# Industrial data science: a review of machine learning applications for the chemical and process industries

**Max Mowbray**

*max.mowbray@manchester.ac.uk*

**Department of Chemical Engineering**

**The University of Manchester**

# 1. Overview

# Overview

**Work conducted collaboratively between The University of Manchester, Imperial College London and Solvay**

- Review paper into major application and challenges of Machine Learning (ML) in process industries [1]

**Considered contributions and challenges in application across the hierarchical control structure**

- Focus on process level

- Paper provides discussion on upper-level decision functions
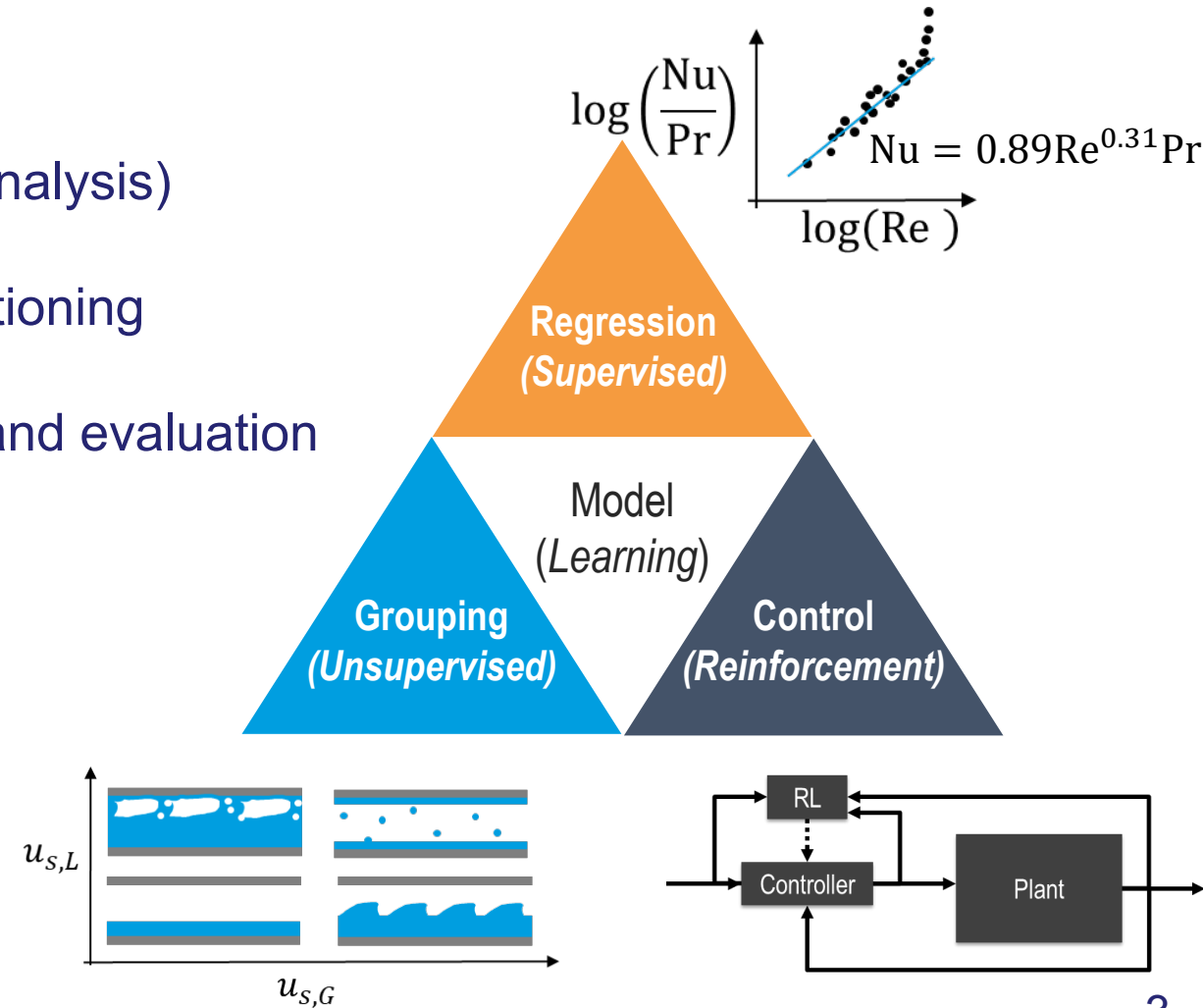


[1] Mowbray *et al.* (2022a)    2

# Overview

**Machine Learning has been widely applied within process systems**, historically under a different disguise

- Feature selection and engineering (dimensional analysis)

- Data pre-processing (signal processing) and partitioning

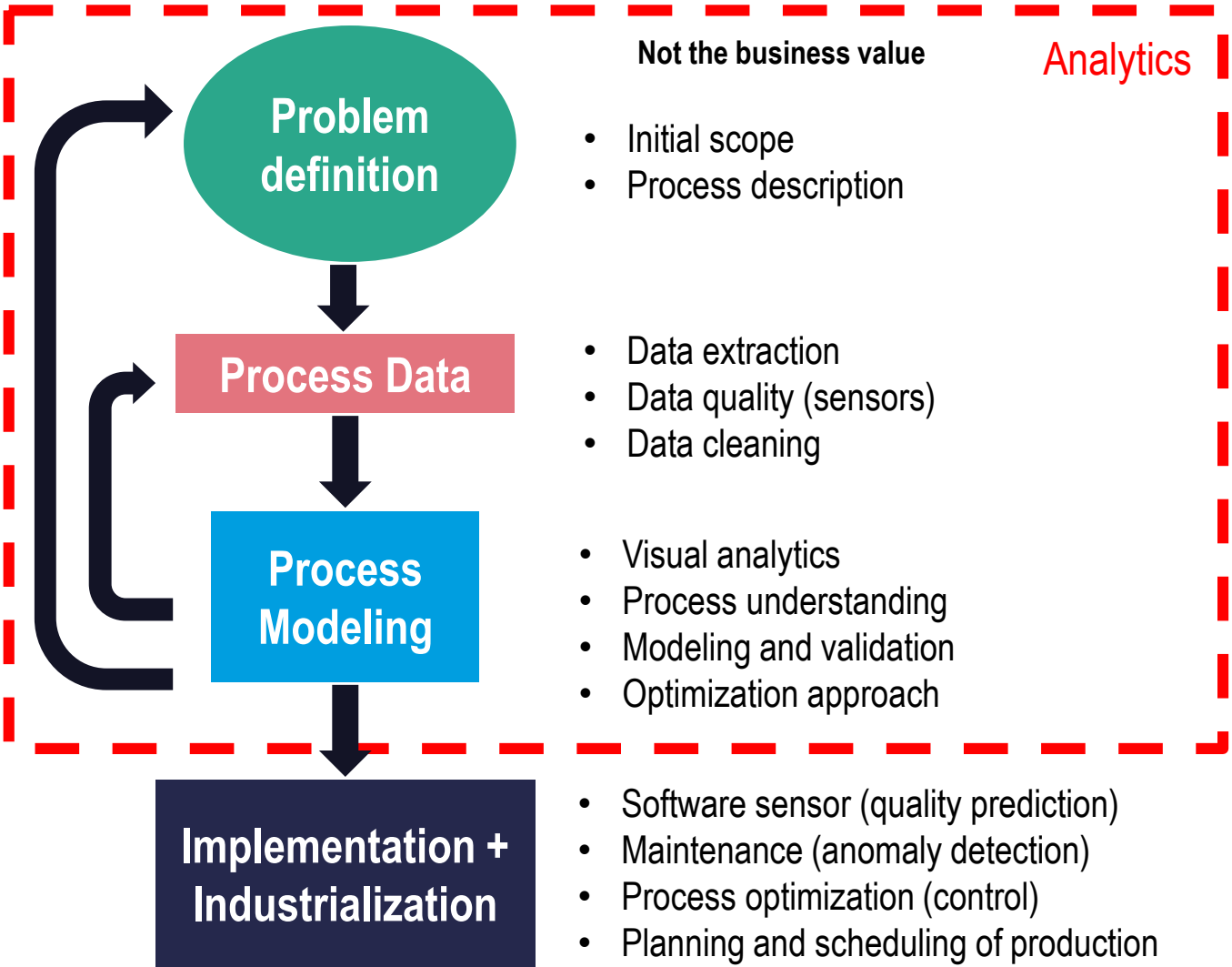- Learning (statistical estimation and optimisation) and evaluation

**Process engineering cover all forms of analytics:**

- Descriptive, diagnostic, predictive and prescriptive

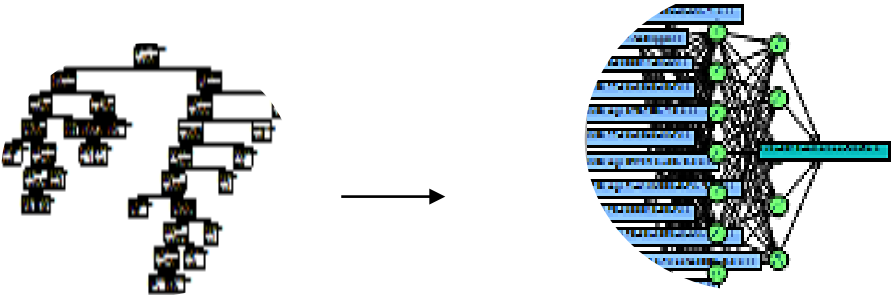- Increasingly flexible model classes, handling uncertainty and data visualisation techniques

# Overview

## Industrial data analytics and science work flow



Not the business value — Analytics

**Problem definition**
- Initial scope
- Process description

**Process Data**
- Data extraction
- Data quality (sensors)
- Data cleaning

**Process Modeling**
- Visual analytics
- Process understanding
- Modeling and validation
- Optimization approach

**Implementation + Industrialization**
- Software sensor (quality prediction)
- Maintenance (anomaly detection)
- Process optimization (control)
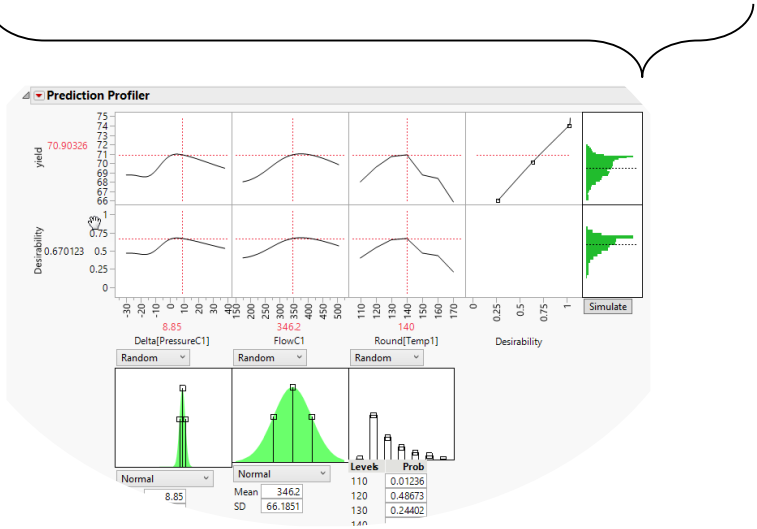- Planning and scheduling of production

## ML tools are key when knowledge is limited, but correlation is present



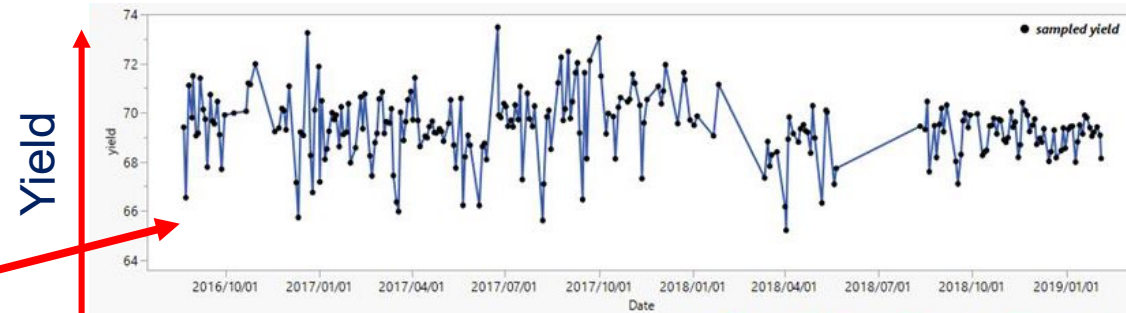Decision tree model(s)    Neural Network

2. Case Study 1: Diagnosing yield variation in a distillation column
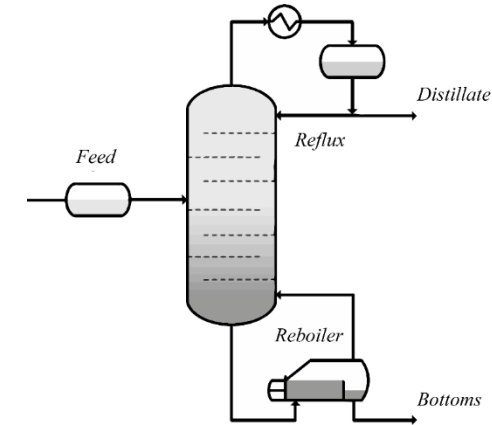
# Diagnosing yield variation

**Problem definition first consider trends in yield: variation driven by season or dynamics**

- No clear seasonal trend

  <span style="color:red">High frequency, noisy variation</span>



Time (months)

- Existing data collected at supposed steady state



$\tau:$  $x_{t-k}$  $x_{t-k+1}$  $\cdots$  $x_t$

History uncorrelated

History correlated

# Diagnosing yield variation

**Problem definition first consider trends in yield: variation driven by season or dynamics**
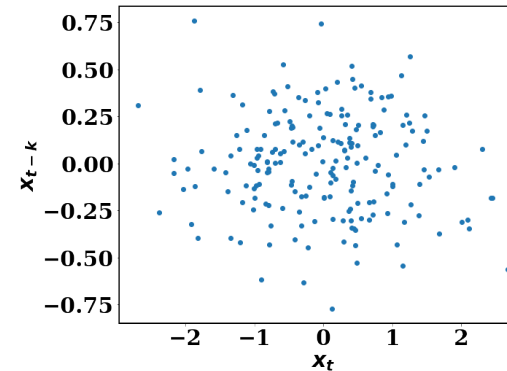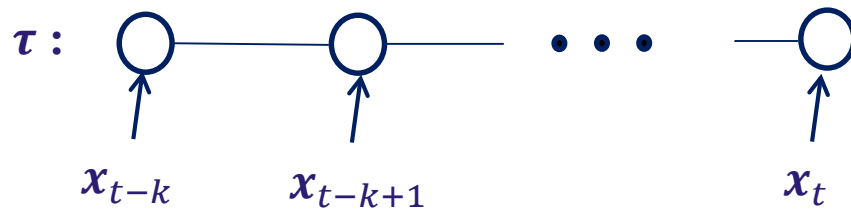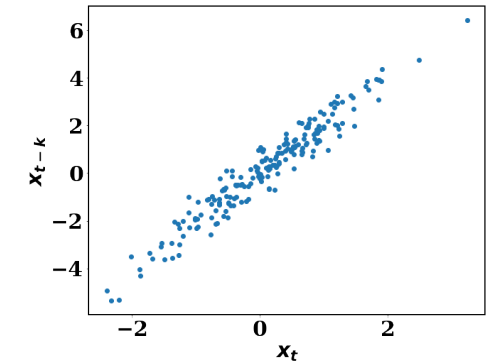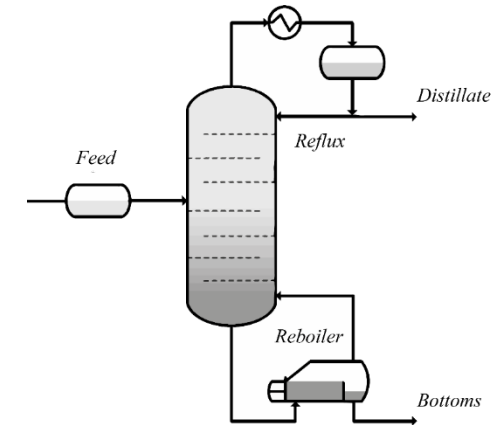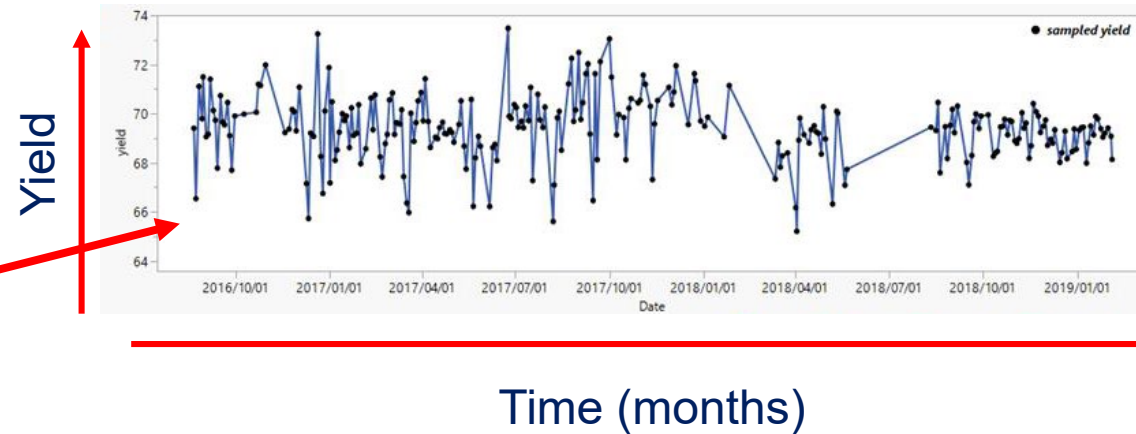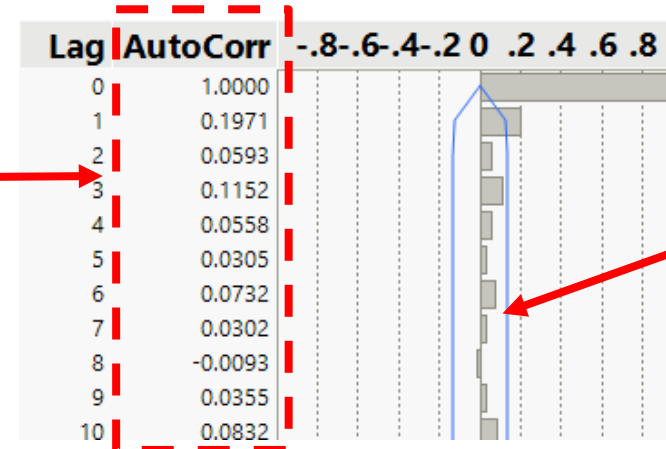
- No clear seasonal trend

High frequency, noisy variation

Yield

Time (months)

- Existing data collected at supposed steady state

Autocorrelation

Dynamics not clearly driving variation

$\tau:$ $x_{t-k}$ $x_{t-k+1}$ $\cdots$ $x_t$

| Lag | AutoCorr | -.8 -.6 -.4 -.2 0 .2 .4 .6 .8 |
|-----|----------|-------------------------------|
| 0 | 1.0000 | |
| 1 | 0.1971 | |
| 2 | 0.0593 | |
| 3 | 0.1152 | |
| 4 | 0.0558 | |
| 5 | 0.0305 | |
| 6 | 0.0732 | |
| 7 | 0.0302 | |
| 8 | -0.0093 | |
| 9 | 0.0355 | |
| 10 | 0.0832 | |

$c << 1$ history not correlated (IID assumption)

# Diagnosing yield variation

**Data pre-processing; first understand the data at hand**

- Data is high dimensional sensor measurements

  - Quantify correlation between two sensors' data

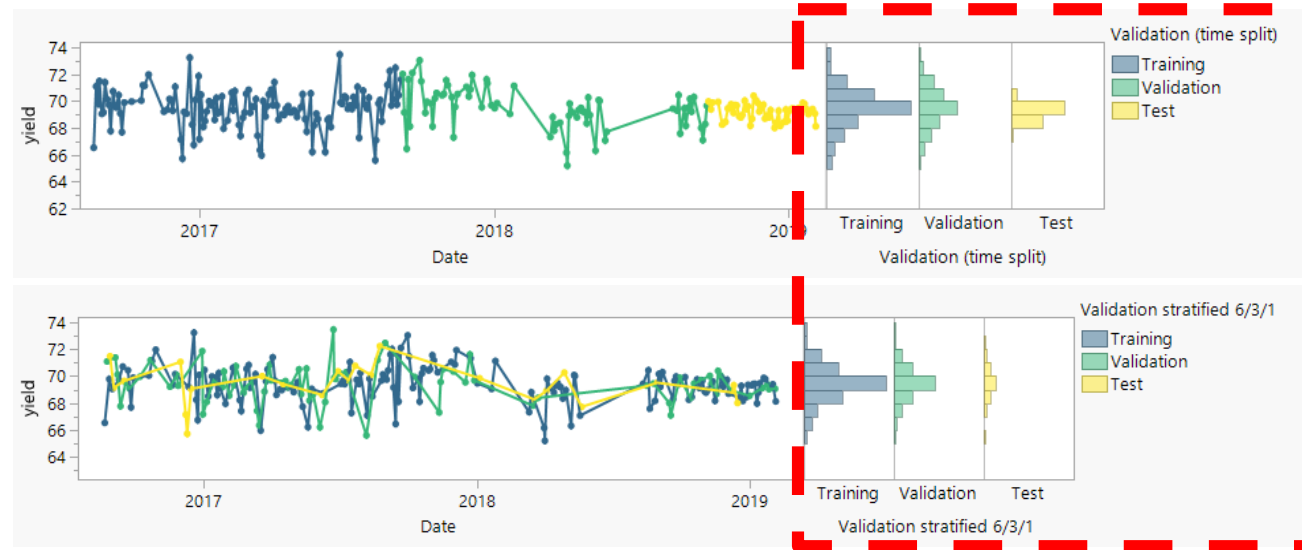**Data visualisation powerful to quickly understand the data**

- Pearson's correlation coefficient

- 2D visualisation of very high dimensional data

  - The more yellow the pixel the higher the correlation

- Mutlicollinear data causes issues in model construction

  - Dimensionality reduction

# Diagnosing yield variation

**Data pre-processing: data partitioning strategy is key for ensuring adequate validation of model**
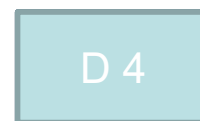


- Training
- Validation
- Test

- Subsampling time-series considers the data to I.I.D. (dynamics are not driving variability, rarely the case!)

- K-fold cross validation designs can ensure against mismatch in data distributions

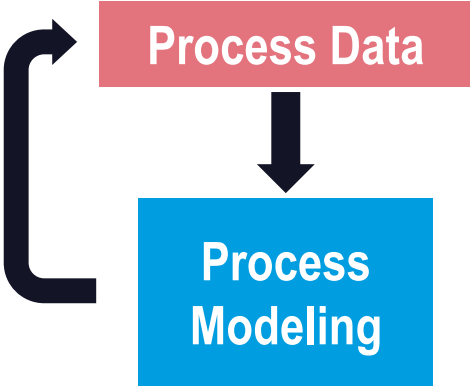**Data processing: k-fold cross validation**

- K-fold cross validation:

| D 1 | D 2 | D 3 | D 4 | D 5 | D 6 |

Training Set

Validation Set

# Diagnosing yield variation

**Generate knowledge by modelling; screen the data to select variables**

Process Data → Process Modeling (iterative loop)

Identifying significant variables may be an iterative process:

2. Repeat and analyse new model

**Updated random forest model: feature importance**

**Column Contributions**

| Term | Number of Splits | SS | | Portion |
|---|---|---|---|---|
| FlowC1 | 2 | 123.846558 | | 0.5811 |
| Temp1 | 1 | 89.2786595 | | 0.4189 |
| Shuffle[yield] | 0 | 0 | | 0.0000 |
| OC1 | 0 | 0 | | 0.0000 |
| PressureC1 | 0 | 0 | | 0.0000 |
| Random Uniform Noise | 0 | 0 | | 0.0000 |
| TempC9 | 0 | 0 | | 0.0000 |
| Temp6 | 0 | 0 | | 0.0000 |
| Random Normal Noise | 0 | 0 | | 0.0000 |
| FlowC2 | 0 | 0 | | 0.0000 |
| TempC1 | 0 | 0 | | 0.0000 |
| TempC3 | 0 | 0 | | 0.0000 |
| Temp4 | 0 | 0 | | 0.0000 |

Flow and Temp. important

Knowledge generation →



a) Partition tree — scatter plot of Temp1 vs FlowC1, with predicted yield colour scale (67.5–71.0). Low yield region (high Temp1) and High yield region (high FlowC1) marked.

11

# Diagnosing yield variation

**Generate knowledge by modelling; screen the data to select variables**

Process Data

Process Modeling

Identifying significant variables may be an iterative process:

2. Repeat and analyse new model

**Updated random forest model: feature importance**

**Column Contributions**

| Term | Number of Splits | SS | | Portion |
|---|---|---|---|---|
| FlowC1 | 2 | 123.846558 | | 0.5811 |
| Temp1 | 1 | 89.2786595 | | 0.4189 |
| Shuffle[yield] | 0 | 0 | | 0.0000 |
| OC1 | 0 | 0 | | 0.0000 |
| PressureC1 | 0 | 0 | | 0.0000 |
| Random Uniform Noise | 0 | 0 | | 0.0000 |
| TempC9 | 0 | 0 | | 0.0000 |
| Temp6 | 0 | 0 | | 0.0000 |
| Random Normal Noise | 0 | 0 | | 0.0000 |
| FlowC2 | 0 | 0 | | 0.0000 |
| TempC1 | 0 | 0 | | 0.0000 |
| TempC3 | 0 | 0 | | 0.0000 |
| Temp4 | 0 | 0 | | 0.0000 |

Flow and Temp. important

Knowledge generation



Colored by yield (kernel density)

12

# Diagnosing yield variation

**Random forests are interpretable and good at screening, but what about predictive accuracy?**

**Process Modeling**

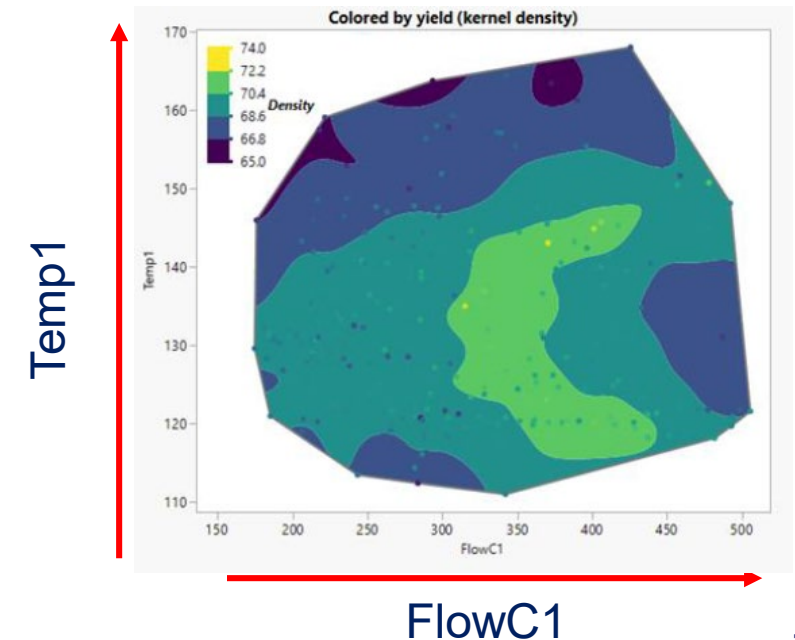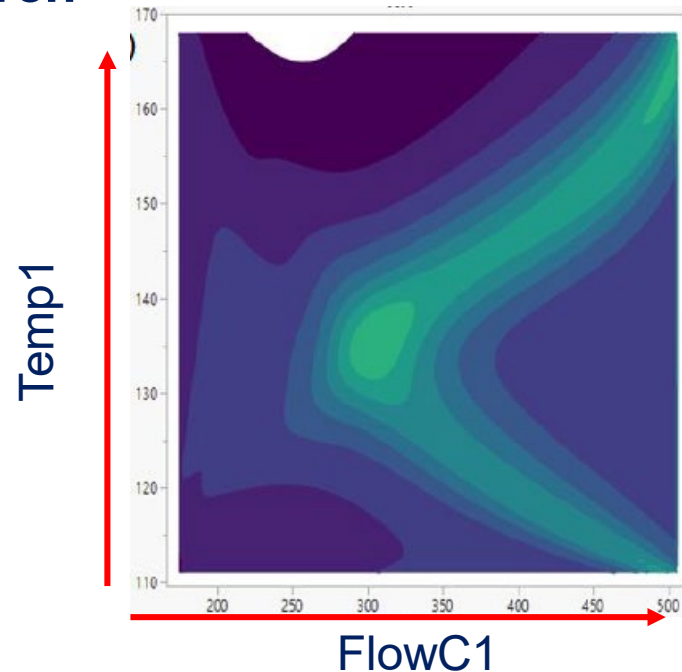- Having identified important sensors, we can search over model classes and structures

**Manual, random or automated search**

- Neural networks are a go-to

- Improved generalisation accuracy

- Smooth relationship identified

# Summary: Case Study 1

**Yield variation: diagnostic summary**

- Two step approach:

Important process variables

Decision tree model(s)

JMP Add ins

Flexible function approximation

Neural Network

**Synthetic dataset with known ground truth and assumed steady state**

- Conceptual demonstration of correlation analysis and screening

- Future work to identify dynamics

- **Real world data key to further exploration of ML**

Implementation + Industrialization

Application 1: demonstration

Truth (Rosenbrock)

Temp1

FlowC1

3. Case Study 2: Data-driven soft-sensing of product quality in batch processing

# Soft sensing

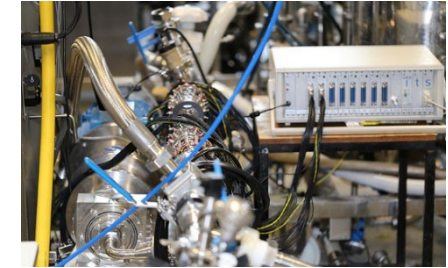**Predictive analytics for estimating viscosity in a batch process handling non-Newtonian fluid** [4]

- Real-time PAT is expensive and difficult to retrofit

- Offline measurement slow and has standard error

**Emulsification process (batch processes) for personal care**

**product manufacturing**

- 28 online sensors (temperature, pressure, flowrate)

- Sensor data recorded once per second (in total 2 hours)

- Available datasets: 30 batches × 28 sensors × 7673 time steps

[4] Mowbray *et al.* (2022b)

# Soft sensing

**Challenge 1:** Identifying the spatio-temporal trajectories that influence product quality



Critical time regions

Critical sensors

## Projection to latent structures (PLS)

Importance of each variable

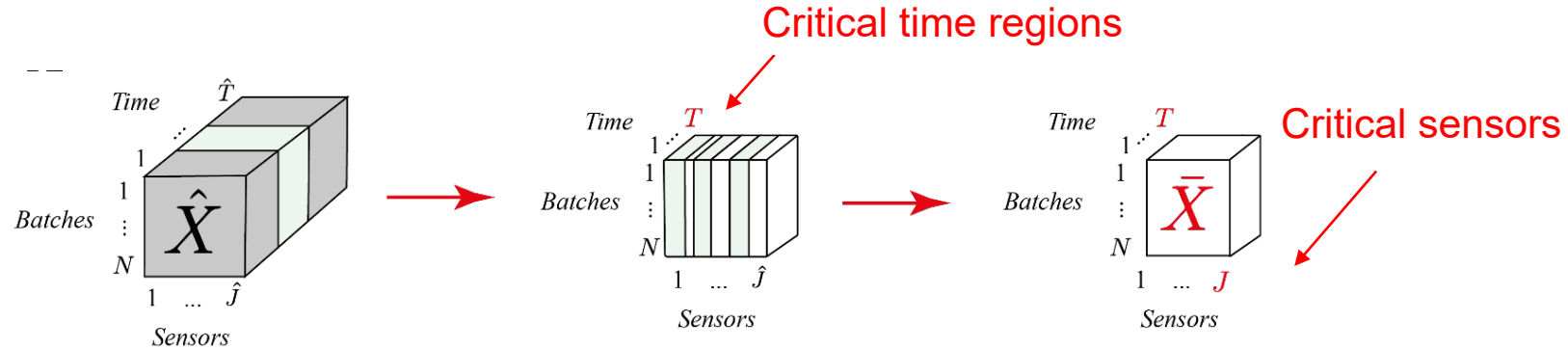$$z_1 = \theta_{1,1} \cdot x_1 + \theta_{1,2} \cdot x_2 + \cdots + \theta_{1,n} \cdot x_n$$

$$y = a_1 \cdot z_1 + e$$

- Screen data to find correlations relevant to prediction

## Solution: Loadings analysis of PLS models [2]

- Construct a PLS model for each sensor in time, $\theta^j \in \mathbb{R}^{n_z \times \hat{T}}$



- Construct a PLS model for the sensors at each timestep, $\theta^t \in \mathbb{R}^{n_z \times \hat{J}}$



17

# Soft sensing

**Challenge 1:** Identifying the spatio-temporal trajectories that influence viscosity



- Analyse the loadings of the PLS models
- Heat maps are interpretable

*Critical time region: 5-15 minute period from 2 hours ; Critical sensors: 8 from 28 sensors*

## Challenge 2: Feature extraction and dimensionality reduction

Data for soft-sensor construction

Targets to predict

$$f(\cdot)$$

Product qualities

## Partial Solution: Multiway methods

- Representation that allows for extraction of important information

# Soft sensing

## Challenge 2: Feature extraction and dimensionality reduction



**Solution:** Feature extraction using multiway PLS or autoencoders*

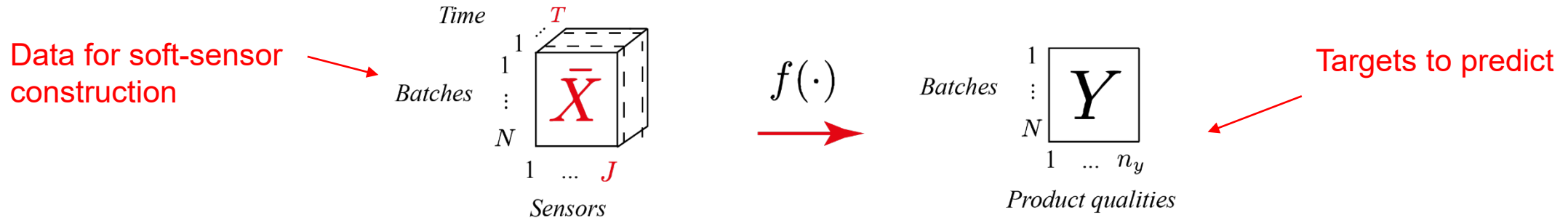- Reduced dimension representation of spatio-temporal trajectories correlated to end-product quality

- Can now identify a map to product quality, $f: \mathbb{R}^{n_z} \to \mathbb{R}$

Number of latent variables enters model structure selection problem

Extracted features expressed in linear subspace

## Challenge 3: Expressing nonlinearity and prediction uncertainty



- Classical regression practice identifies deterministic models
- Our data has moderate noise in the measurement

# Soft sensing

## Challenge 3: Expressing nonlinearity and prediction uncertainty



- Classical regression practice identifies deterministic models

- Our data has moderate noise in the measurement

## Challenge 3: Expressing nonlinearity and prediction uncertainty

### Gaussian Processes (GPs):

- $f_{GP}(\boldsymbol{z}) \sim GP\left(m(\boldsymbol{z}), k(\boldsymbol{z}, \boldsymbol{z}')\right)$

- Exploit statistical relationships in data

    - Bayes' Rule: $y_i \sim p(y_i | \boldsymbol{z}_i^*, \mathcal{D}) = \mathcal{N}(\bar{\mu}, \Sigma)$

- Uncertainty reflects data variation and lack of information

# Soft sensing

## Model structure selection via cross-validation and subsequent model testing

| Dataset | Season | Use | Batches |
|---------|--------|-------|---------|
| A | 1 | Train | 30 |
| B | 2 | Test | 16 |

1. **Cross validation on Dataset A** to identify model structure for each class

→

2. **Test on Dataset B** to evaluate prediction **same variant** but **different season**

## Cross-validation and test results

- GP performs well in validation with average error of 10%

- Uncertainty estimate covers residual – indicating it is reliable

Prediction plots for GP on dataset B

# Summary: Case Study 2

**Data visualisation is an effective step to analyse historical datasets**

- Screen critical time region and sensors (***knowledge informed*** dimensionality reduction)

**Probabilistic machine learning methods are excellent for soft-sensor design**

- A high accuracy soft-sensor provides avenue to monitor

- Reliable uncertainty estimates to guide process engineers

**Potential for industrial application**

- Fast prediction online of critical product quality vs slow offline measurement

- Investigating methodologies to transfer soft-sensors between processes.

# 4. Conclusions

# Conclusion

**Here we present an intuitive data focused framework for problem solving**

- Problem definition; data processing; modelling; implementation

**Machine Learning tools can be used for descriptive, diagnostic and predictive analysis**

- Correlation analysis for quick screening of tags (sensors)

- Mapping identification for steady state behaviour as well as spatio-temporal trajectories to final product qualities

- Future work will consider methodology for identification of dynamic behaviour

**The litmus test of Machine Learning is practical implementation to real processes and data**

# References

[1] Mowbray, M., Vallerio, M., Perez-Galvan, C., Zhang, D., Chanona, A. D. R., & Navarro-Brull, F. J. (2022a). Industrial data science–a review of machine learning applications for chemical and process industries. Reaction Chemistry & Engineering.

[2] Beck, D. A., Carothers, J. M., Subramanian, V. R., & Pfaendtner, J. (2016). Data science: Accelerating innovation and discovery in chemical engineering. *AIChE Journal*, *62*(5), 1402-1416.

[3] Shang, C., & You, F. (2019). Data analytics and machine learning for smart process manufacturing: recent advances and perspectives in the big data era. Engineering, 5(6), 1010-1016.

[4] Mowbray, M., Kay, H., Kay, S., Caetano, P. C., Hicks, A., Mendoza, C., ... & Zhang, D. (2022b). Probabilistic machine learning based soft-sensors for product quality prediction in batch processes. Chemometrics and Intelligent Laboratory Systems, 228, 104616.

## Thank you for listening!

@MaxMowbray3     @DongdaZhang     @AntonioE89     Phil Martin     Francisco Navarro-Brull

28