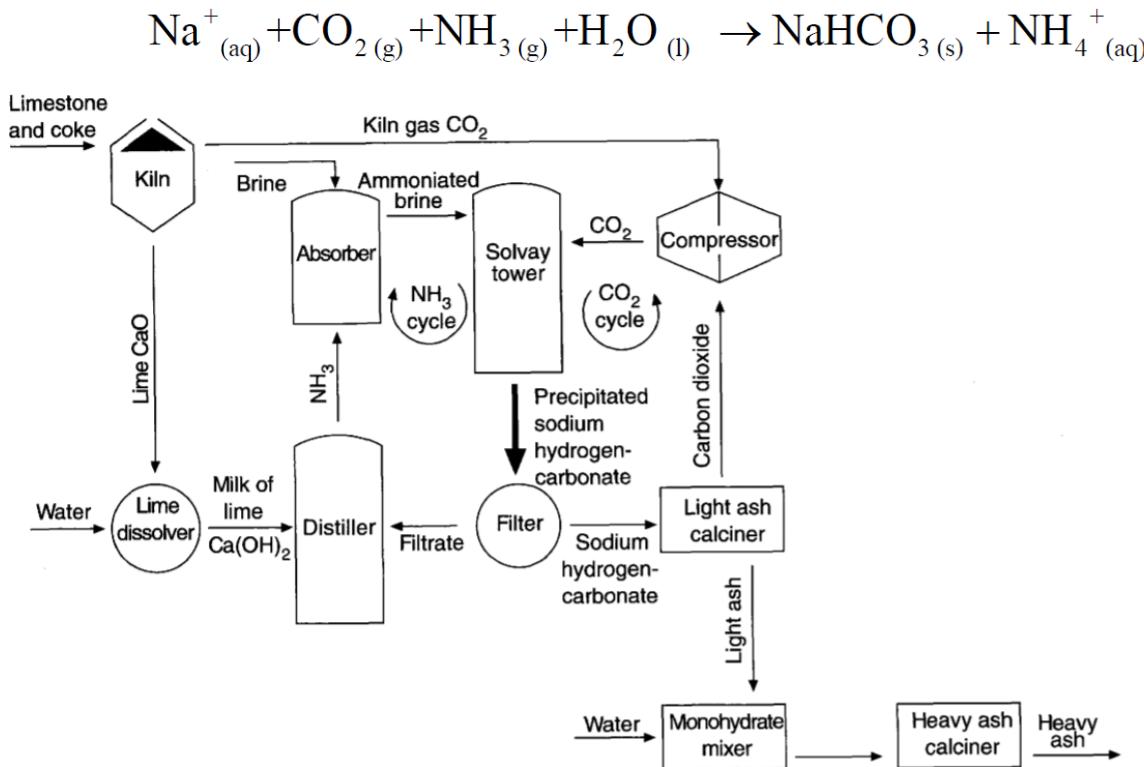


Fig 2 A simplified flow diagram for the Solvay process showing the ammonia and carbon dioxide cycles

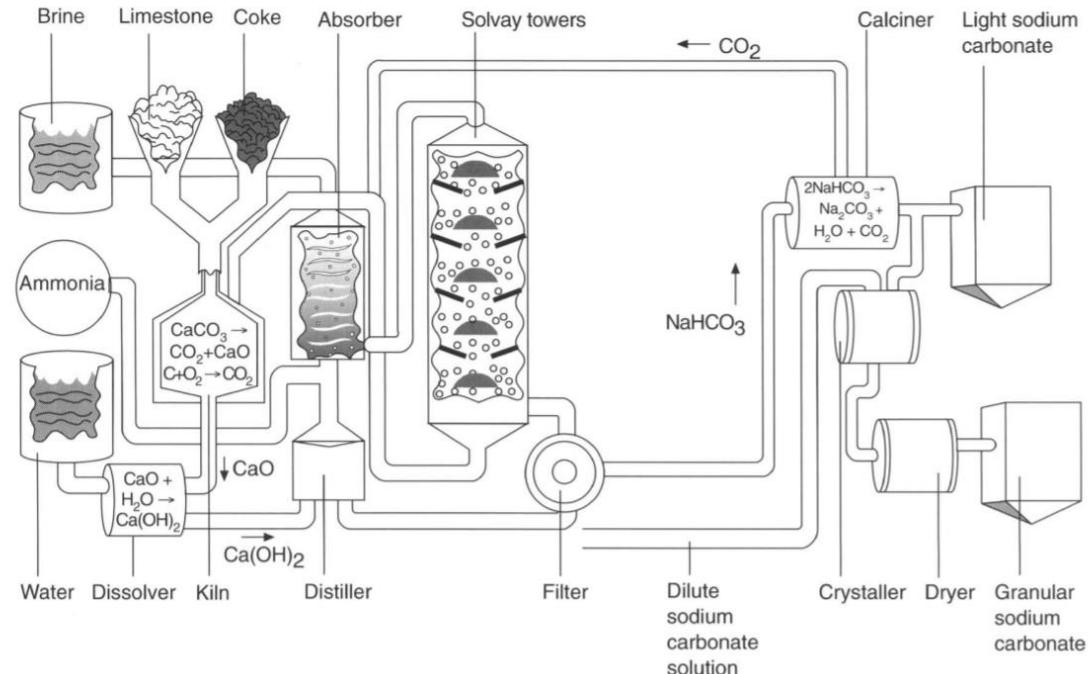


Fig 3 The Solvay process

Three different products are formed from this process, light sodium carbonate (light ash), granular sodium carbonate (heavy ash) and refined sodium hydrogencarbonate.

The main by-product is calcium chloride. A little of this can be sold for use in refrigeration, curing concrete and as a suspension in oil drilling. The bulk of it is disposed of in the nearby river Weaver.

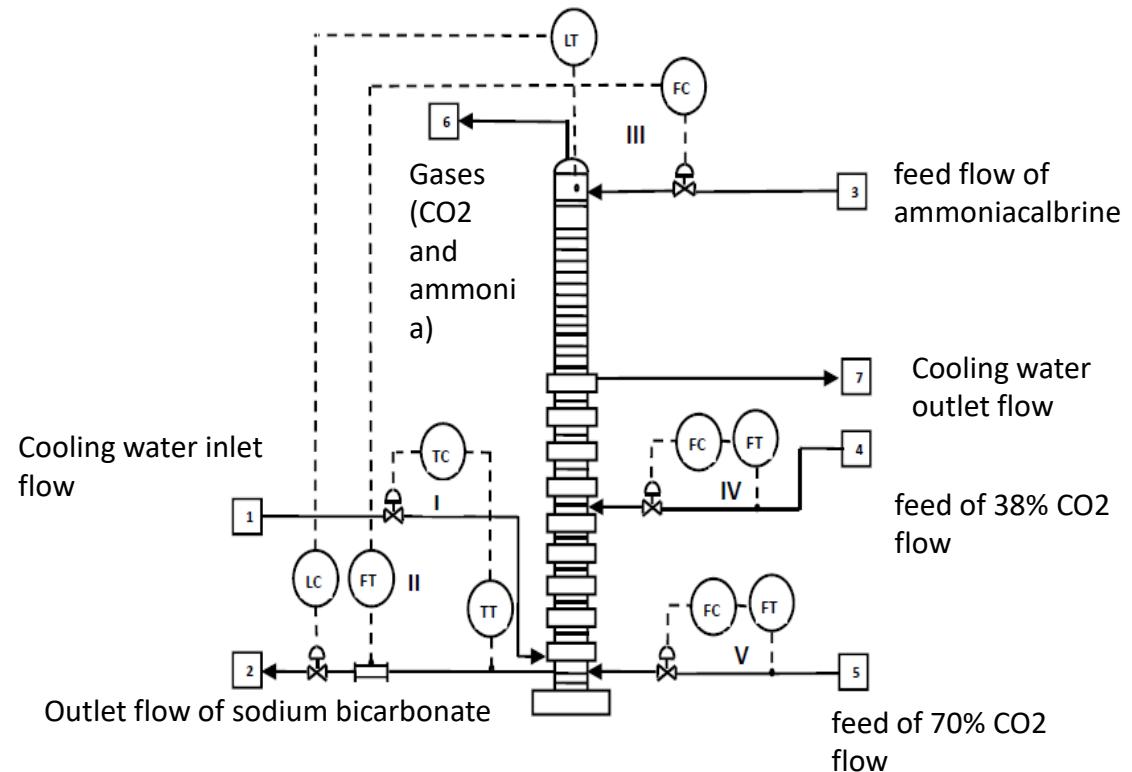
Data-driven approach improve plant performance

Context

- Plant structure across sector makes difficult to track the performance of individual equipment to track root-cause
- There is a lack of measurements for key process variables

What was done?

- Screening of multiple variables and select the most important ones to explain the variability of the target (yield)
- To visualize the long-term variability of the target and its relationship between the most important variables
- Data-driven predictive model on the ratios of key variables
- Optimization on data-driven model to recommend best setpoint across equipment



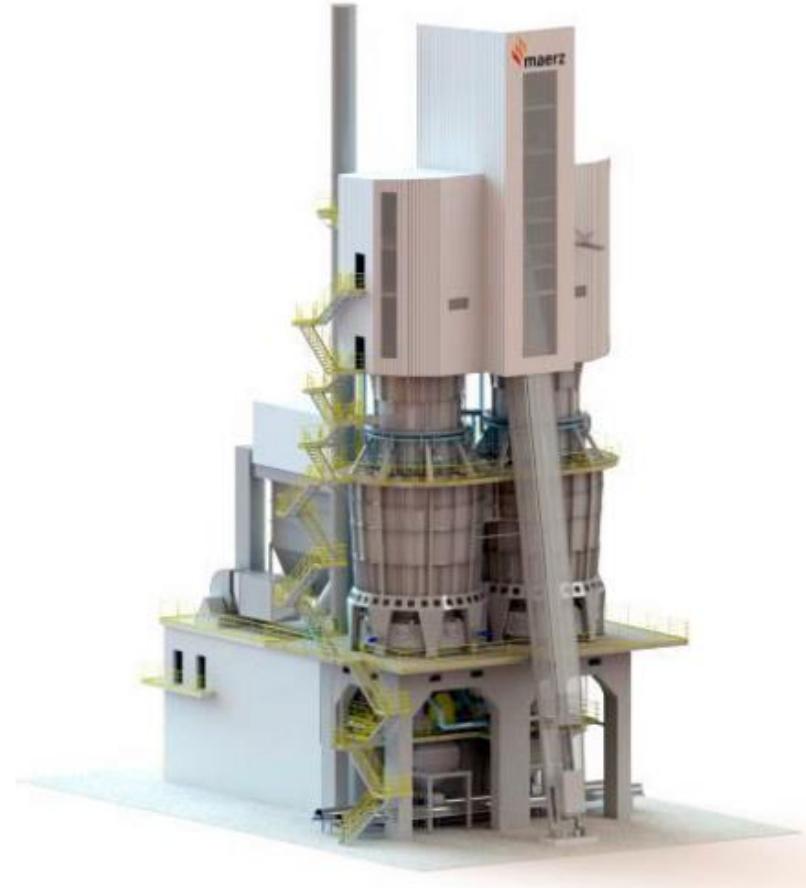
Source:

http://www.chem.ubbcluj.ro/~studiachemia/issues/chemia2017_4/tom2/01Cristea_221_229.pdf

Debottlenecking of limekilns

Problem

Energy-intensive process cannot effectively respond to the changes in raw materials, fuel variability, and setpoint of the plant leading to losses in productivity



Source:

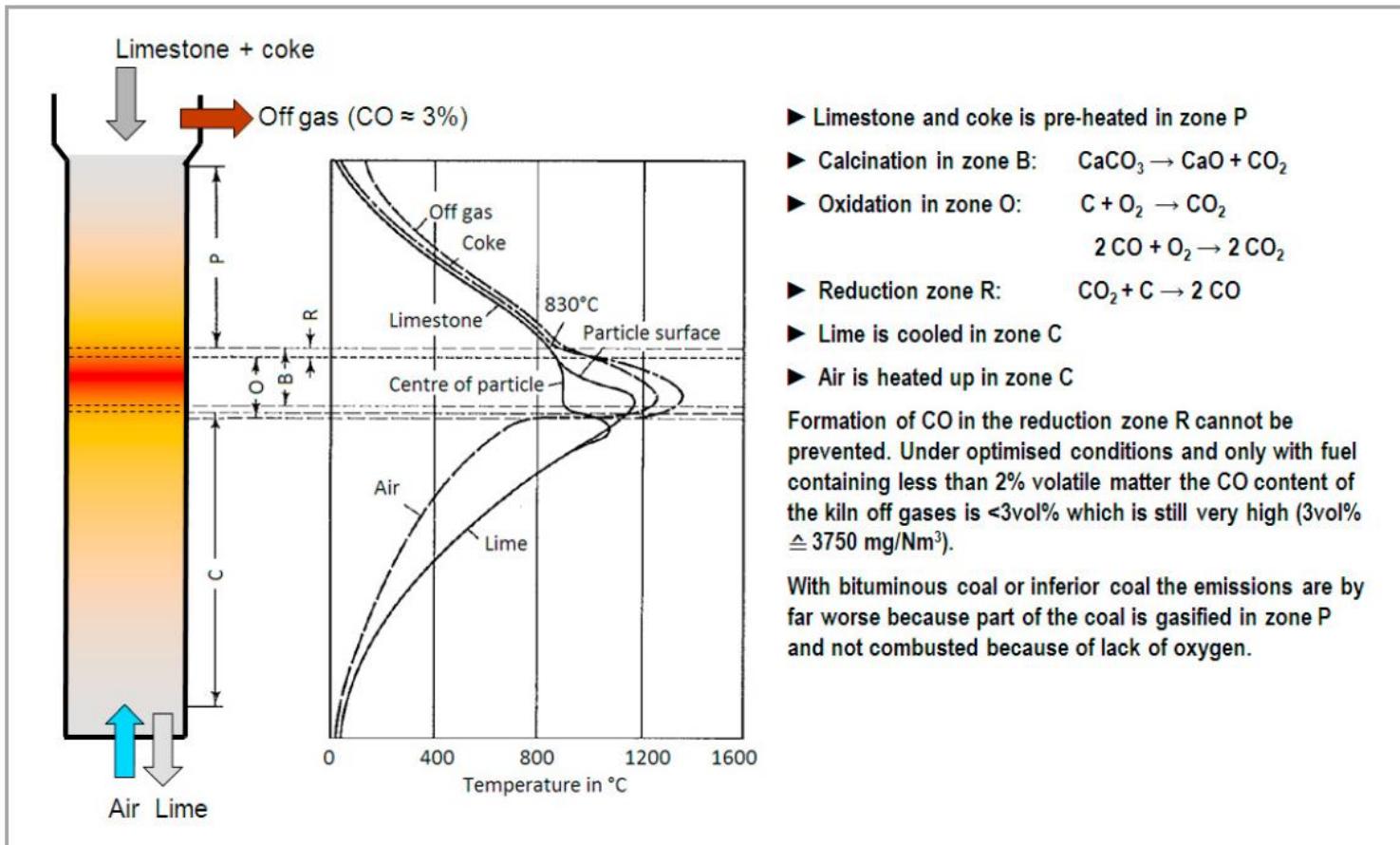
<https://www.sciencedirect.com/science/article/pii/S1876610217327248>

137

Debottlenecking of limekilns

Context

- The cooking of the limestone has a temperature profile that is optimal when stoichiometry is satisfied
- Measurement of the key process variables is challenging.



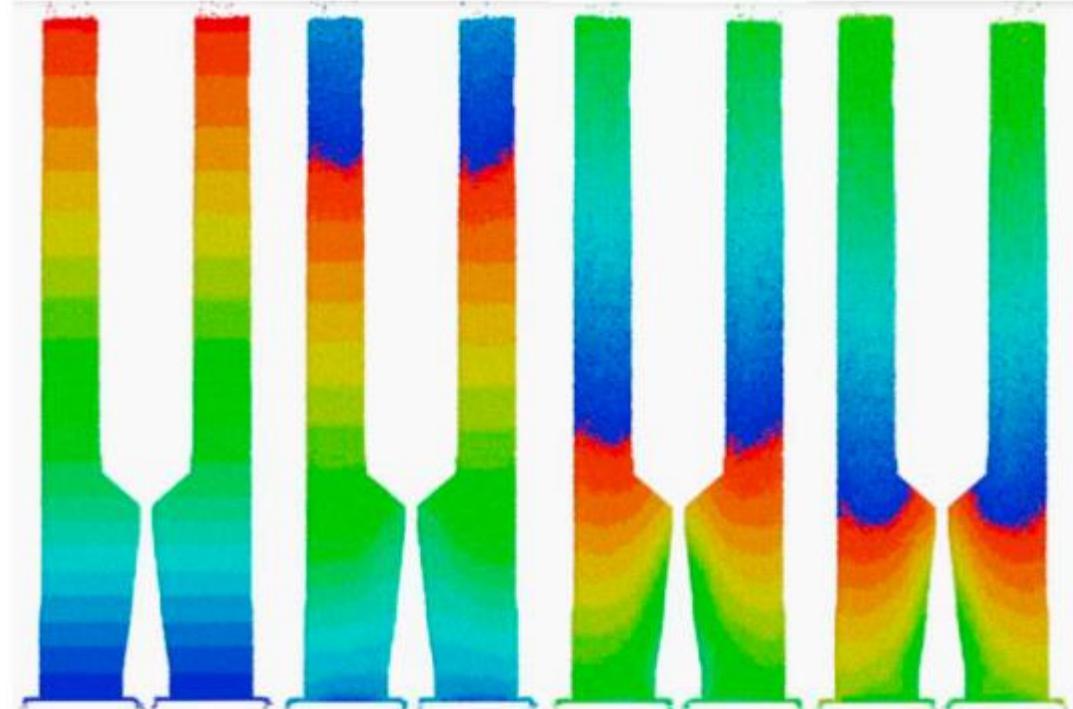
Source:

<https://www.sciencedirect.com/science/article/pii/S1876610217327248>

Debottlenecking of limekilns

What was done?

- New measurement of key process variables allow us to understand behaviour of kiln
- Screening of process variables confirms main contributions of variability
- Soft sensor provides useful KPI
- New control system design takes into account the insights from the data



Source:

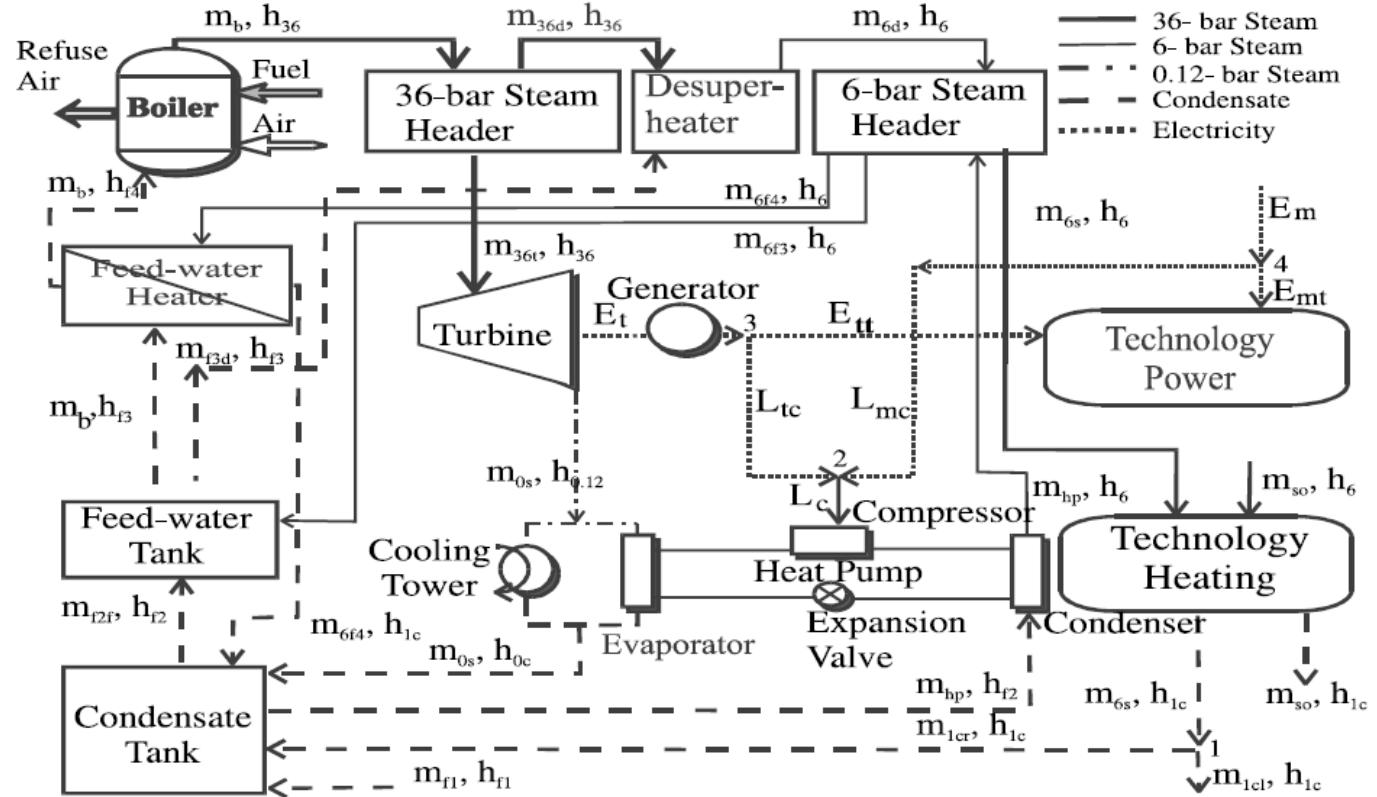
<https://www.sciencedirect.com/science/article/pii/S1876610217327248>

139

Leveraging data-driven model in optimization of steam networks

Problem

Energy-intensive steam network cannot effectively respond to the changes in prices and demand causing financial losses and negative environmental impact.



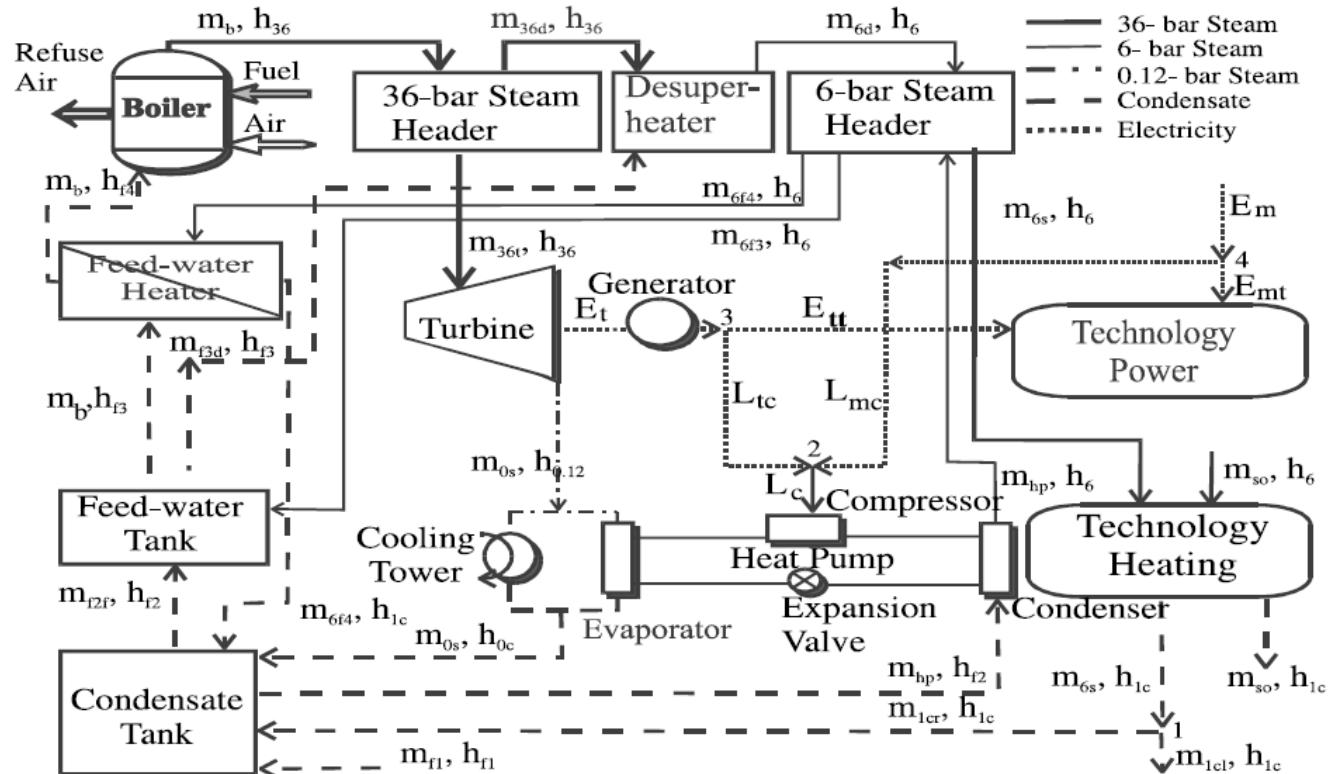
Source:

<https://www.sciencedirect.com/science/article/pii/S0196890401000280>

Leveraging data-driven model in optimization of steam networks

Context

- Steam networks are widely used in the manufacturing environment
- Producers, consumers, electricity generators, letdown valves and cogeneration are some of the elements that need to work together to obtain the best network configuration.
- Equipment behavior can be approximated using data-driven models to account for main effects



Source:

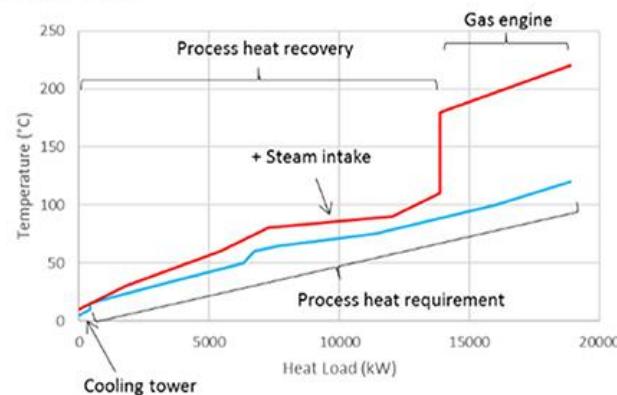
<https://www.sciencedirect.com/science/article/pii/S0196890401000280>

Leveraging data-driven model in optimization of steam networks

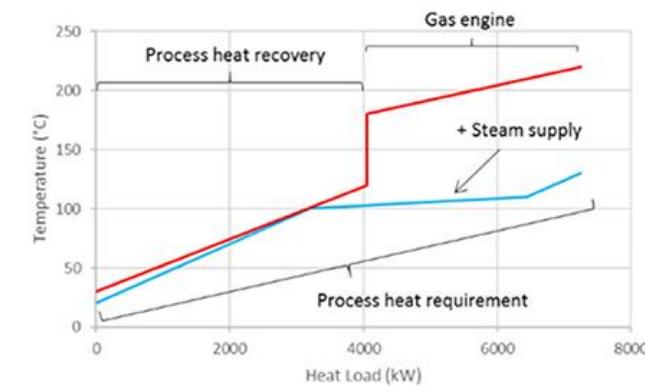
What was done?

- Full description of objective and constraints was done to define an optimization problem as in the literature
- Plant gets recommendations to minimize cost and environmental impact

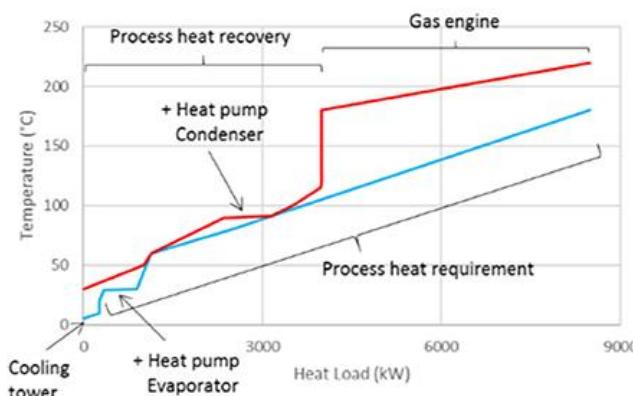
A Process A



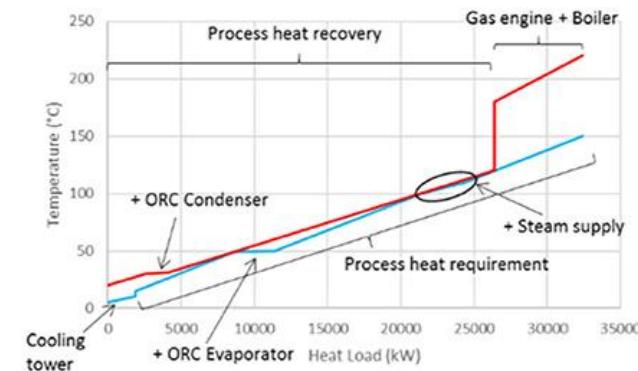
B Process B



C Process C



D Process E



Source:

<https://www.frontiersin.org/articles/10.3389/fenrg.2020.00049/full>

142

Troubleshooting in evaporator

Problem

Production losses in evaporators with chemical additives not helping

Recurring events, main KPI subject to operator's opinion



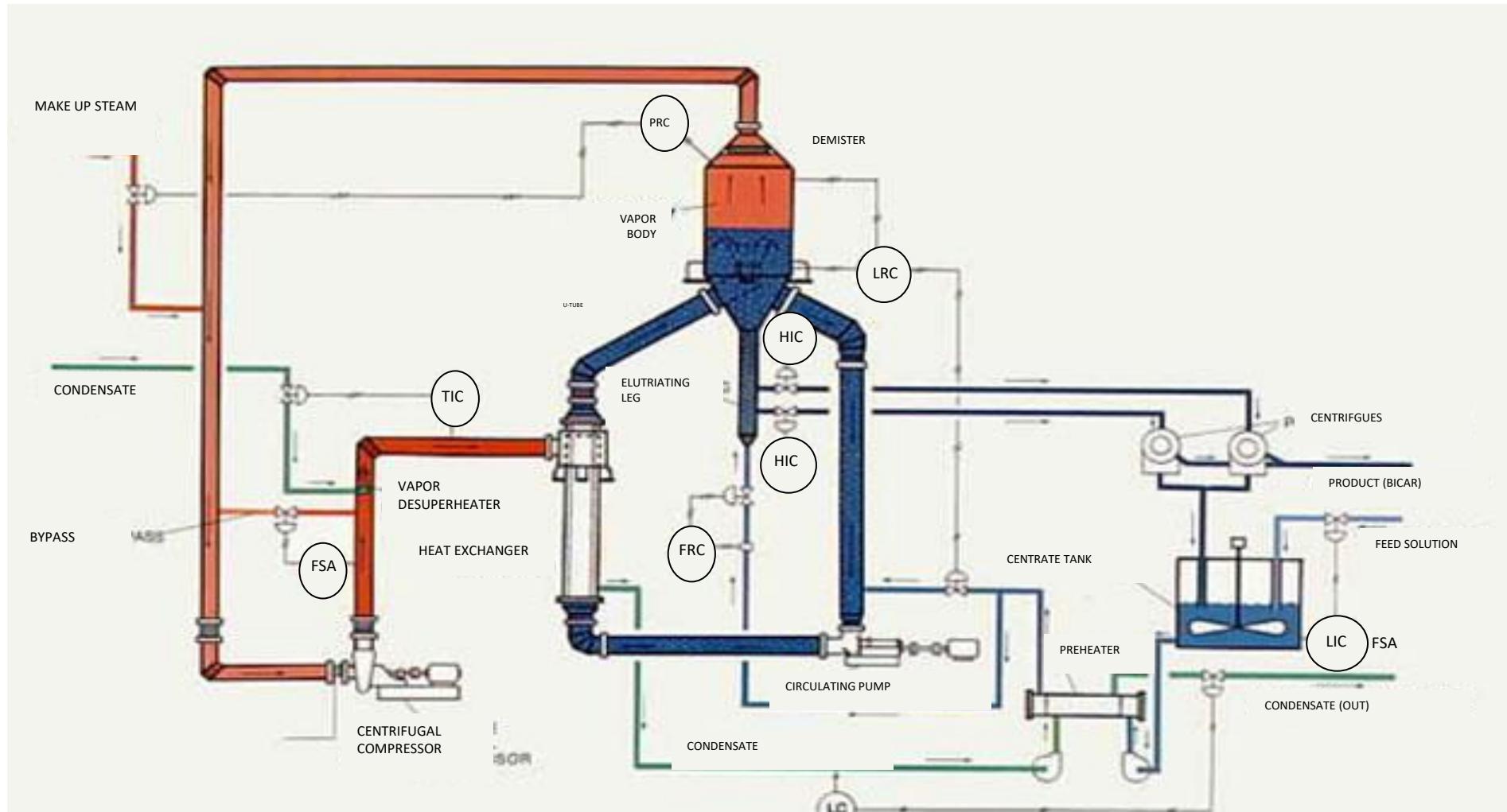
Source:

<https://www.whiting.ca/evaporators-whiting-equipment-canada-inc/>
<https://www.veoliawatertech.com/en/expertise/industries-we-cover/soda-ash>

Troubleshooting in evaporator

Solution

- Screening found key contributors to process anomalies leading the elimination of events
- Control strategy was implemented to automate the addition of additives



Source:

<https://www.whiting.ca/evaporators-whiting-equipment-canada-inc/>

Thank you!

Dr. Mattia Vallerio

Manufacturing Excellence Site Manager at Solvay and
Advanced Process Control Specialist [\[in\]](#)

Dr. Carlos Perez

Industrial Data Scientist at Solvay and
Optimization Specialist [\[in\]](#)

Dr. Francisco Navarro

Sr. Data Science Training Lead at IFF and
Visiting Researcher at Imperial College London [\[in\]](#)