# Predicting Quality-relevant Variables in Industrial Process via Causality Analysis

Reporter: Liang Cao

Supervisor: Bhushan Gopaluni

Department of Chemical and Biological Engineering,
University of British Columbia

2020.02.14

# Situation

✓ Chemical process usually has thousands of process variables (X) available to predict quality-relevant variables (Y).

✓ How to find the features of X that are important for predicting Y (which features of X helps predict Y) is one of the most important problems in ML, but very messy.

| Egg | Milk | Fish | Wheat | Shellfish | Peanuts | ... |  | Sick? |
|-----|------|------|-------|-----------|---------|-----|--|-------|
| 0 | 0.7 | 0 | 0.3 | 0 | 0 | | → | 1 |
| 0.3 | 0.7 | 0 | 0.6 | 0 | 0.01 | | → | 1 |
| 0 | 0 | 0 | 0.8 | 0 | 0 | | → | 0 |
| 0.3 | 0.7 | 1.2 | 0 | 0.10 | 0.01 | | → | 1 |

✓ We want to know which foods are important for predicting "sick"

# Problem

✓ A common way to do feature selection is compute the correlation between feature values $X_i$ and Y, if the correlation is above certain value, take these features.

✓ Usually gives unsatisfactory results as it ignores variable interactions:
  – Includes irrelevant variables: "Taco Tuesdays".
    • If tacos make you sick, and you often eat tacos on Tuesdays, it will say "Tuesday" is relevant.
  – Excludes relevant variables: "Diet Coke + Mentos Eruption".
    • Diet coke and Mentos don't make you sick on their own, but *together* they make you sick.

✓ To build simpler, more powerful, more interpretable model, we use causality analysis to find causal features.

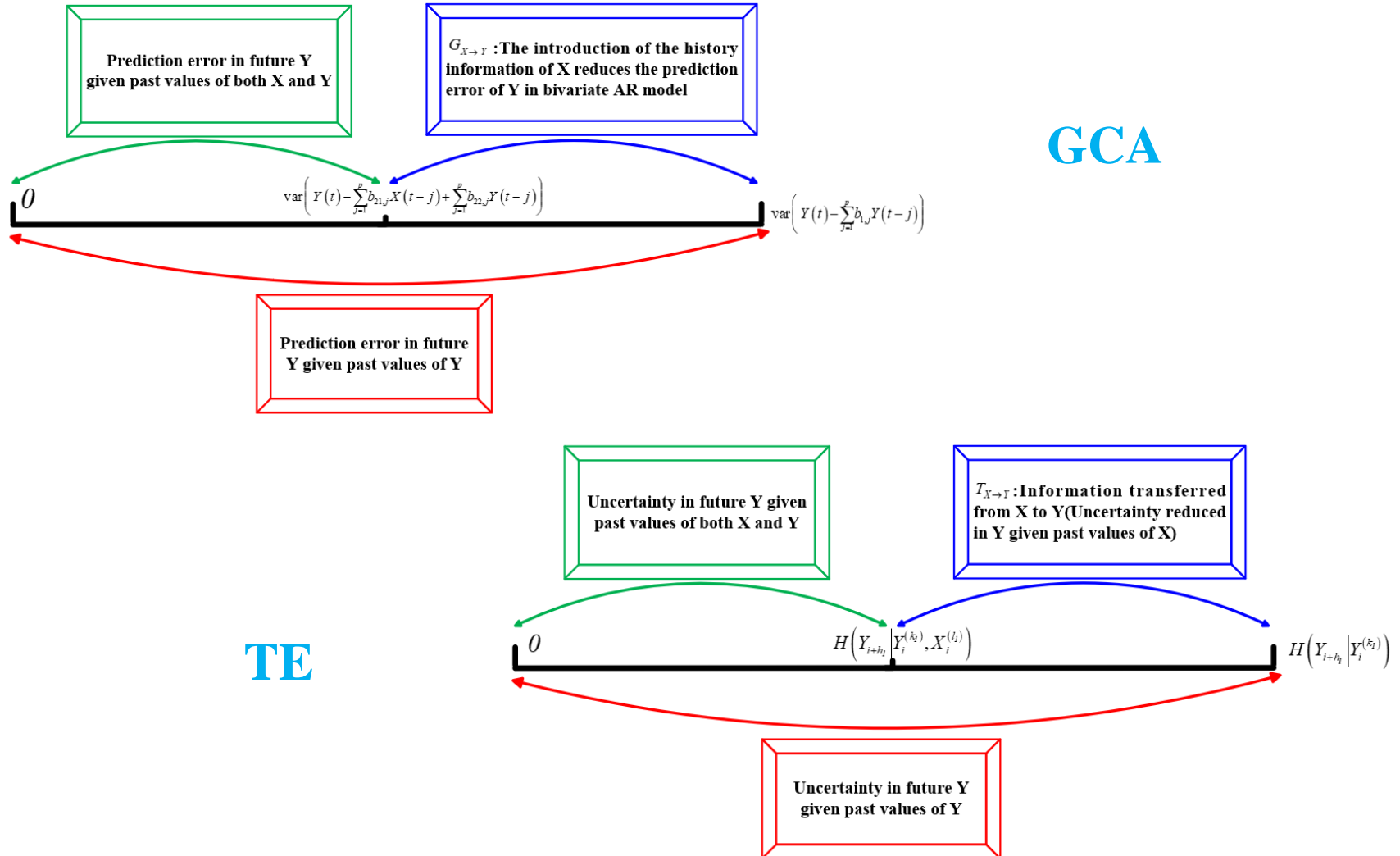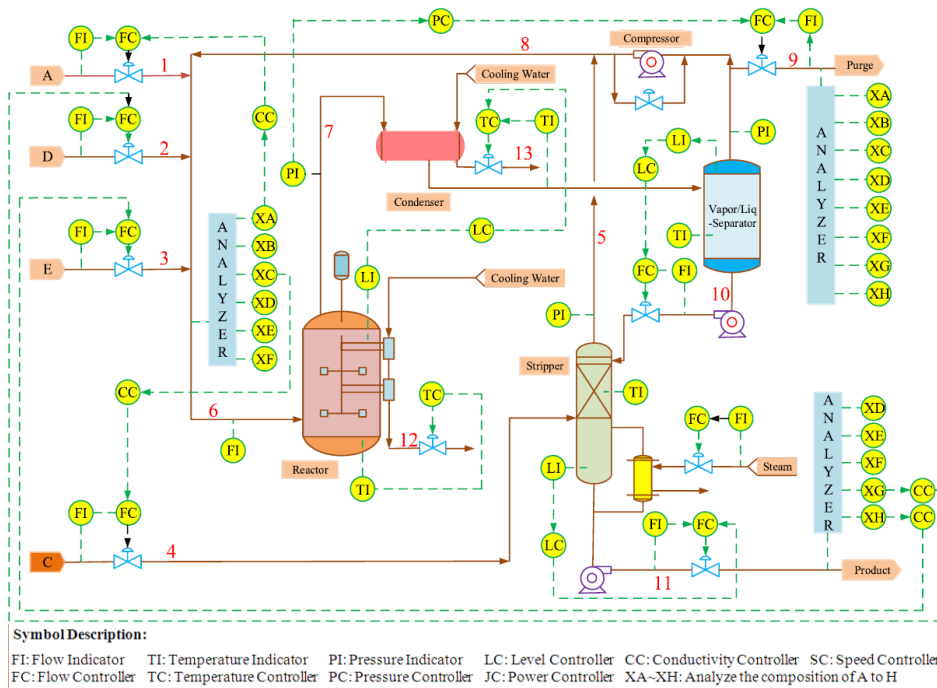Fig.1 Flowchart of two kinds of causality-based inferential sensors

# How


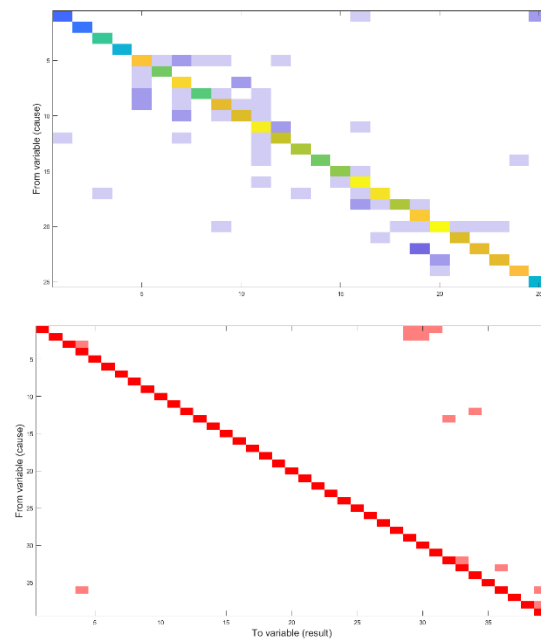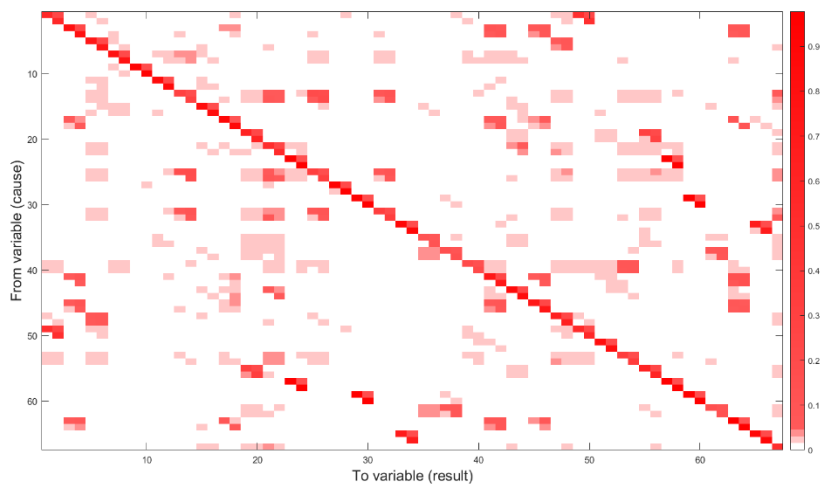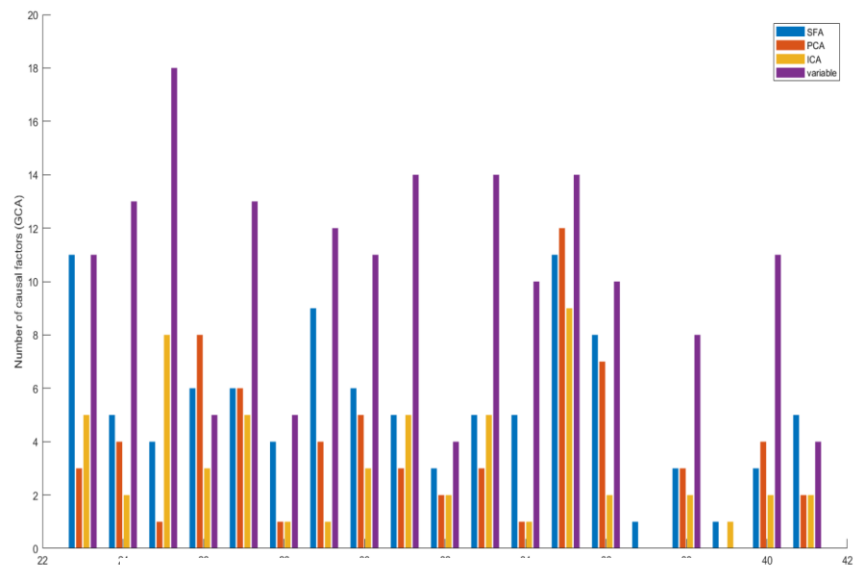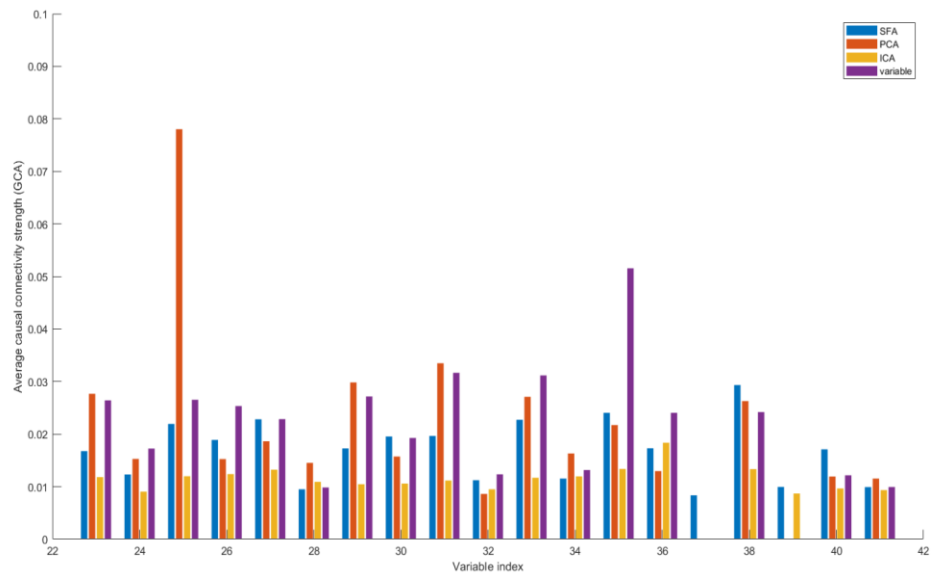
Fig.1 Flowchart of two kinds of causality-based inferential sensors

Causality analysis methods:

Granger Causality analysis
Transfer Entropy
…..

Latent feature extraction methods:

PCA ICA SFA CVA
…..

Regression methods:

LS, Bayesian, Decision tree, SVM, Neural networks
…..

Prediction error in future Y given past values of both X and Y

$G_{X \to Y}$ : The introduction of the history information of X reduces the prediction error of Y in bivariate AR model

GCA

$\mathrm{var}\left( Y(t) - \sum_{j=1}^{p} b_{21,j} X(t-j) + \sum_{j=1}^{p} b_{22,j} Y(t-j) \right)$

$\mathrm{var}\left( Y(t) - \sum_{j=1}^{p} b_{1,j} Y(t-j) \right)$

Prediction error in future Y given past values of Y

Uncertainty in future Y given past values of both X and Y

$T_{X \to Y}$ : Information transferred from X to Y(Uncertainty reduced in Y given past values of X)

TE

$H\left( Y_{i+h_l} \middle| Y_i^{(k_l)}, X_i^{(l_l)} \right)$

$H\left( Y_{i+h_l} \middle| Y_i^{(k_l)} \right)$

Uncertainty in future Y given past values of Y

**Symbol Description:**

FI: Flow Indicator    TI: Temperature Indicator    PI: Pressure Indicator    LC: Level Controller    CC: Conductivity Controller    SC: Speed Controller
FC: Flow Controller    TC: Temperature Controller    PC: Pressure Controller    JC: Power Controller    XA~XH: Analyze the composition of A to H

There are 52 different variables in this process, among which 33 variables can be measured in real time while another 19 variables need to be analyzed respectively. Hence, 33 variables are chosen as the process data and 19 variables are seen as the quality-relevant variables to be predicted (only use normal data, no fault data). We choose and 33 process variables as X and the variable 31 as Y.

# Results

# Results

# Results