



# Strengthened Regression Models through Response Variable Bounds

Alison Cozad<sup>1,2</sup>, Nick Sahinidis<sup>1,2</sup>, David Miller<sup>1</sup>

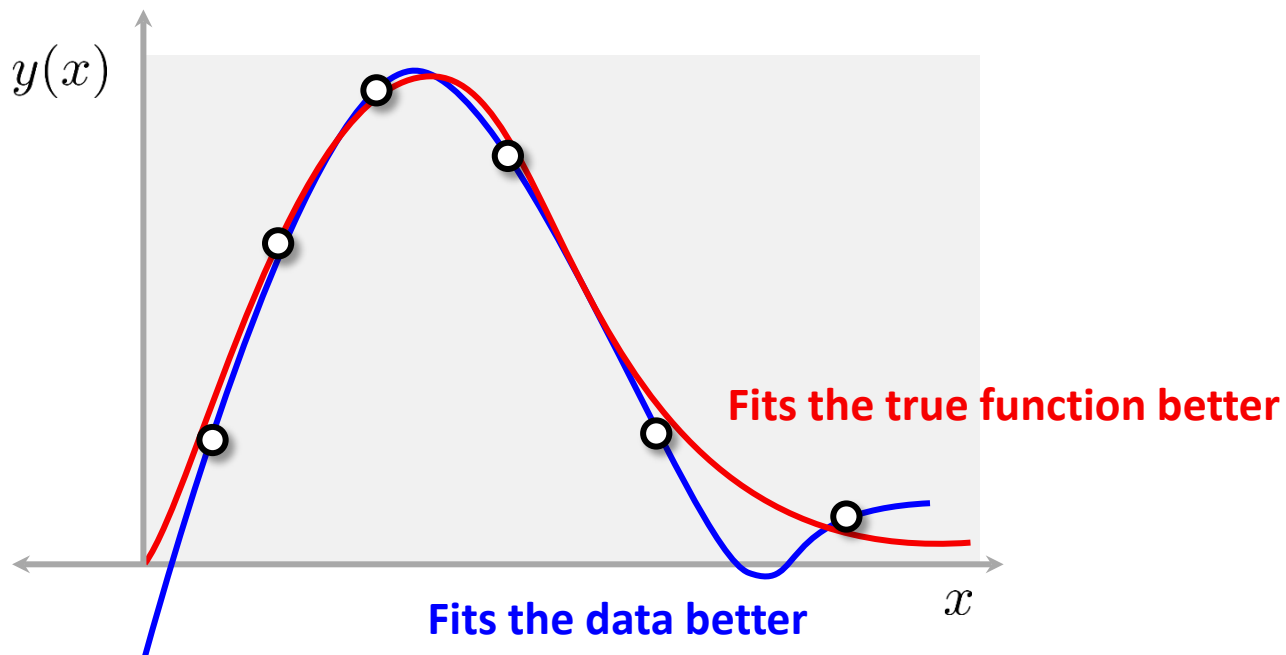
<sup>1</sup>National Energy Technology Laboratory, Pittsburgh, PA, USA

<sup>2</sup>Department of Chemical Engineering, Carnegie Mellon University, Pittsburgh, PA, USA

*This report was prepared as an account of work sponsored by an agency of the United States Government. Neither the United States Government nor any agency thereof, nor any of their employees, makes any warranty, express or implied, or assumes any legal liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States Government or any agency thereof. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States Government or any agency thereof.*

# MOTIVATION

- Leverage more information than just sampled data when building a surrogate model



# LEAST SQUARES REGRESSION

- Ordinary least squares regression
  - Chooses regression coefficients based on a set of data points
- Generate a model for response,

$$\hat{y}(x) = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 \exp(x) + \dots$$

- Ordinary least squares regression problem

$$\min_{\beta} \sum_{i=1}^N (y_i - \hat{y}(x_i))^2$$

OR

$$\min_{\beta} \sum_{i=1}^N (y_i - [\beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \beta_3 \exp(x_i) + \dots])^2$$

Optimization variables

# CONSTRAINED REGRESSION

- What if more information is known of a less exact nature?
  - **Leverage all information available to the modeler**
- **Explicit restrictions placed on regressors**
  - Often times logical bounds on regressors can be found by inspection and/or analysis
    - *Ex: Physical constants*

$$k = \beta_0 \exp\left(\frac{\beta_1}{T}\right), \quad \beta_1 \leq 0 \rightarrow \text{Activation energy, positive}$$

- **Relationships between parameters can be found by inspection or analysis**
  - *Ex: Intuitive relationships*

$$H = \beta_0 + \beta_1 T^{\text{in}} + \beta_2 T^{\text{out}} + \dots$$

$$\beta_2 - \beta_1 \geq 0 \rightarrow \text{Heat capacity, nonnegative}$$

# ADDING EXPLICIT CONSTRAINTS

- Adding in explicit constraints is rather straight forward,

$$\hat{H}(T) = \beta_0 + \beta_1 T^{\text{in}} + \beta_2 T^{\text{out}} + \dots$$

$$\beta_2 - \beta_1 \geq 0 \rightarrow \text{Heat capacity, nonnegative}$$

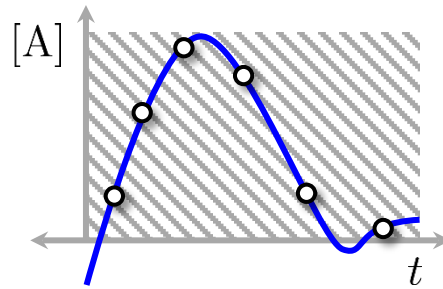
$$\min_{\beta} \sum_{i=1}^N (H_i - [\beta_0 + \beta_1 T_i^{\text{in}} + \beta_2 T_i^{\text{out}} + \dots])^2$$

$$\text{s.t. } \beta_2 - \beta_1 \geq 0$$

# CONSTRAINED REGRESSION

- What if more information is known of a less exact nature?
  - **Leverage all information available to the modeler**
- Restrictions implied by constraints on dependent variables
  - Implied by bounds on dependent variable

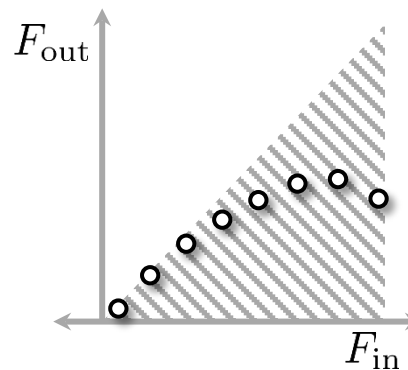
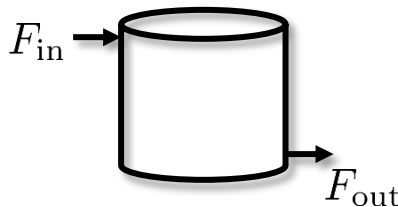
• *Ex:*



$$0 \leq [A]_t \leq [A]^{\max}$$

- Implied by constraints on dependent variable

• *Ex:*



$$F_{\text{out}}(F_{\text{in}}) \leq F_{\text{in}}$$

# ADDING IMPLIED CONSTRAINTS

- Adding in constraints implied by dependent variable bounds is less straight forward,

Generate a model for  $y$  given that  $y^l \leq y \leq y^u$

$$\hat{y}(x) = \beta_0 + \beta_1 x + \beta_2 x^2 + \dots$$

$$\min_{\beta} \sum_{i=1}^N (y_i - [\beta_0 + \beta_1 x + \beta_2 x^2 + \dots])^2$$

$$\text{s.t. } y^l \leq \beta_0 + \beta_1 x + \beta_2 x^2 + \dots \leq y^u \quad \forall x$$

# ADDING IMPLIED CONSTRAINTS

- Adding in constraints implied by dependent variable bounds is less straight forward,

Generate a model for  $y$  given that  $y^l \leq y \leq y^u$

$$\hat{y}(x) = \beta_0 + \beta_1 x + \beta_2 x^2 + \dots$$

$$\min_{\beta} \sum_{i=1}^N (y_i - [\beta_0 + \beta_1 x + \beta_2 x^2 + \dots])^2$$

$$\text{s.t. } y^l \leq \beta_0 + \beta_1 x + \beta_2 x^2 + \dots \leq y^u$$

$\forall x$

Semi-infinite  
programming



# ADDING IMPLIED CONSTRAINTS

- Adding in constraints implied by dependent variable bounds is less straight forward,

Generate a model for  $y$  given that  $y^l \leq y \leq y^u$

$$\hat{y}(x) = \beta_0 + \beta_1 x + \beta_2 x^2 + \dots$$

$$\min_{\beta} \sum_{i=1}^N (y_i - [\beta_0 + \beta_1 x + \beta_2 x^2 + \dots])^2$$

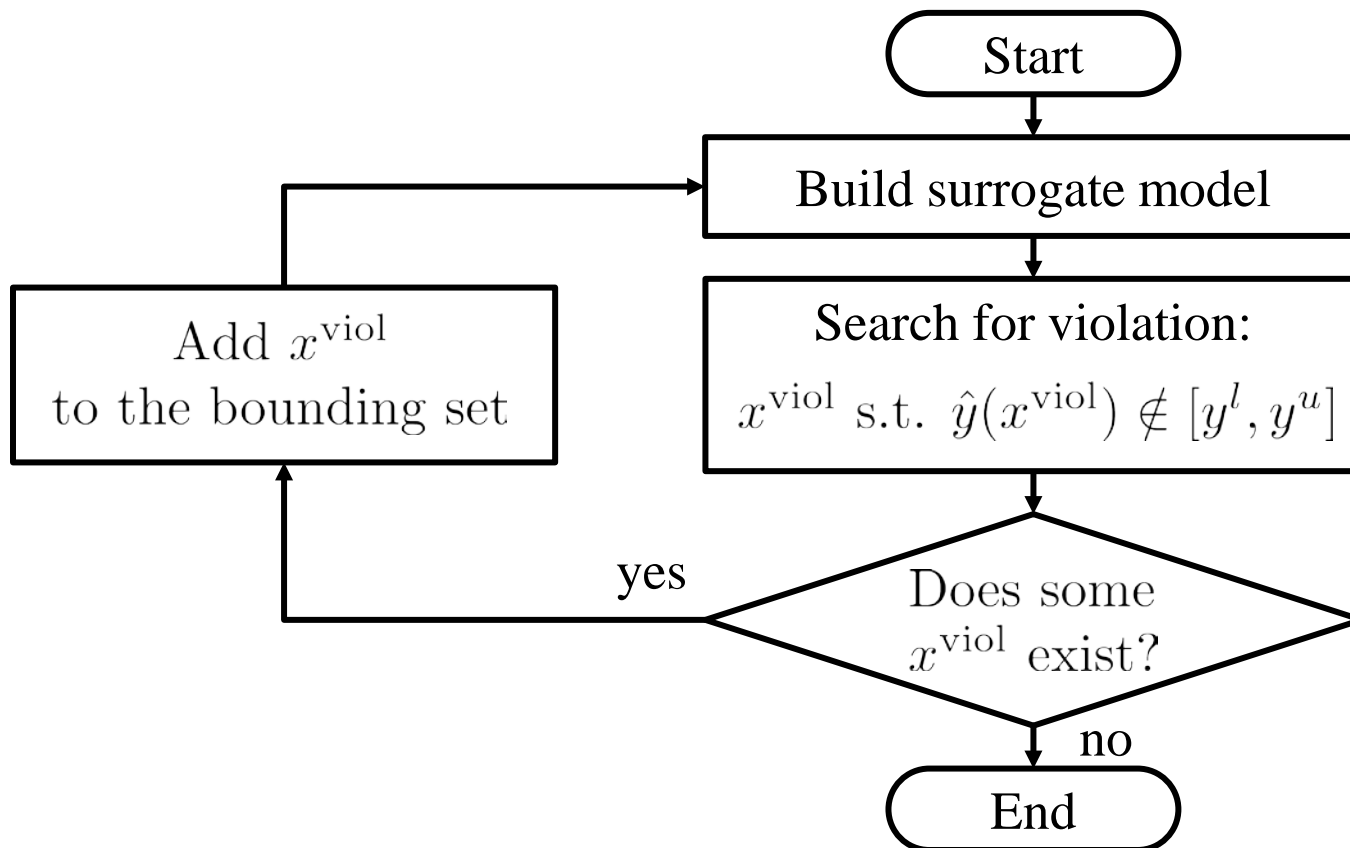
$$\text{s.t. } y^l \leq \beta_0 + \beta_1 x_j + \beta_2 x_j^2 + \dots \leq y^u$$

$\forall j \in \text{Bounding set}$

# IMPLEMENTATION

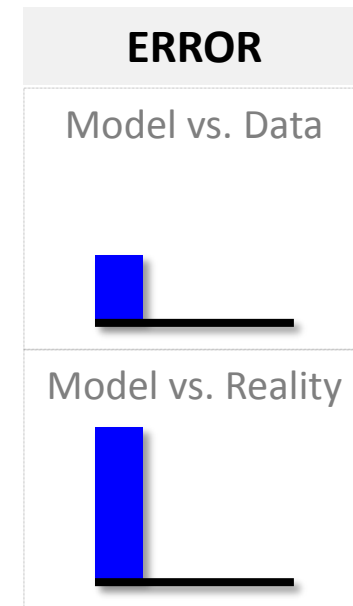
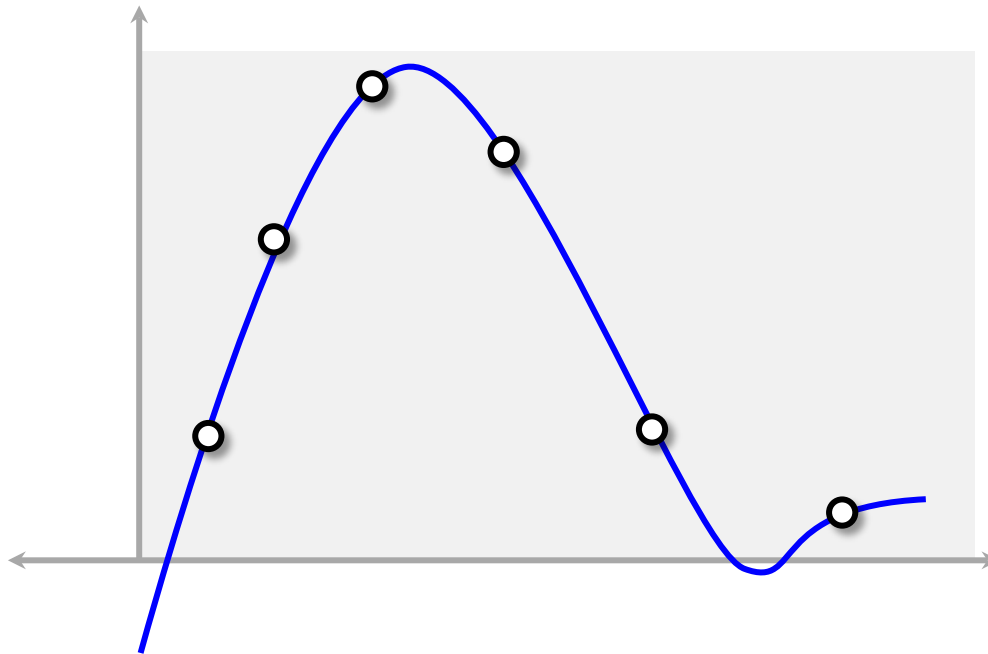
Generate a model for  $y$  given that  $y^l \leq y \leq y^u$

$$\hat{y}(x) = \beta_0 + \beta_1 x + \beta_2 x^2 + \dots$$



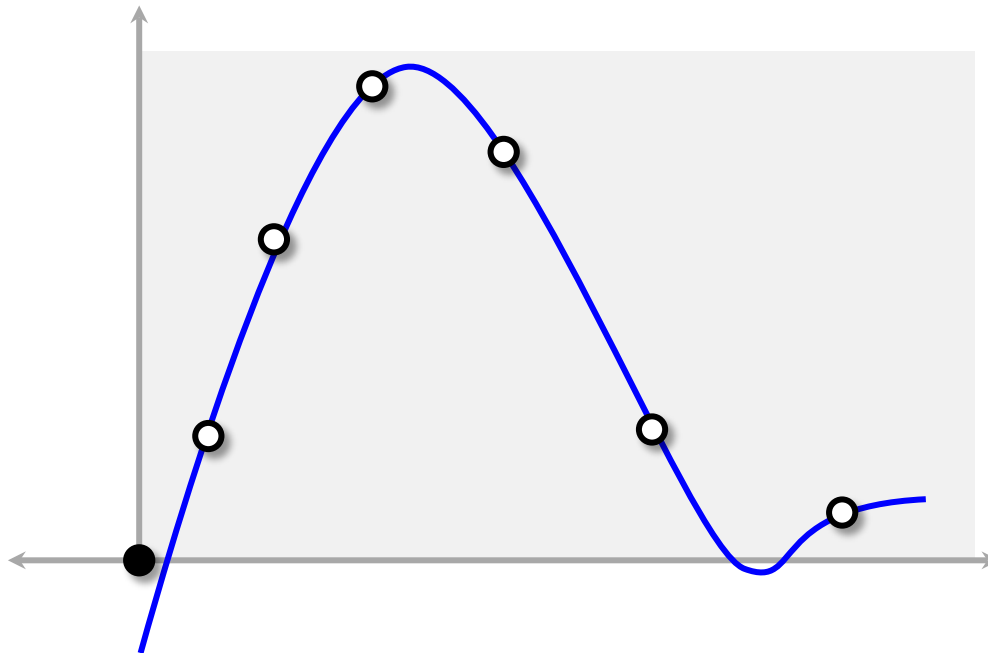
# ILLUSTRATIVE EXAMPLE

- Build an initial model

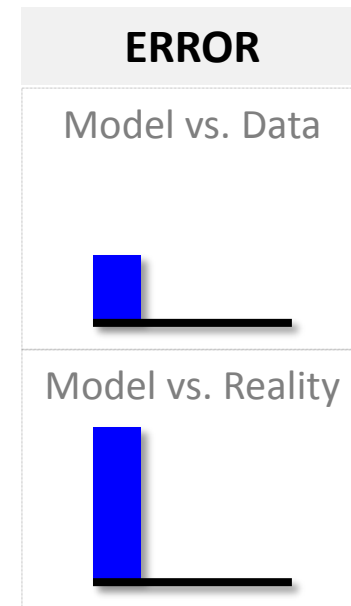


# ILLUSTRATIVE EXAMPLE

- Build an initial model
- Locate areas that violate output bounds

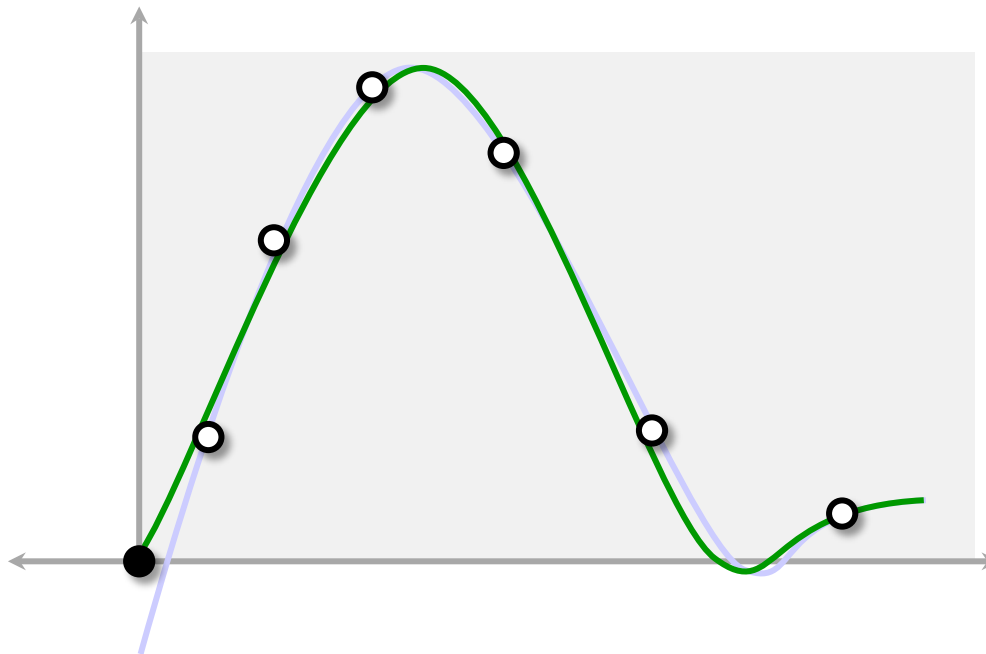


- Data points
- Bounding points

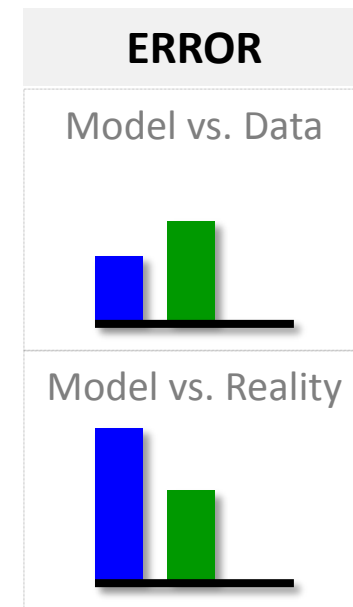


# ILLUSTRATIVE EXAMPLE

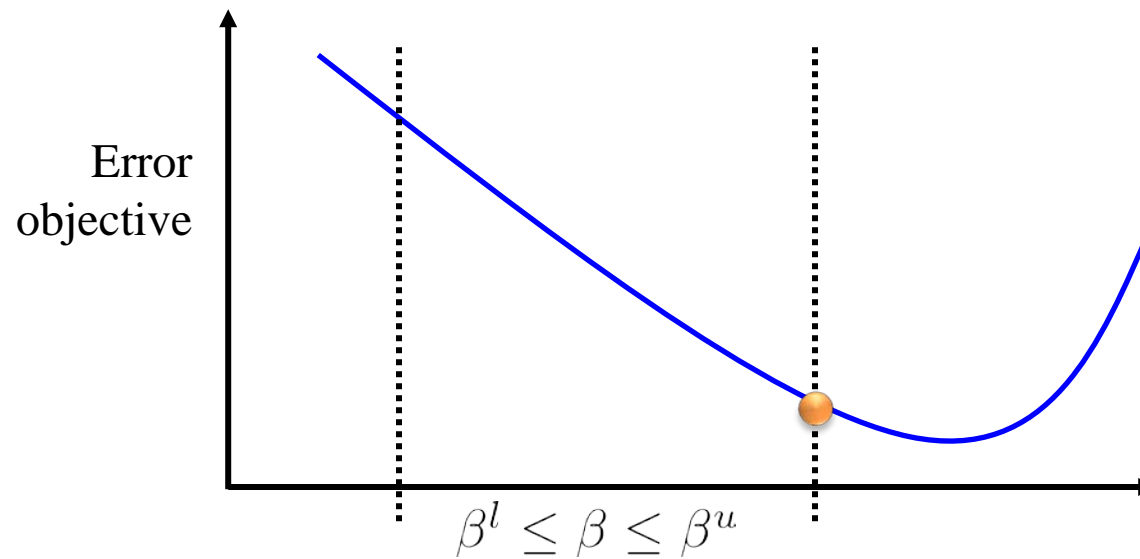
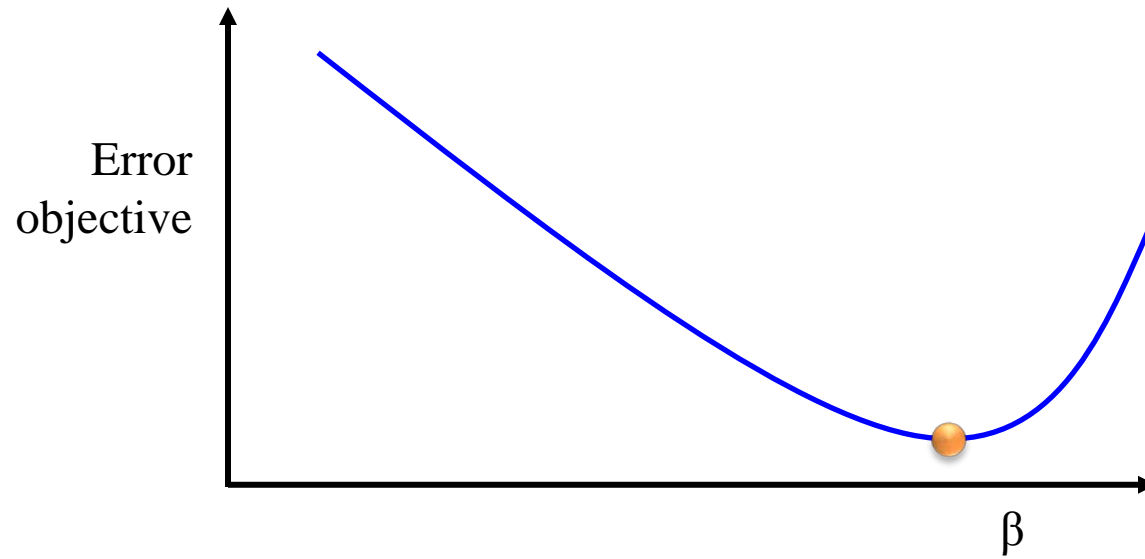
- Rebuild the model ensuring the output is within bounds at the bounding points



- Data points
- Bounding points

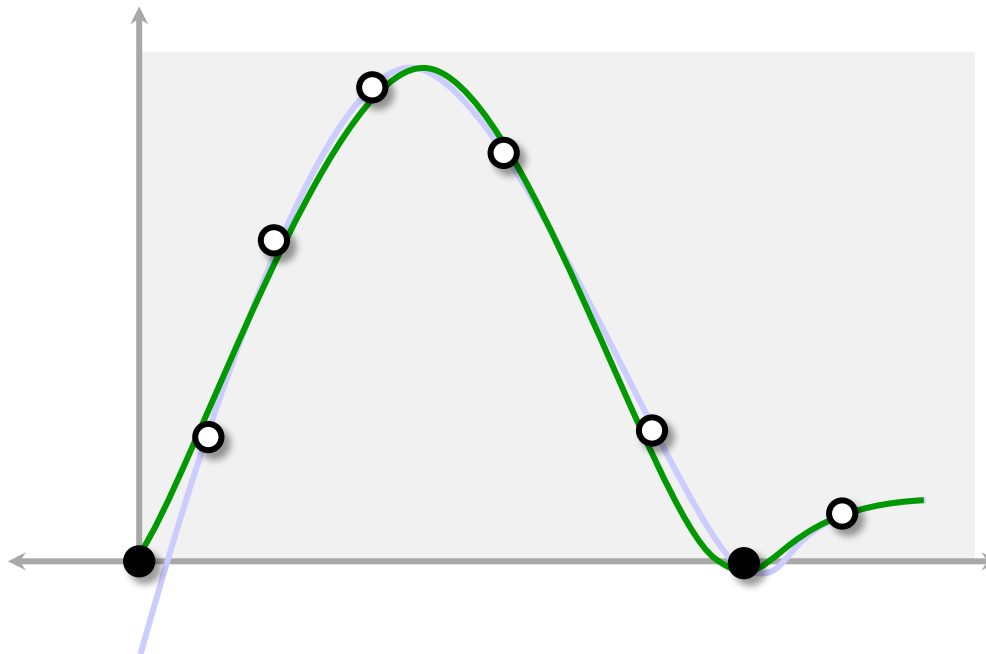


# NON STATIONARITY

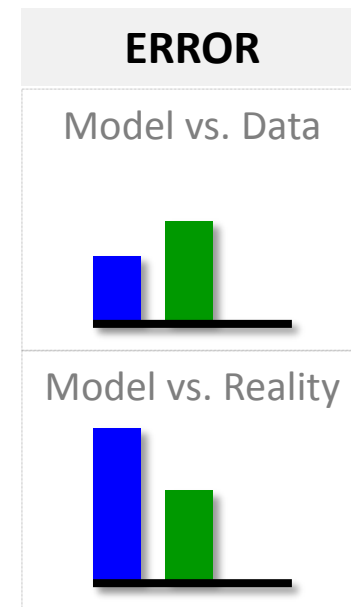


# ILLUSTRATIVE EXAMPLE

- Rebuild the model ensuring the output is within bounds at the bounding points
- Search for additional violation points

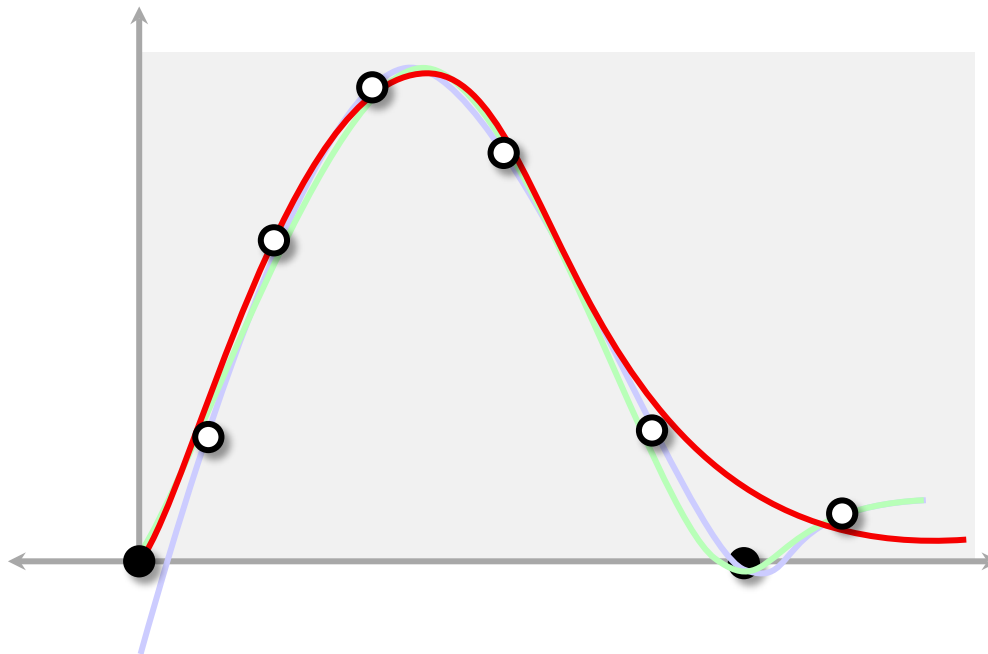


- Data points
- Bounding points

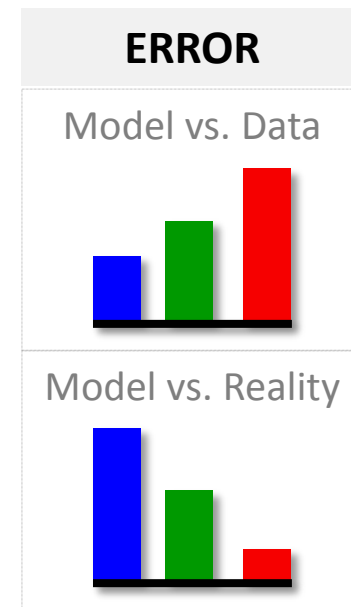


# ILLUSTRATIVE EXAMPLE

- Rebuild the model ensuring the output is within bounds at the bounding points
- Ensure that no violation points remain



- Data points
- Bounding points

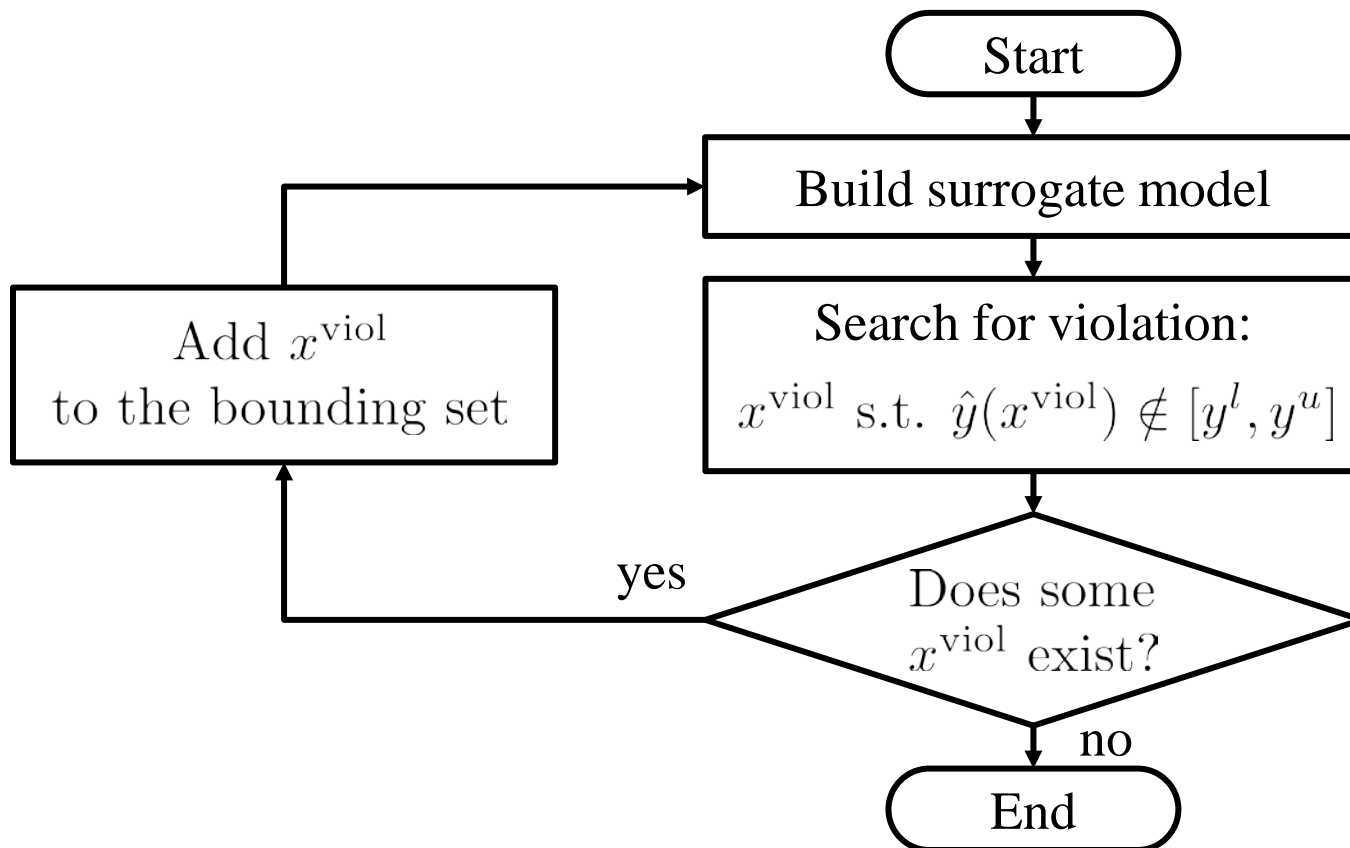




# IMPLEMENTATION

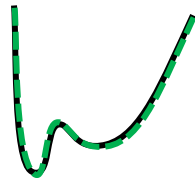
Generate a model for  $y$  given that  $y^l \leq y \leq y^u$

$$\hat{y}(x) = \beta_0 + \beta_1 x + \beta_2 x^2 + \dots$$



# TESTING PLATFORM

- We will test this implementation on an existing software: **ALAMO**
- **ALAMO (Automated Learning of Algebraic Models for Optimization)**
  - Iteratively sample and model black box systems as algebraic model that
    - *Accurate*
      - We want to reflect the true nature of the simulation
    - *Simple*
      - Low-complexity models



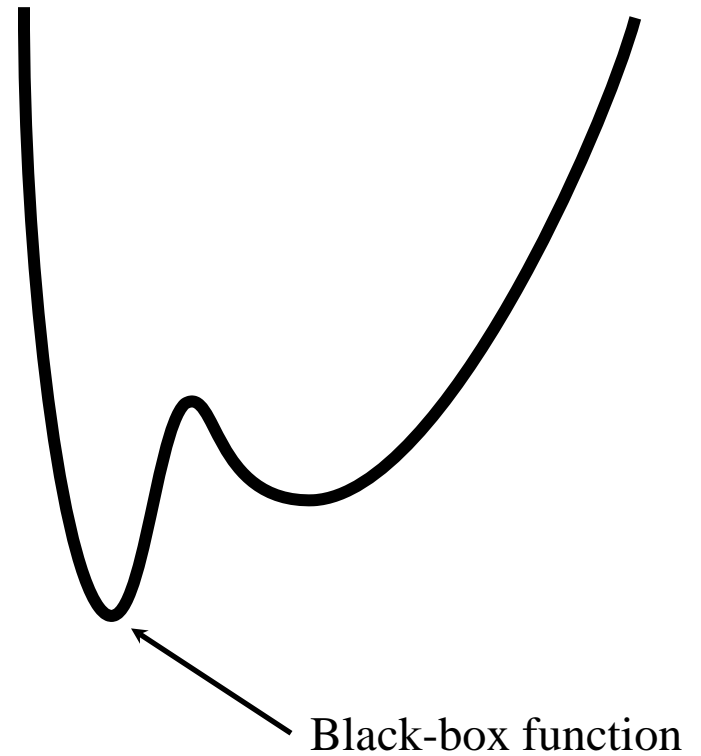
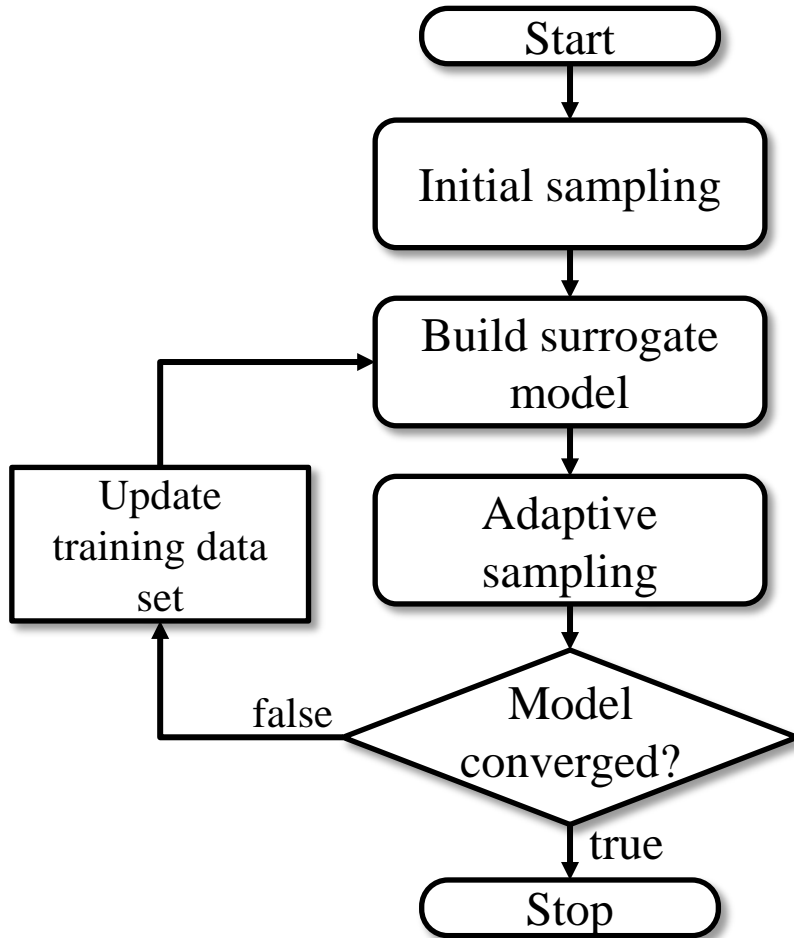
$$\hat{f}(x) = \sum_{i=1}^n \gamma_i \exp\left(\frac{\|x\|}{\sigma^2}\right) + \beta_0 + \beta_1 x + \dots$$

$$\hat{f}(x) = \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + \beta_4 e^x$$

- *Generated from a minimal data set*
  - Reduce experimental and simulation requirements

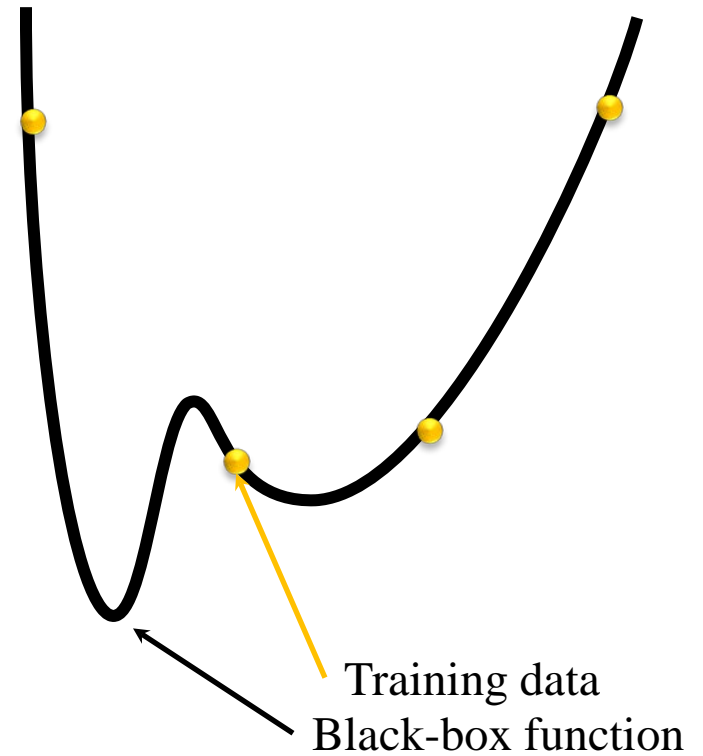
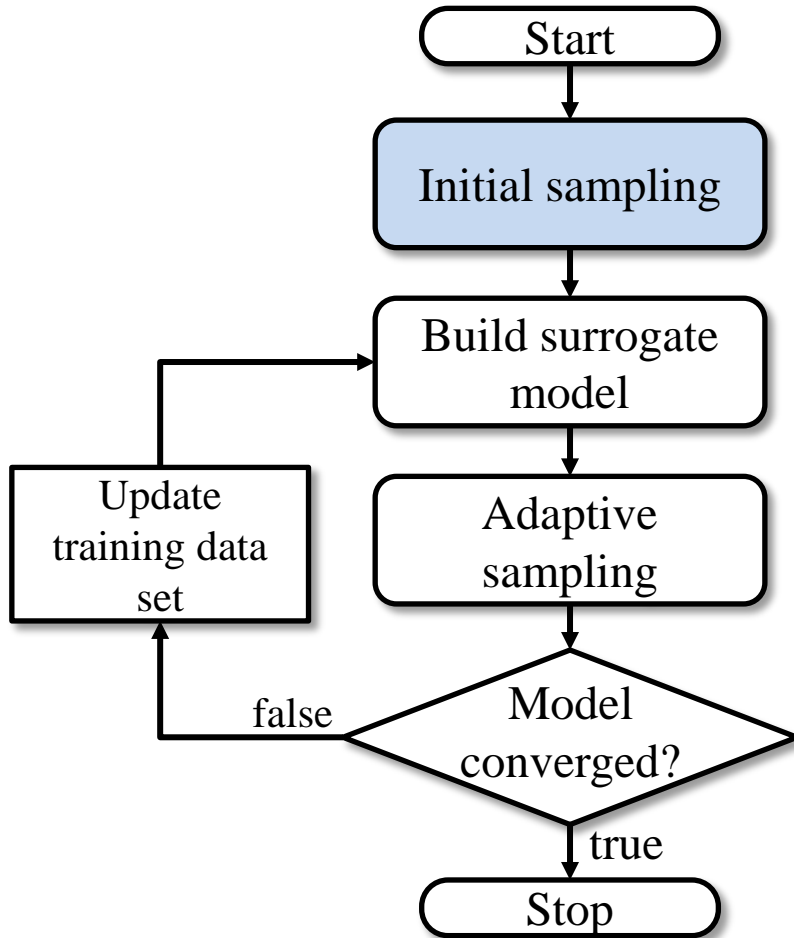
# ALAMO

## Automated Learning of Algebraic Models for Optimization



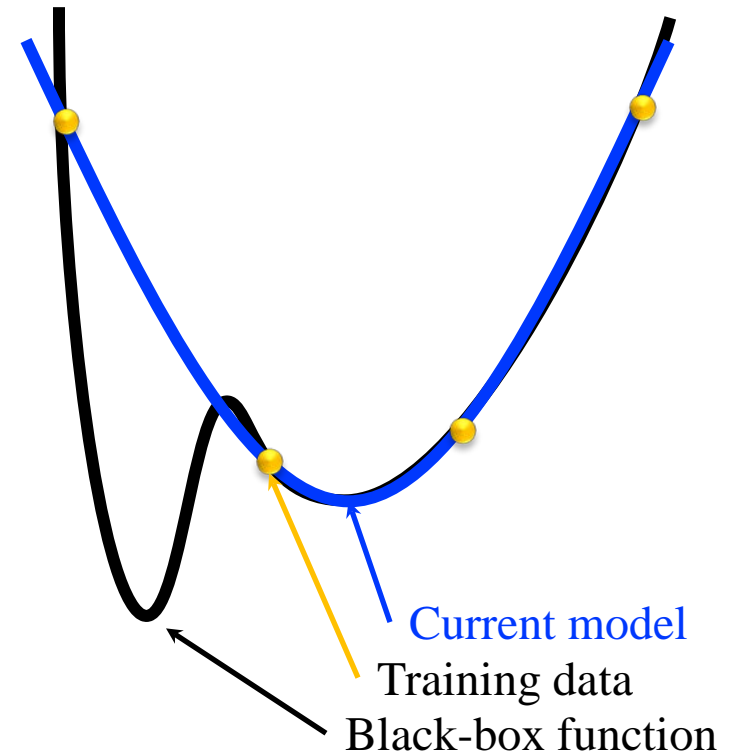
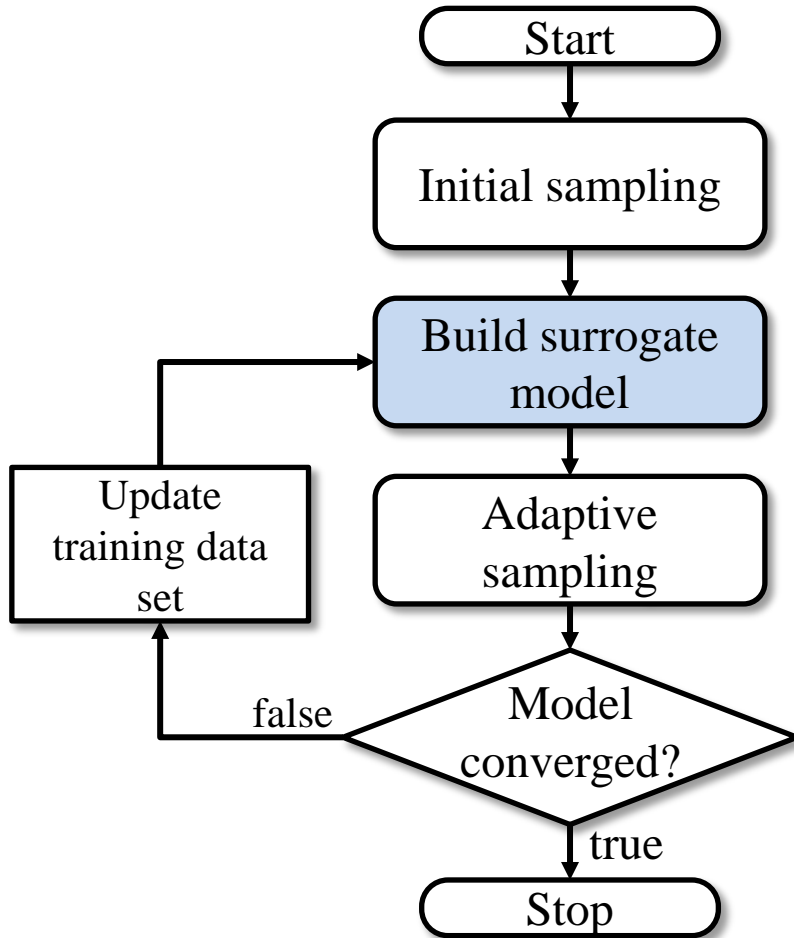
# ALAMO

## Automated Learning of Algebraic Models for Optimization



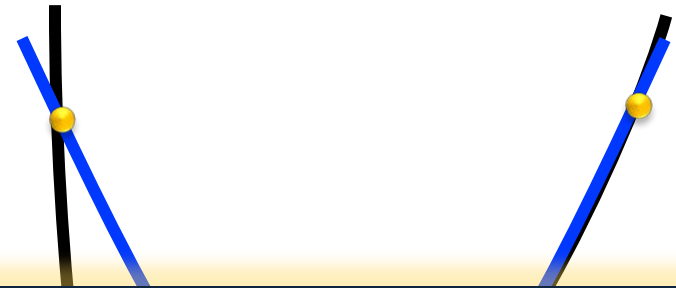
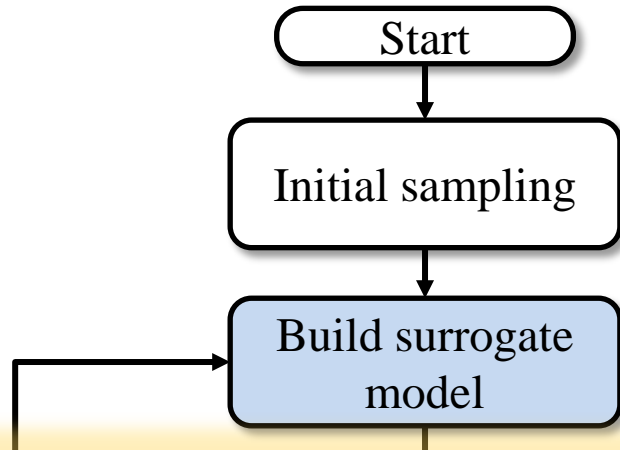
# ALAMO

## Automated Learning of Algebraic Models for Optimization



# ALAMO

## Automated Learning of Algebraic Models for Optimization



### Determining the unknown functional form

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

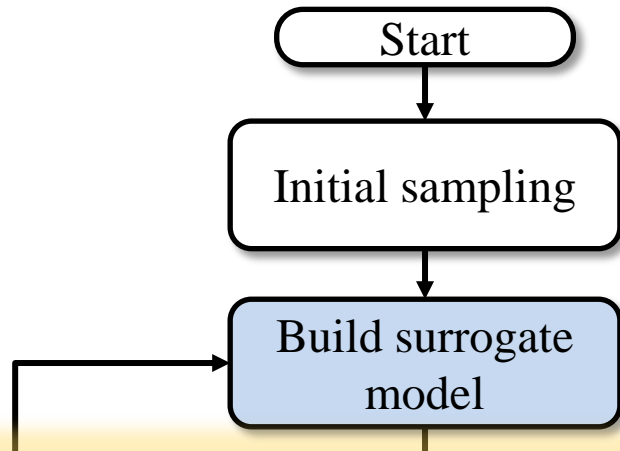
$$\hat{z}(x_1) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2 + \beta_4 \frac{x_1}{x_2} + \beta_5 \frac{x_2}{x_1} + \beta_6 e^{x_1} + \beta_7 e^{x_2} + \dots$$

**Step 2: Model reduction**

$$\hat{z}(x) = \beta_0 + \beta_2 x_2 + \beta_5 \frac{x_2}{x_1} + \beta_7 e^{x_2}$$

# ALAMO

## Automated Learning of Algebraic Models for Optimization



Using the new method, we **guarantee** that the model does not violate output variable bounds

### Determining the unknown functional form

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

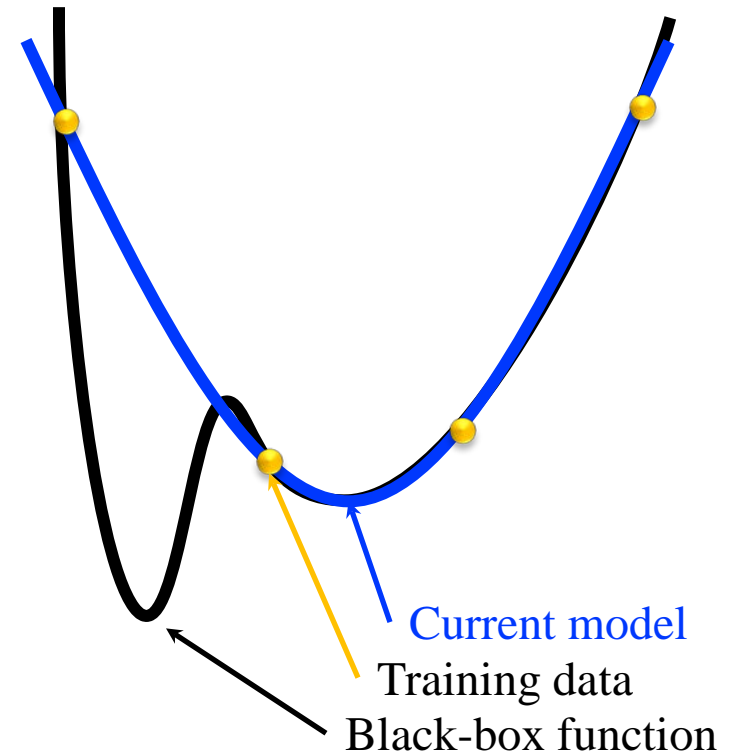
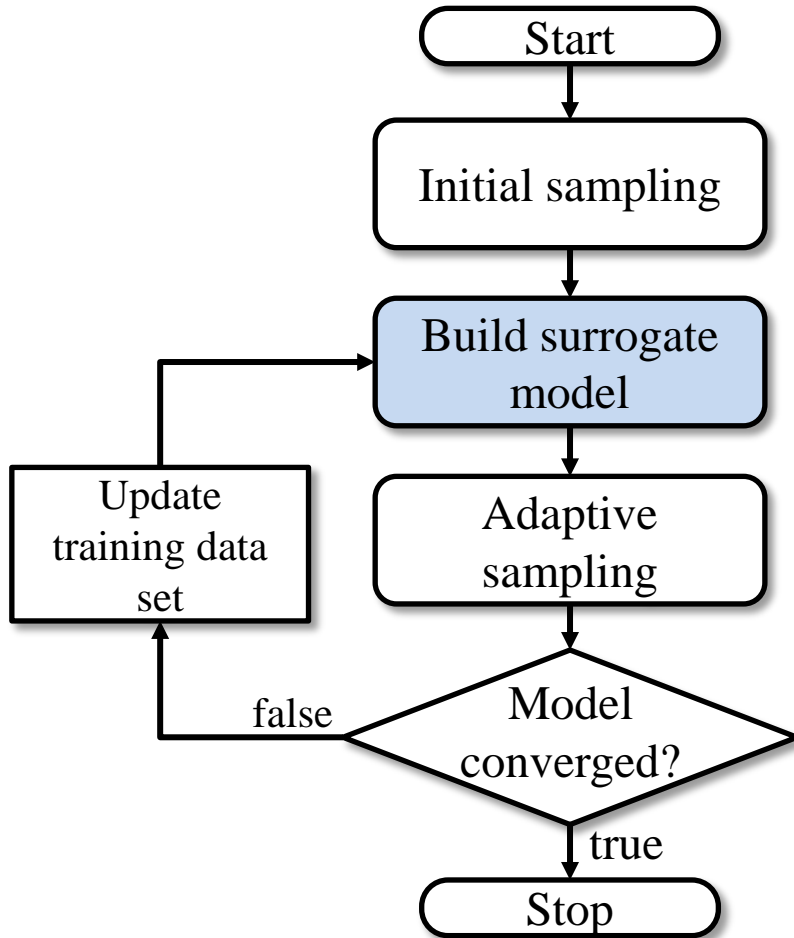
$$\hat{z}(x) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2 + \beta_4 \frac{x_1}{x_2} + \beta_5 \frac{x_2}{x_1} + \beta_6 e^{x_1} + \beta_7 e^{x_2} + \dots$$

**Step 2: Model reduction**

$$\hat{z}(x) = \beta_0 + \beta_2 x_2 + \beta_5 \frac{x_2}{x_1} + \beta_7 e^{x_2}$$

# ALAMO

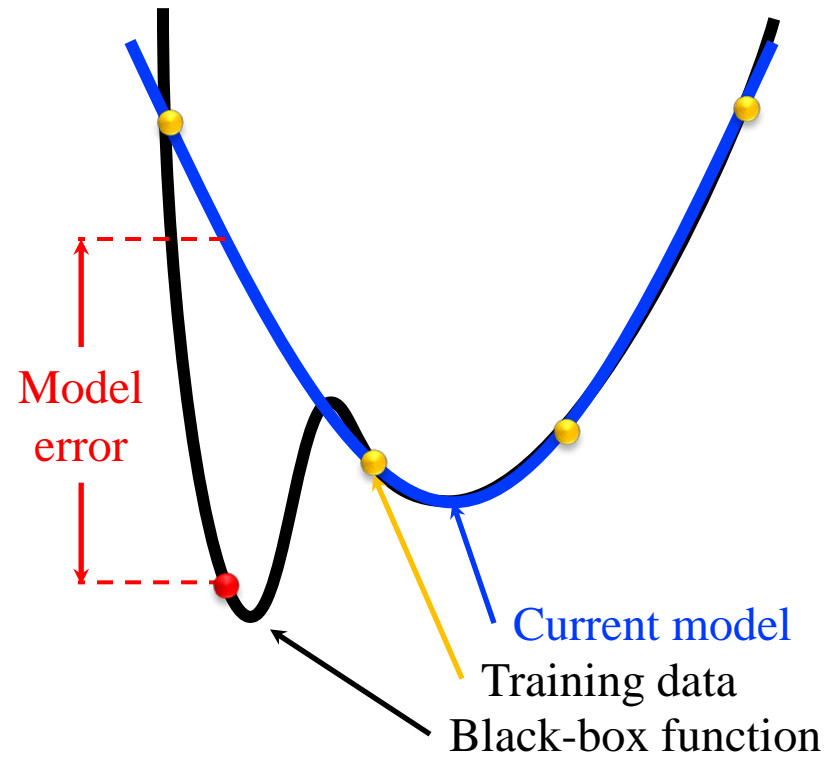
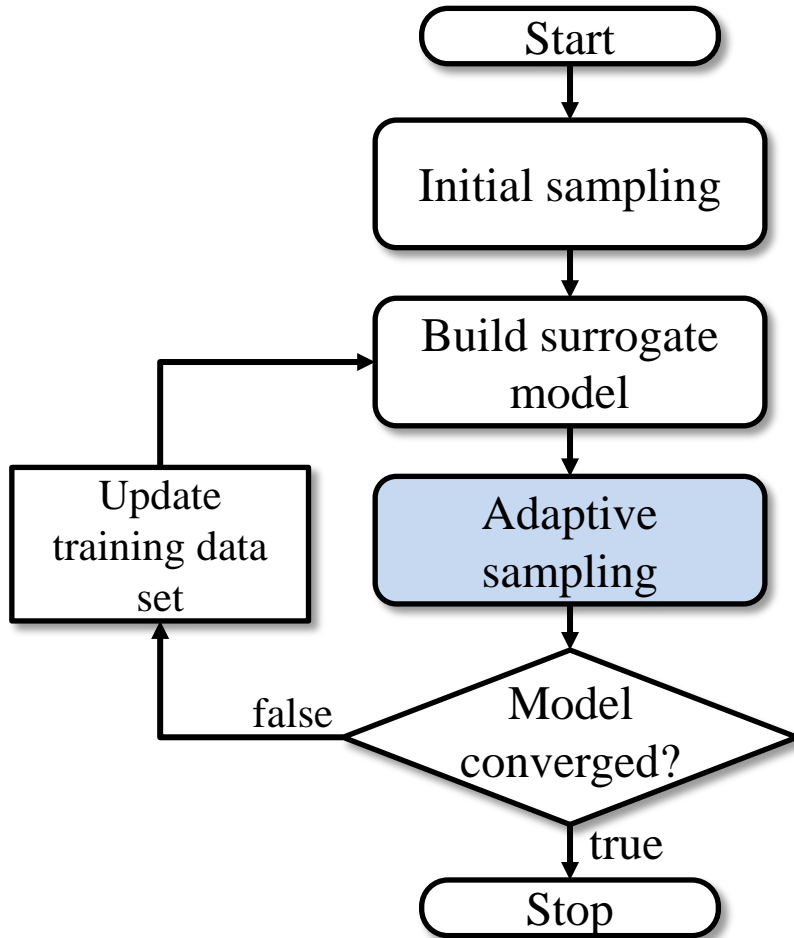
## Automated Learning of Algebraic Models for Optimization





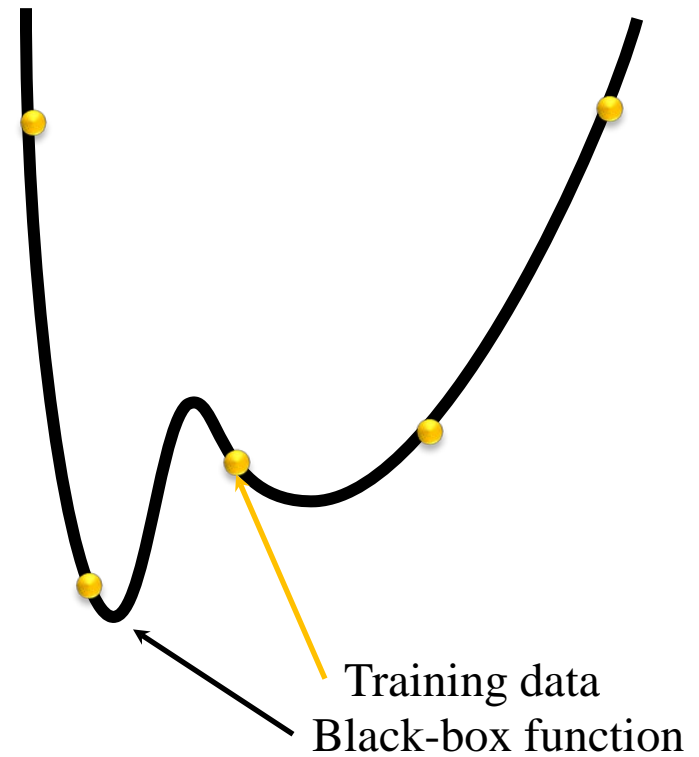
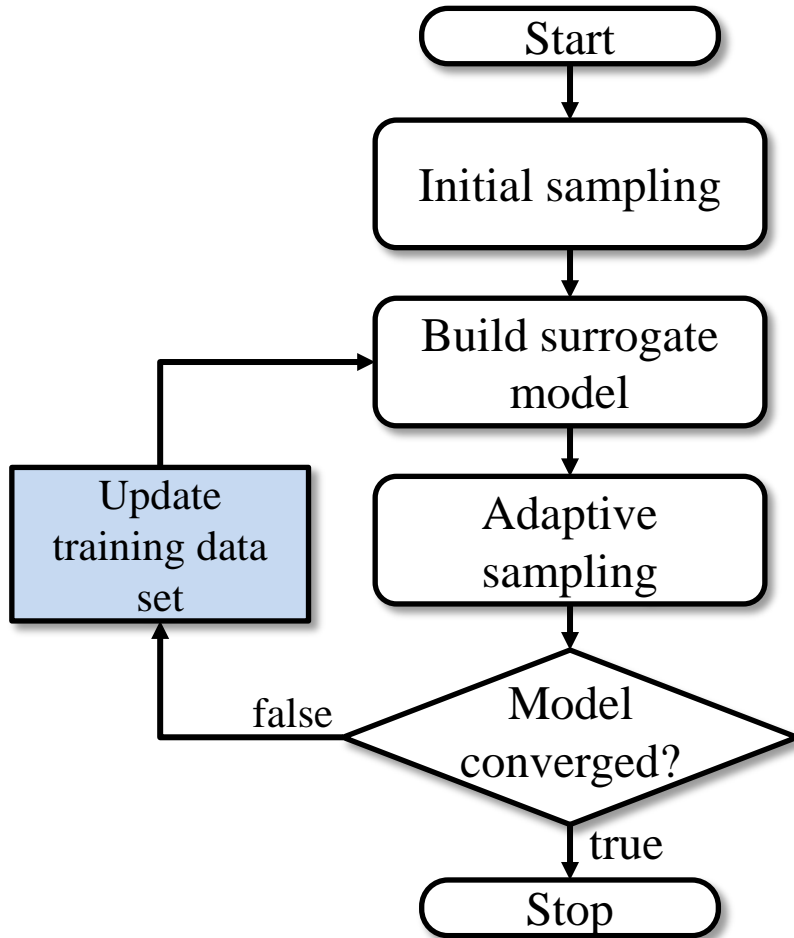
# ALAMO

## Automated Learning of Algebraic Models for Optimization



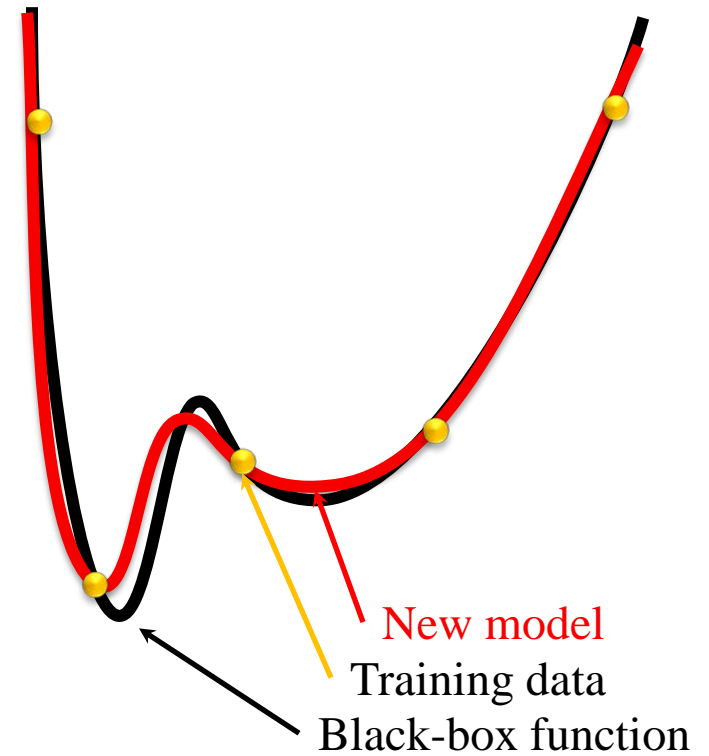
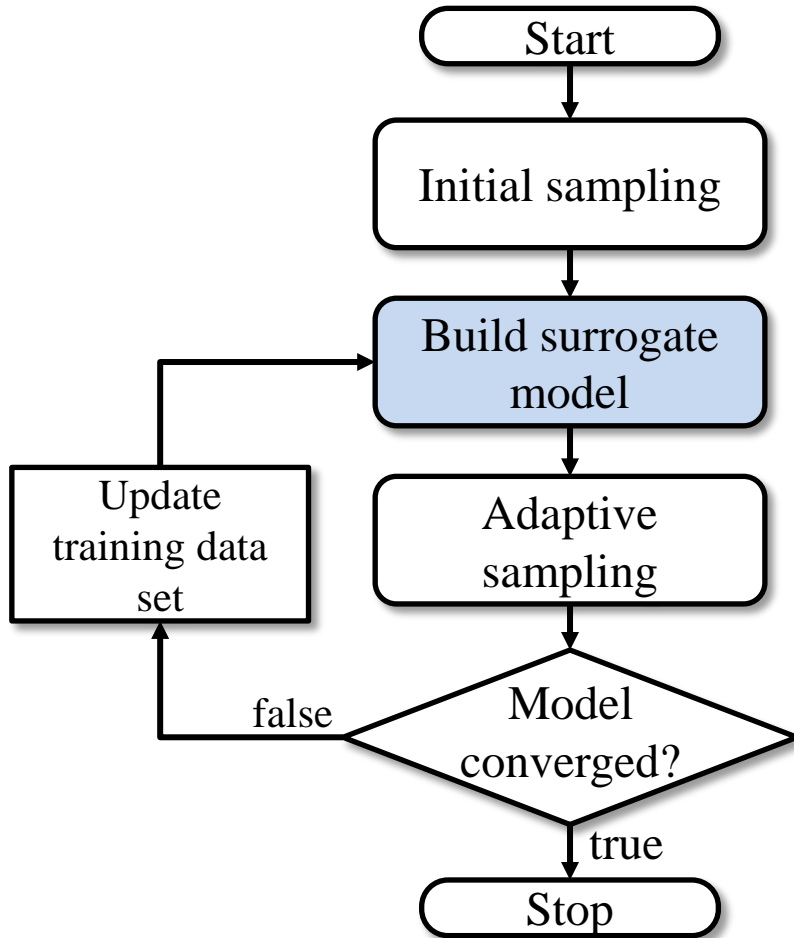
# ALAMO

## Automated Learning of Algebraic Models for Optimization

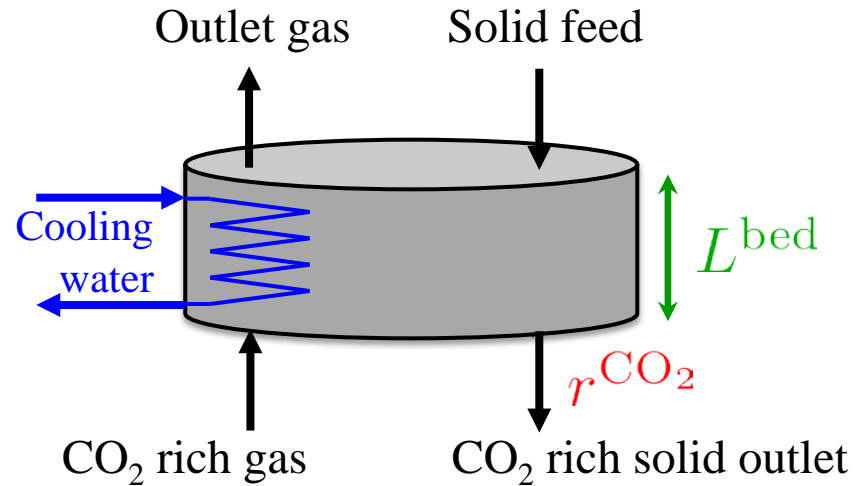


# ALAMO

## Automated Learning of Algebraic Models for Optimization



# SMALL EXMAPLE – CARBON CAPTURE



- **Dependent variable and range**
  - Fraction of CO<sub>2</sub> remove from the gas stream

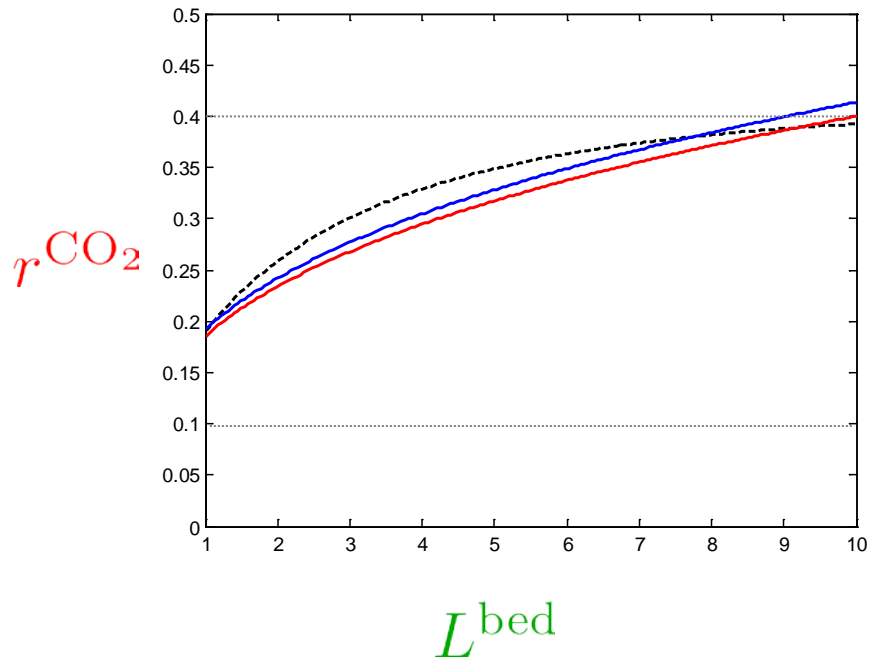
$$r^{\text{CO}_2}(L^{\text{bed}}, F^{\text{cw}}) = f_1(L^{\text{bed}}, F^{\text{cw}}) \quad 0.10 \leq r^{\text{CO}_2} \leq 0.4$$

- **Independent variable and range**

$$1 \text{ m} \leq L^{\text{bed}} \leq 10 \text{ m}$$

# COMPARISON – ITER 1

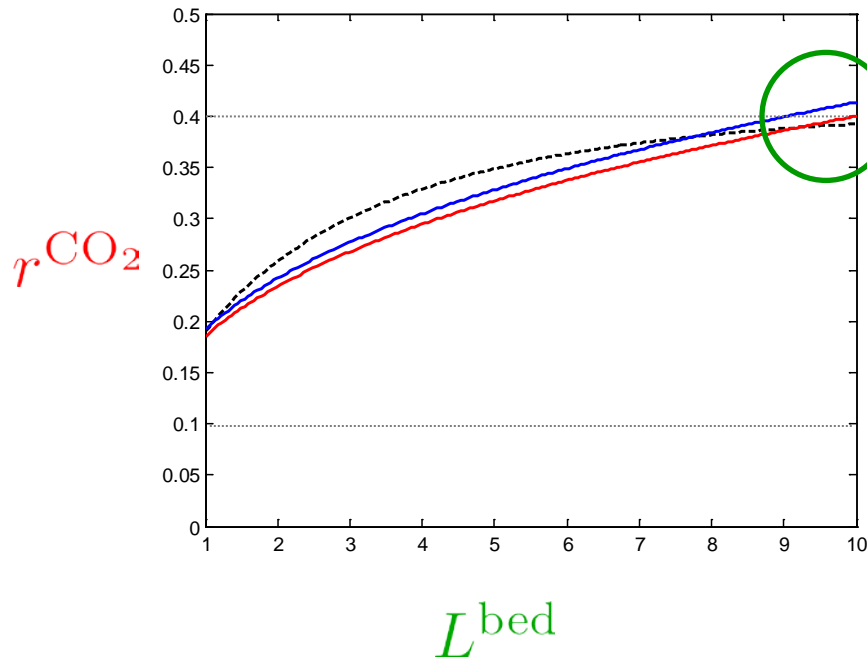
Model comparison



- ALAMO
- Constrained ALAMO

# COMPARISON – ITER 1

Model comparison

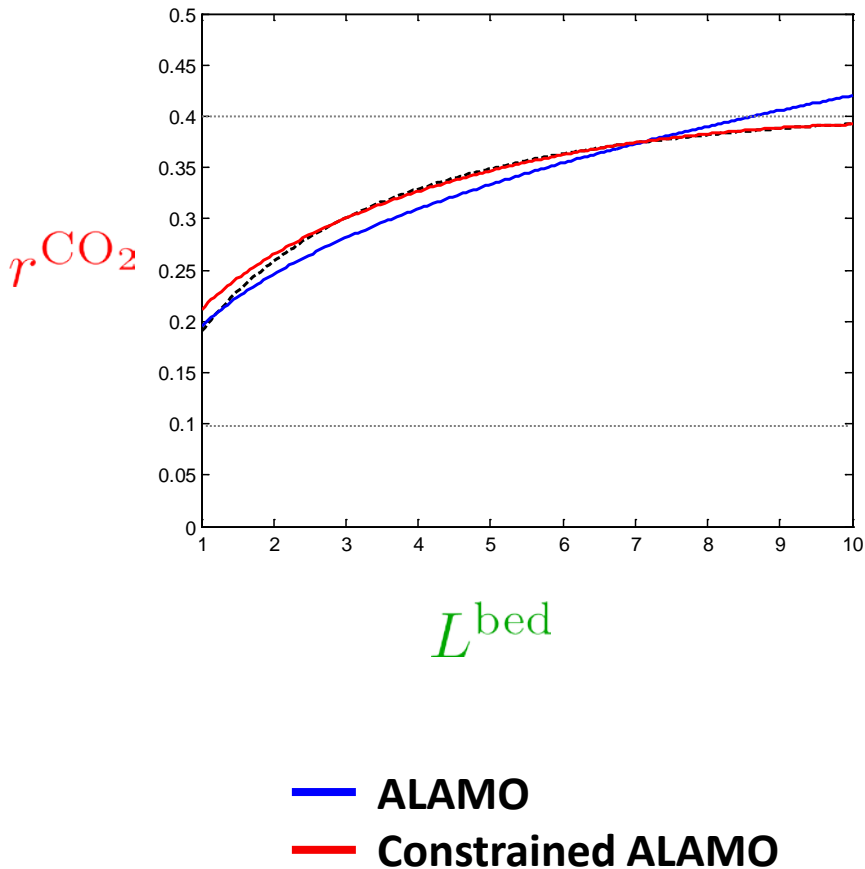


The constrained ALAMO model is able to stay within the bounds

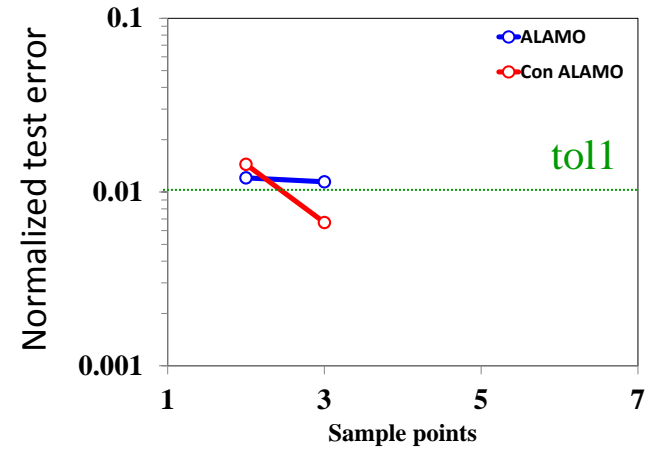
- ALAMO
- Constrained ALAMO

# COMPARISON – ITER 2

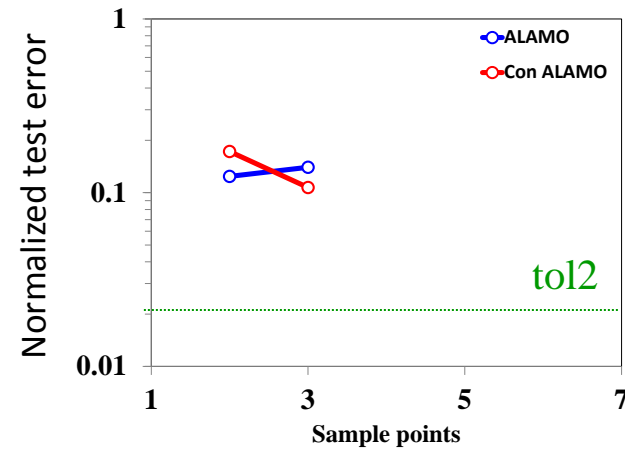
## Model comparison



## Mean

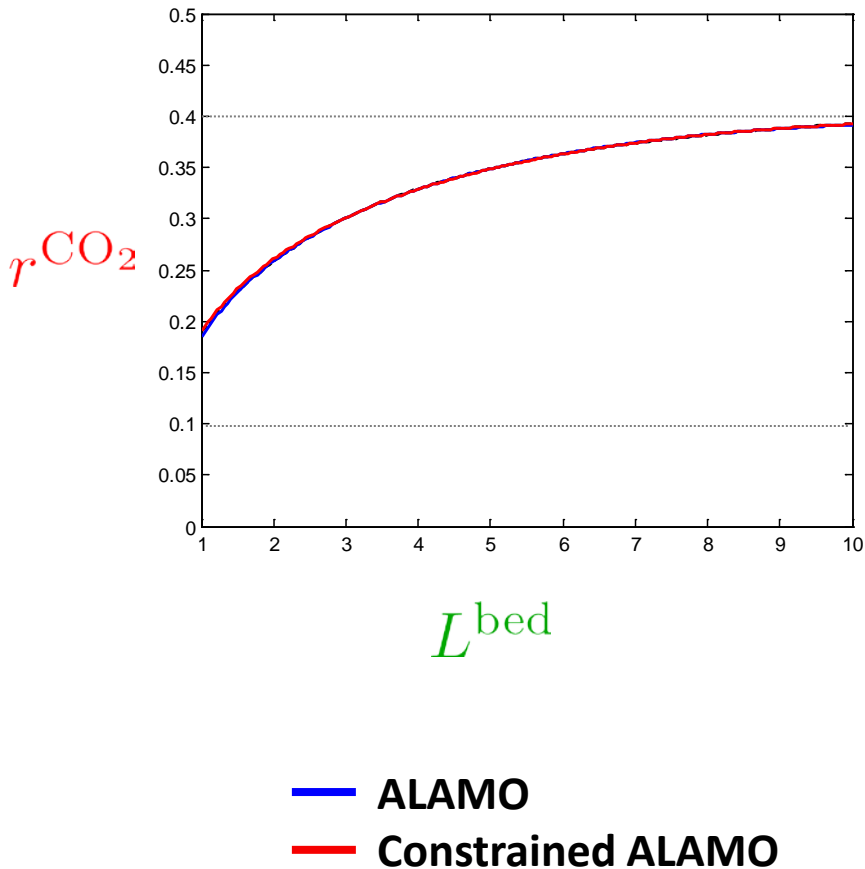


## Maximum

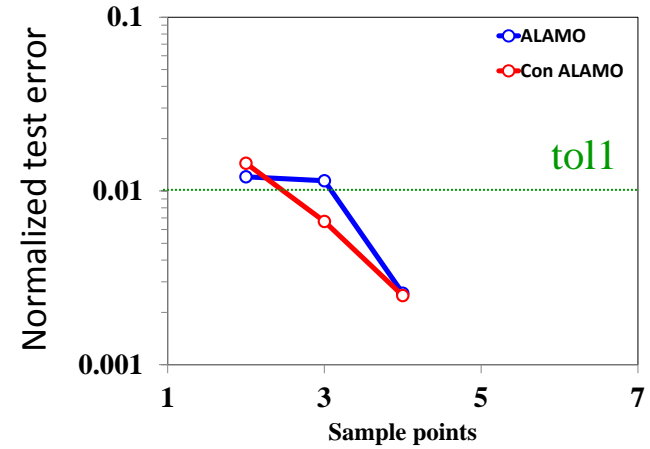


# COMPARISON – ITER 3

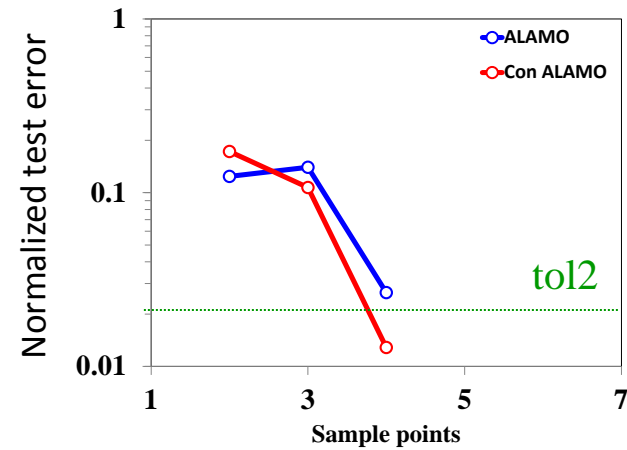
## Model comparison



## Mean



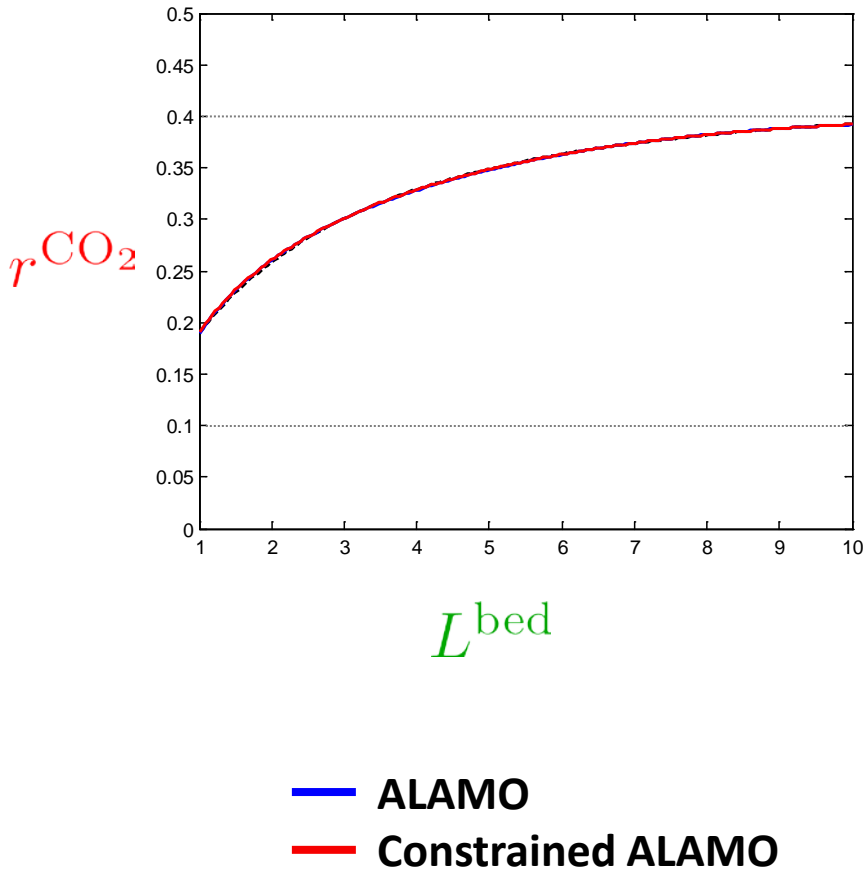
## Maximum



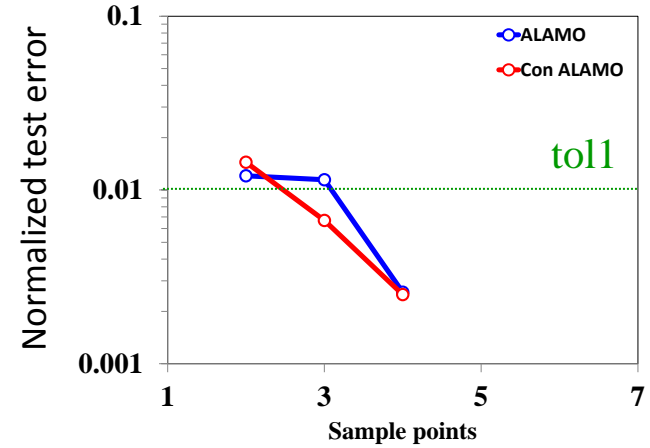


# COMPARISON – ITER 4

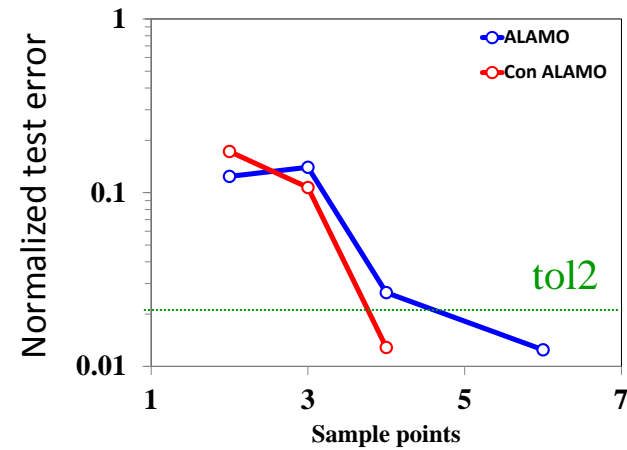
## Model comparison



## Mean

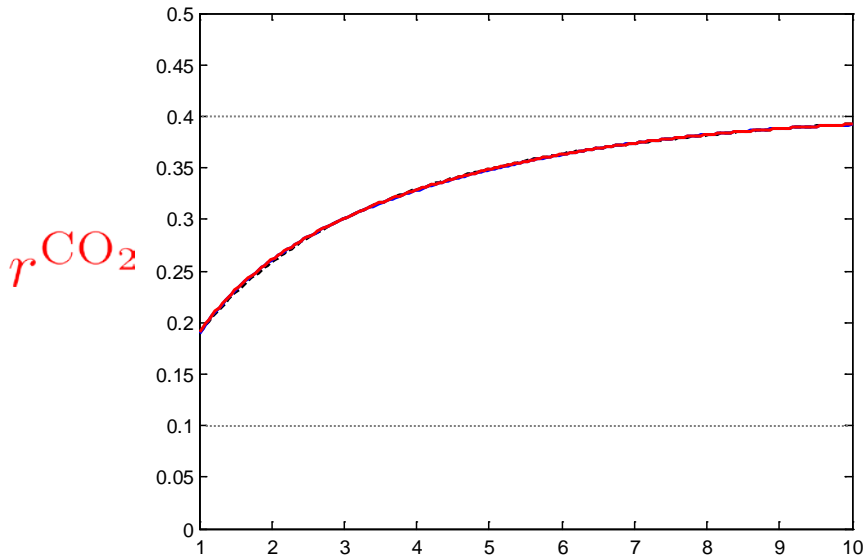


## Maximum



# COMPARISON – ITER 4

## Model comparison

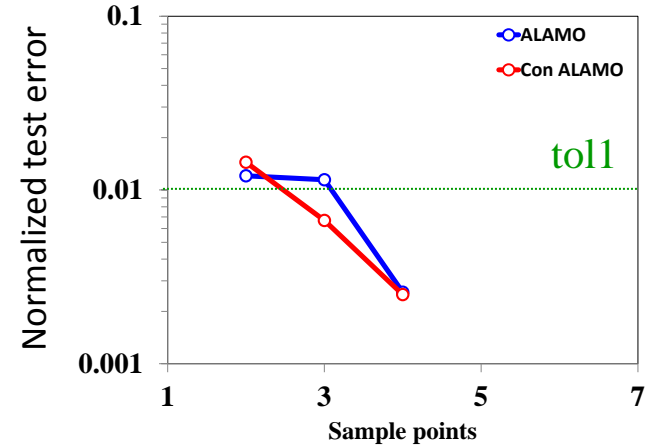


$L^{bed}$

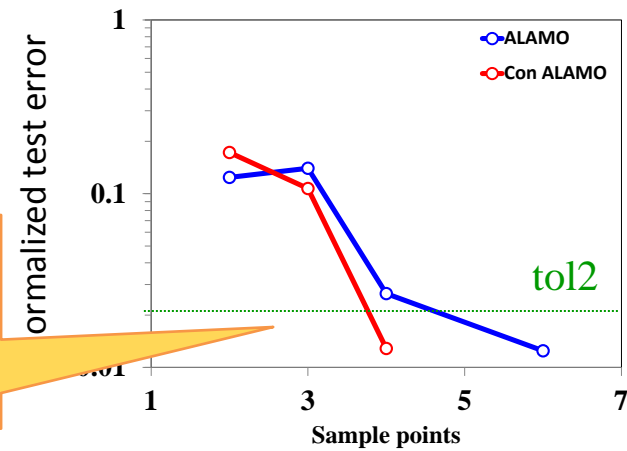
- ALAMO
- Constrained

**Standard ALAMO required 50% more sample points**

## Mean



## Maximum



# CONCLUSIONS

- **Constrained optimization provides a new avenue to provide a model with a priori information**
  - Include more information in your model without additional sampling
  - Reduced the sampling required for an accurate model
- **Ensure a more robust model by using output bounds as a “reality” check on the model**
- **Future work,**
  - More complex solution manifolds
    - *Nonlinear constraints on regressors*
    - *Nonlinear feasible domain for output variables*
  - Simultaneous model generation and constraints
    - *Restrictions implied by constraints on multiple outputs*
      - Ex: Sum-to-one constraints

# STANDARD BASIS FUNCTION SELECTION

$$\min SE = \sum_{i=1}^N \left| z_i - \sum_{j \in \mathcal{B}} \beta_j X_{ij} \right|$$

Find the model with the least error

$$\text{s.t. } \sum_{j \in \mathcal{B}} y_j = T$$

$$-U(1 - y_j) \leq \sum_{i=1}^N X_{ij} \left( z^i - \sum_{j \in \mathcal{B}} \beta_j X_{ij} \right) \leq U(1 - y_j) \quad j \in \mathcal{B}$$

$$\beta^l y_j \leq \beta_j \leq \beta^u y_j \quad j \in \mathcal{B}$$

$$y_j = \{0, 1\} \quad j \in \mathcal{B}$$

# BASIS FUNCTION SELECTION

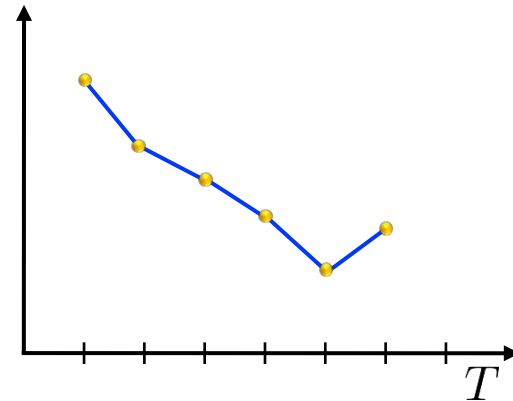
$$\min SE = \sum_{i=1}^N \left| z_i - \sum_{j \in \mathcal{B}} \beta_j X_{ij} \right|$$

$$\text{s.t. } \sum_{j \in \mathcal{B}} y_j = T$$

$$-U(1 - y_j) \leq \sum_{i=1}^N X_{ij} \left( z_i - \sum_{j \in \mathcal{B}} \beta_j X_{ij} \right) \leq U(1 - y_j) \quad j \in \mathcal{B}$$

$$\beta^l y_j \leq \beta_j \leq \beta^u y_j$$

$$y_j = \{0, 1\}$$



We will solve this model for increasing  $T$  until we determine a model

# BASIS FUNCTION SELECTION

$$\min \quad SE = \sum_{i=1}^N \left| z_i - \sum_{j \in \mathcal{B}} \beta_j X_{ij} \right|$$

$$\text{s.t.} \quad \sum_{j \in \mathcal{B}} y_j = T$$

$$-U(1 - y_j) \leq \sum_{i=1}^N X_{ij} \left( z^i - \sum_{j \in \mathcal{B}} \beta_j X_{ij} \right) \leq U(1 - y_j) \quad j \in \mathcal{B}$$

$$\beta^l y_j \leq \beta_j \leq \beta^u y_j \quad j \in \mathcal{B}$$

$$y_j = \{0, 1\} \quad j \in \mathcal{B}$$

