# Subset Selection in Multiple Linear Regression
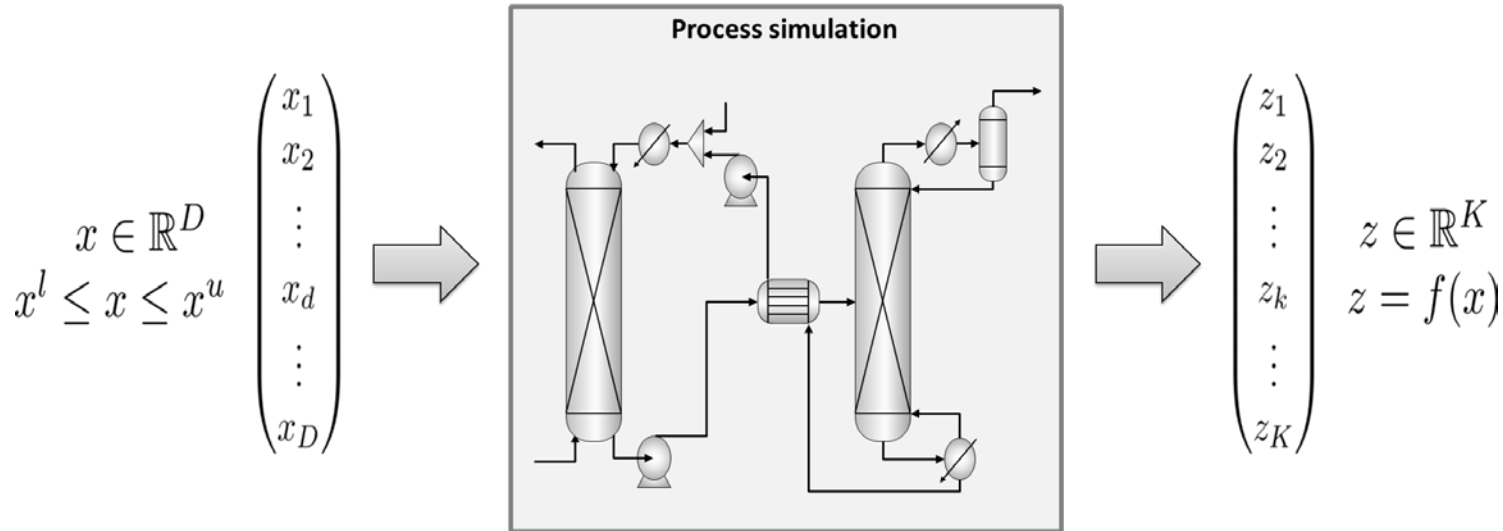
## Zachary Wilson

ztw@andrew.cmu.edu

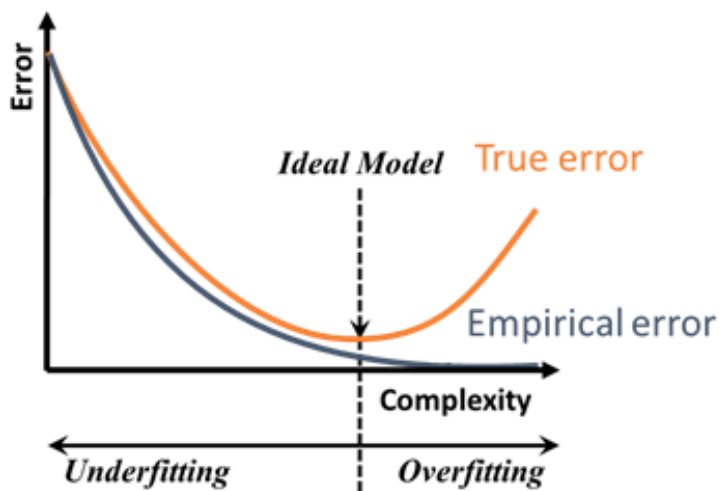## Nick Sahinidis

Sahinidis@cmu.edu

# Subset Selection in Multiple Linear Regression

**Step 1: Define a large set of potential basis functions**

$$\hat{z}(x) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2 + \beta_4 \frac{x_1}{x_2} + \beta_5 \frac{x_2}{x_1}$$

**Step 2: Model reduction**

$$\hat{z}(x) = \beta_0 + \beta_2 x_2 + \beta_5 \frac{x_2}{x_1}$$

Error vs Complexity:
- Ideal Model
- True error
- Empirical error
- Underfitting / Overfitting

**Subset Selection** is used to build surrogate models that are
- **Accurate** representations of higher order functions or black-box simulations
- **Simple** in functional form, tailored for algebraic optimization

**Fitness Criterion**
- **Balances** model complexity with reduction in empirical error
- **Penalize directly** for the number of explanatory variables in the regression model

# IP Formulations of Fitness Criterion

$$\min \quad \frac{1}{2} x^T Q x + c^T x$$

$$\text{s.t.} \quad -M z_j \leq \beta_j \leq M z_j \quad (j = 1, 2..., k)$$

$$z_j \in \{0, 1\}$$

**MIQP formulations**
- Solved **directly** (Cp, BIC)
- Solved in **nested optimization problem** (AIC, MSE)

**Alternative Model Selection Techniques**
- Regularization – LASSO, Ridge Regression
- Stepwise Heuristics

$$\min_{K \in \{1,...,K^u\}} \quad [\phi_{\beta,y}(\beta, y)|_K] + \phi_K(K)$$

$$\text{s.t.}$$

$$\min_{\beta,y} \quad [\phi_{\beta,y}(\beta, y)|_K]$$

$$\text{s.t.} \quad -M z_j \leq \beta_j \leq M z_j$$

$$\sum_{j \in J} z_j \leq K$$

$$z_j \in \{0, 1\}$$