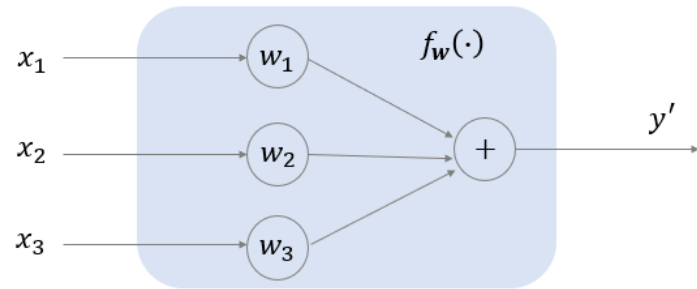# CS 4501/6501 Interpretable Machine Learning

## Interpretable Generalized Additive Models

Hanjie Chen, Yangfeng Ji
Department of Computer Science
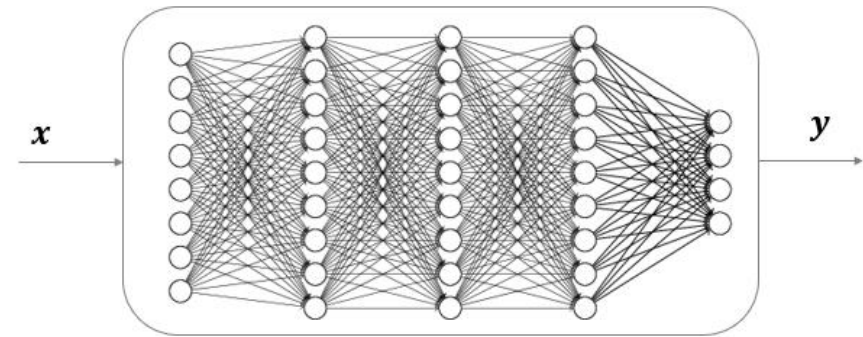University of Virginia
{hc9mx, yangfeng}@virginia.edu

# Interpretability

Bad performance
Good interpretability

Good performance
Bad interpretability



- Three parameters $(w_1, w_2, w_3)$

- $y' = w_1 x_1 + w_2 x_2 + w_3 x_3$

- Contributions:
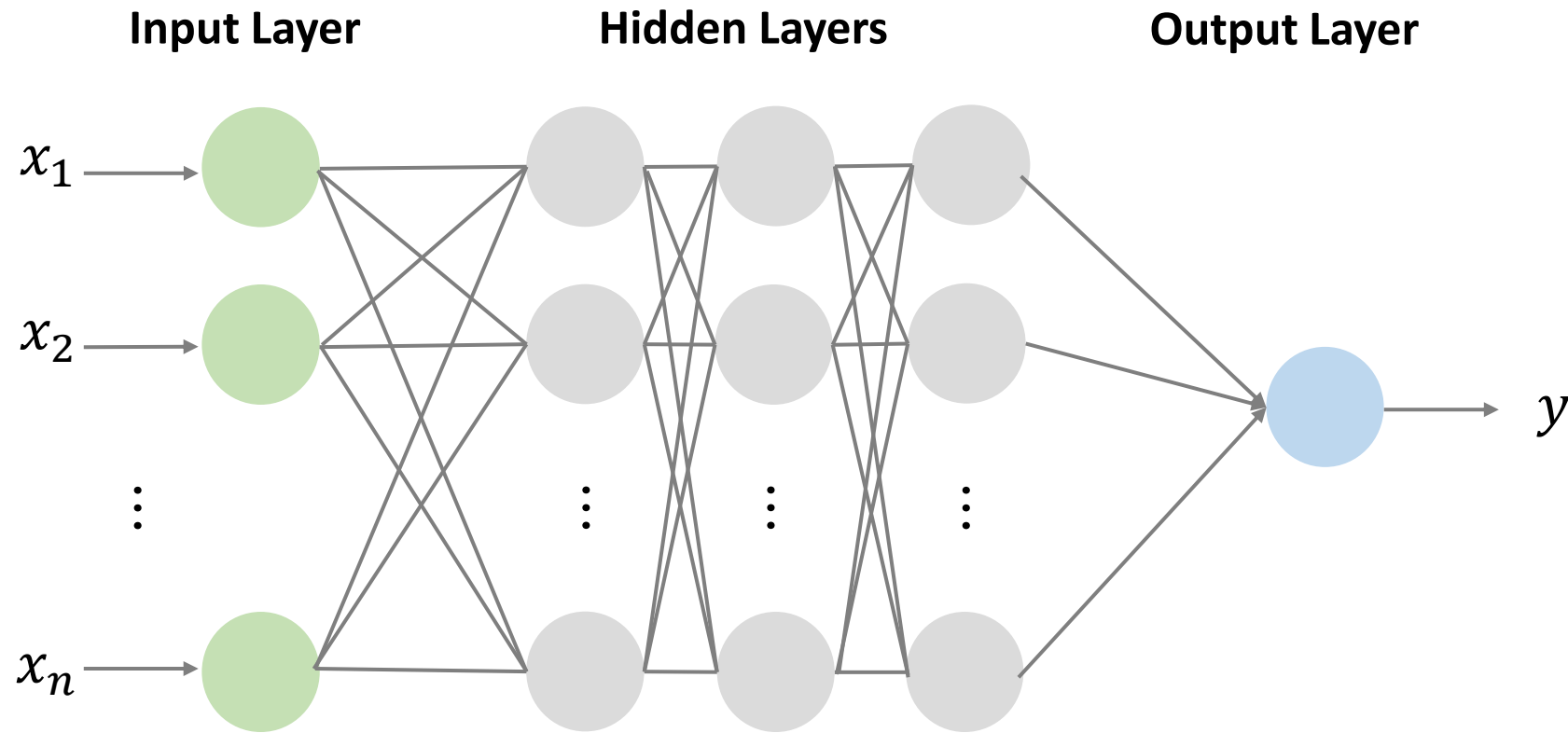
$$x_1 : w_1 x_1$$
$$x_2 : w_2 x_2$$
$$x_3 : w_3 x_3$$

- Millions of parameters

- $y' = f_w(x)$ (complex transformations)

- Model decision-making and feature attributions are unclear
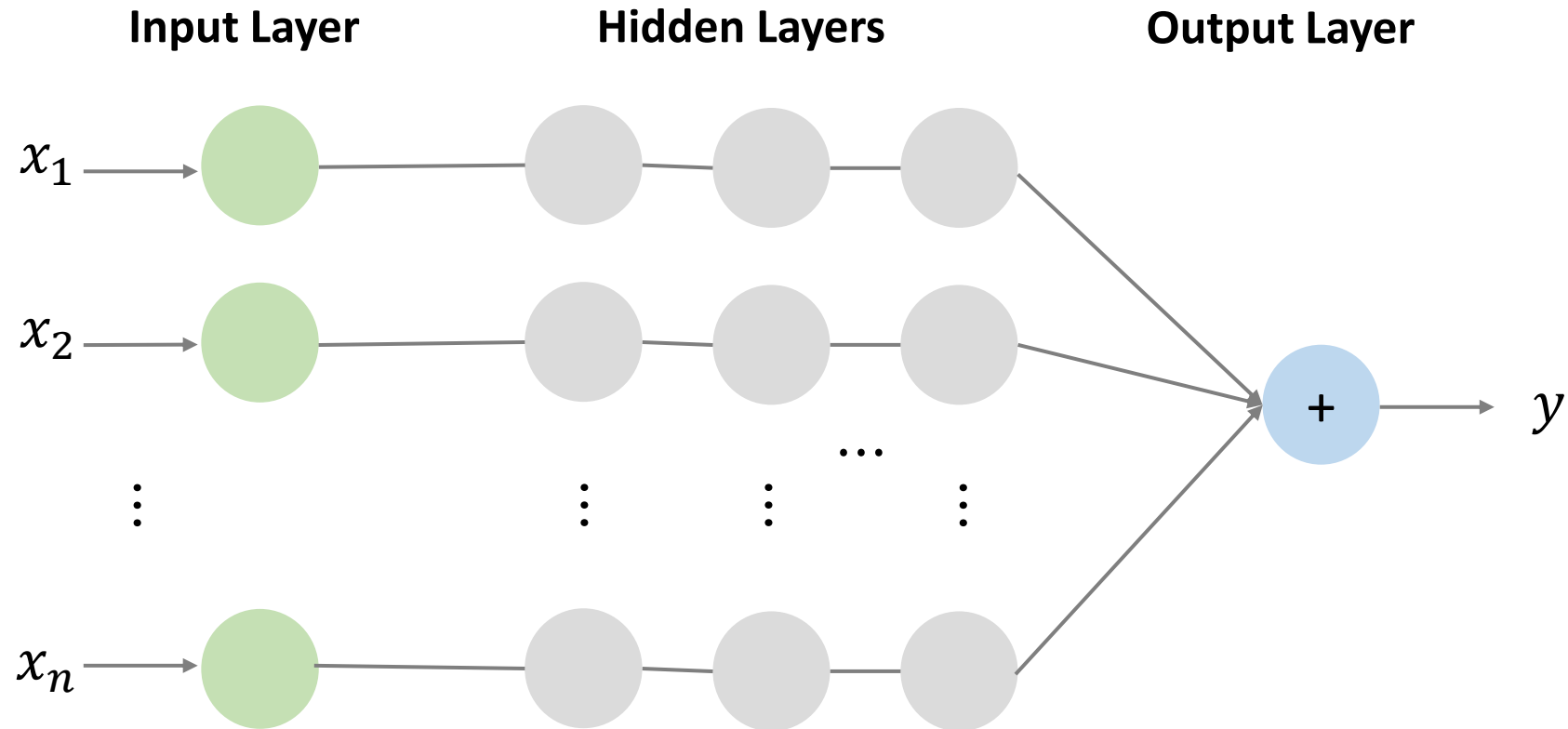
# Trade-off

The information of input features is mixed

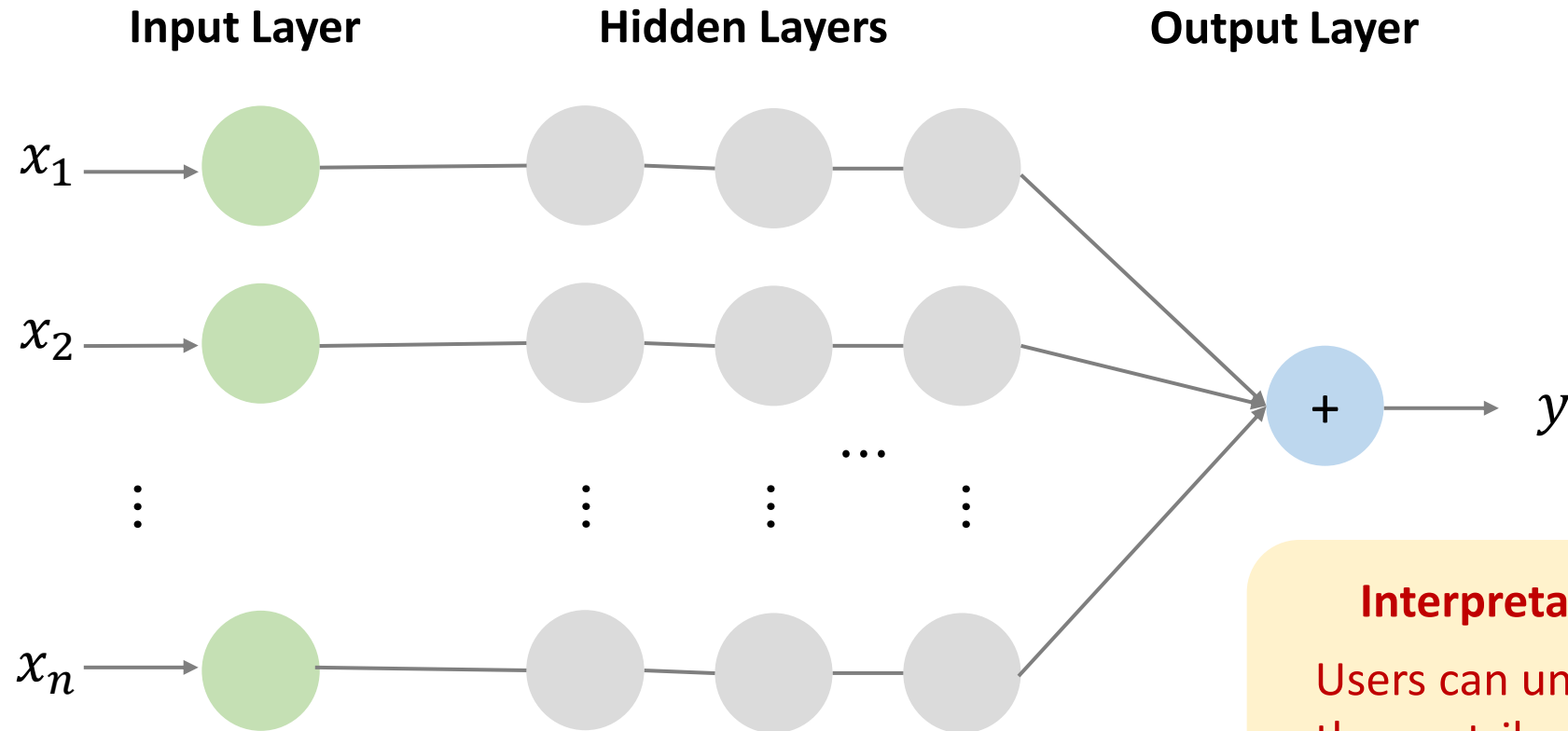# Trade-off

Keep the information of individual features "locally"

# Trade-off

Keep the information of individual features "locally"

# Trade-off

Generalized additive models (GAMs)

$$g(y) = f_1(x_1) + f_2(x_2) + \cdots + f_n(x_n)$$

- Permit complex relationships between individual features ($x_i$) and the

  target ($g(y)$)

- Exclude complex interactions between features

# GAM

Generalized additive models (GAMs)

$$g(y) = f_1(x_1) + f_2(x_2) + \cdots + f_n(x_n)$$

- $g(\cdot)$: link function

  - Identity: $g(y) = y$ $\longrightarrow$ Regression

  - Logistic function: $g(y)$ represents the probability on a class $\longrightarrow$ Classification

$$\frac{L}{1 + e^{-k(x-x_0)}}$$

$$(L = 1, k = 1, x_0 = 0)$$

# GAM

Generalized additive models (GAMs)

$$g(y) = f_1(x_1) + f_2(x_2) + \cdots + f_n(x_n)$$

- $f_i(\cdot)$: shape function

  - Splines

# GAM

Generalized additive models (GAMs)

$$g(y) = f_1(x_1) + f_2(x_2) + \cdots + f_n(x_n)$$

- $f_i(\cdot)$: shape function

  - Binary Trees

# GAM

Generalized additive models (GAMs)

$$g(y) = f_1(x_1) + f_2(x_2) + \cdots + f_n(x_n)$$

- $f_i(\cdot)$: shape function

  - Binary Trees

# GAM

Generalized additive models (GAMs)

$$g(y) = f_1(x_1) + f_2(x_2) + \cdots + f_n(x_n)$$

- $f_i(\cdot)$: shape function

  - Binary Trees

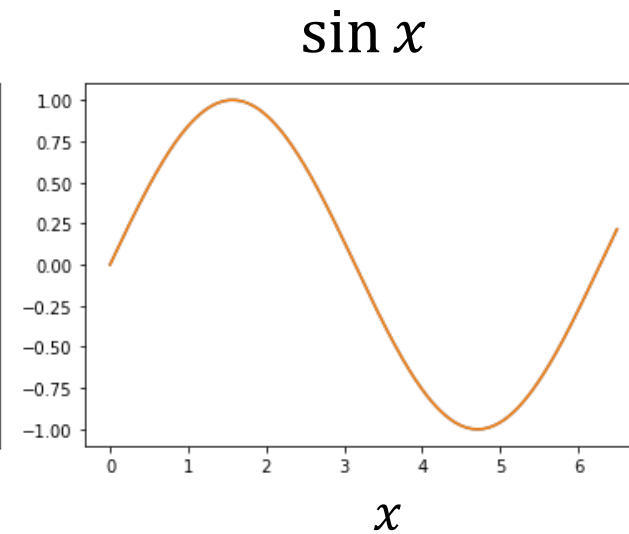For interpretability, we control tree complexity (nodes, leaves, depth)

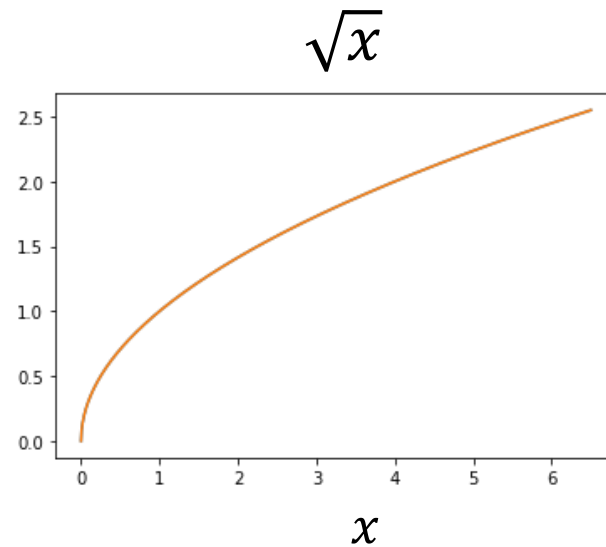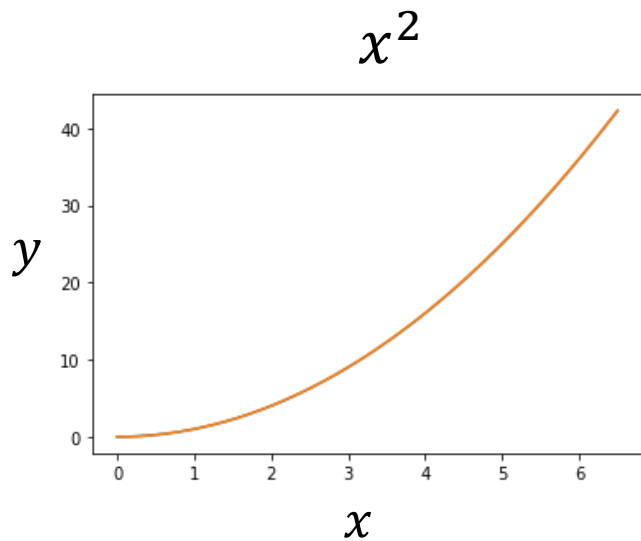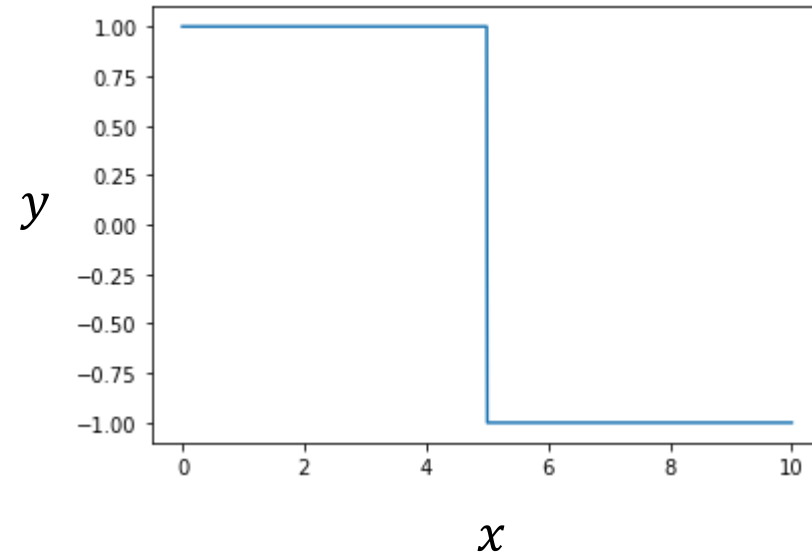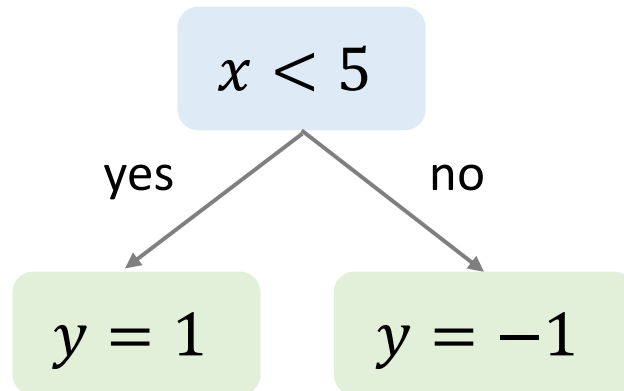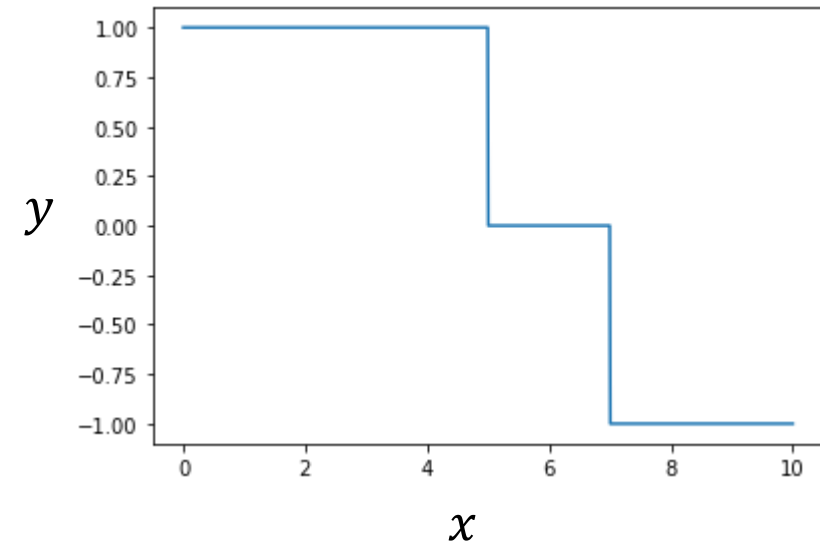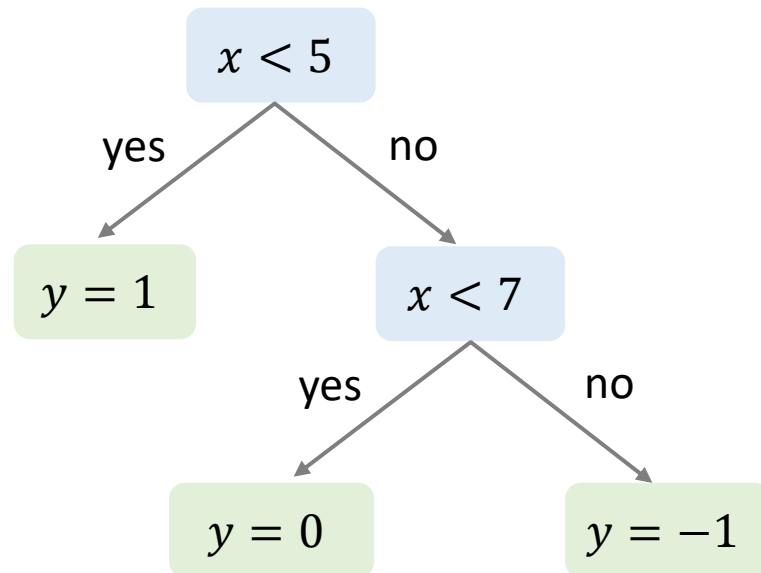# GAM

Generalized additive models (GAMs)

$$g(y) = f_1(x_1) + f_2(x_2) + \cdots + f_n(x_n)$$

- $f_i(\cdot)$: shape function

    - Bagged Trees (reduce the variance)

# GAM

Generalized additive models (GAMs)

$$g(y) = f_1(x_1) + f_2(x_2) + \cdots + f_n(x_n)$$

- Training

  - Shape functions for individual features

  - Learning methods

# Learning GAM

- Gradient Boosting

  - Learning tree or tree ensemble shape functions

---

**Algorithm** Gradient Boosting for GAM

---

1. $f_j \leftarrow 0, j = 1, \cdots, n$    Initialize all shape functions as zero

2. **for** $m = 1, \cdots, M$ **do**    Loop over M iterations

3.    **for** $j = 1, \cdots, n$ **do**    Loop over all features

4.    $\mathcal{R} \leftarrow \left\{ x_{ij}, y_i - \sum_k f_k \right\}_{i=1}^{N}$    Calculate residuals

5.    Learning shape function $S: x_j \to y$ using $\mathcal{R}$ as training data    Learn the one-dimensional function to predict the residuals

6.    $f_j \leftarrow f_j + S$    Update the shape function

---

# Learning GAM

- Gradient Boosting

  - Learning tree or tree ensemble shape functions

---

**Algorithm** Gradient Boosting for GAM

---

1.  $f_j \leftarrow 0, j = 1, \cdots, n$    Initialize all shape functions as zero

2.  **for** $m = 1, \cdots, M$ **do**    Loop over M iterations

3.     **for** $j = 1, \cdots, n$ **do**    Loop over all features

4.  $\mathcal{R} \leftarrow \left\{ x_{ij}, y_i - \sum_k f_k \right\}_{i=1}^{N}$    Calculate residuals

5.  Learning shape function $S: x_j \to y$ using $\mathcal{R}$ as training data    Learn the one-dimensional function to predict the residuals

6.  $f_j \leftarrow f_j + S$    Update the shape function

---

# Learning GAM

- Gradient Boosting

Training data $\{(\boldsymbol{x}_i, y_i)\}_{i=1}^{N}$

$\boldsymbol{x}_2$

| $i$ | $x_1$ | $x_2$ | $\cdots$ | $x_j$ | $\cdots$ | $x_n$ | $y$ |
|---|---|---|---|---|---|---|---|
| 1 | $x_{11}$ | $x_{12}$ | $\cdots$ | $x_{1j}$ | $\cdots$ | $x_{1n}$ | $y_1$ |
| 2 | $x_{21}$ | $x_{22}$ | $\cdots$ | $x_{2j}$ | $\cdots$ | $x_{2n}$ | $y_2$ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |
| $N$ | $x_{N1}$ | $x_{N2}$ | $\cdots$ | $x_{Nj}$ | $\cdots$ | $x_{Nn}$ | $y_N$ |

# Learning GAM

- Gradient Boosting

Training data $\{(\boldsymbol{x}_i, y_i)\}_{i=1}^N$

| $i$ | $x_1$ | $x_2$ | $\cdots$ | $x_j$ | $\cdots$ | $x_n$ | $y$ |
|---|---|---|---|---|---|---|---|
| 1 | $x_{11}$ | $x_{12}$ | $\cdots$ | $x_{1j}$ | $\cdots$ | $x_{1n}$ | $y_1$ |
| 2 | $x_{21}$ | $x_{22}$ | $\cdots$ | $x_{2j}$ | $\cdots$ | $x_{2n}$ | $y_2$ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |
| $N$ | $x_{N1}$ | $x_{N2}$ | $\cdots$ | $x_{Nj}$ | $\cdots$ | $x_{Nn}$ | $y_N$ |

# Learning GAM

- Gradient Boosting

Training data $\{(\boldsymbol{x}_i, y_i)\}_{i=1}^N$

| $i$ | $x_1$ | $x_2$ | $\cdots$ | $x_j$ | $\cdots$ | $x_n$ | $y$ |
|---|---|---|---|---|---|---|---|
| 1 | $x_{11}$ | $x_{12}$ | $\cdots$ | $x_{1j}$ | $\cdots$ | $x_{1n}$ | $y_1$ |
| 2 | $x_{21}$ | $x_{22}$ | $\cdots$ | $x_{2j}$ | $\cdots$ | $x_{2n}$ | $y_2$ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |
| $N$ | $x_{N1}$ | $x_{N2}$ | $\cdots$ | $x_{Nj}$ | $\cdots$ | $x_{Nn}$ | $y_N$ |

# Learning GAM

- Gradient Boosting

Training data $\{(\boldsymbol{x}_i, y_i)\}_{i=1}^N$

$f_j$

| $i$ | $x_1$ | $x_2$ | $\cdots$ | $x_j$ | $\cdots$ | $x_n$ | $y$ |
|---|---|---|---|---|---|---|---|
| 1 | $x_{11}$ | $x_{12}$ | $\cdots$ | $x_{1j}$ | $\cdots$ | $x_{1n}$ | $y_1$ |
| 2 | $x_{21}$ | $x_{22}$ | $\cdots$ | $x_{2j}$ | $\cdots$ | $x_{2n}$ | $y_2$ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |
| $N$ | $x_{N1}$ | $x_{N2}$ | $\cdots$ | $x_{Nj}$ | $\cdots$ | $x_{Nn}$ | $y_N$ |

Residuals

$\longrightarrow y_1 - \sum_k f_k$

$\longrightarrow y_2 - \sum_k f_k$

$\vdots$

$\longrightarrow y_N - \sum_k f_k$

(errors made by the current model)

# Learning GAM

- Gradient Boosting

  Update $f_j$ based on $\left\{\left(\underbrace{x_{ij}}_{x}, \underbrace{y_i - \sum_k f_k}_{y}\right)\right\}_{i=1}^{N}$

  - Learn a shape function S that fits: $x \to y$

  - Update $f_j \leftarrow f_j + S$

# Learning GAM

- Gradient Boosting

**Example**

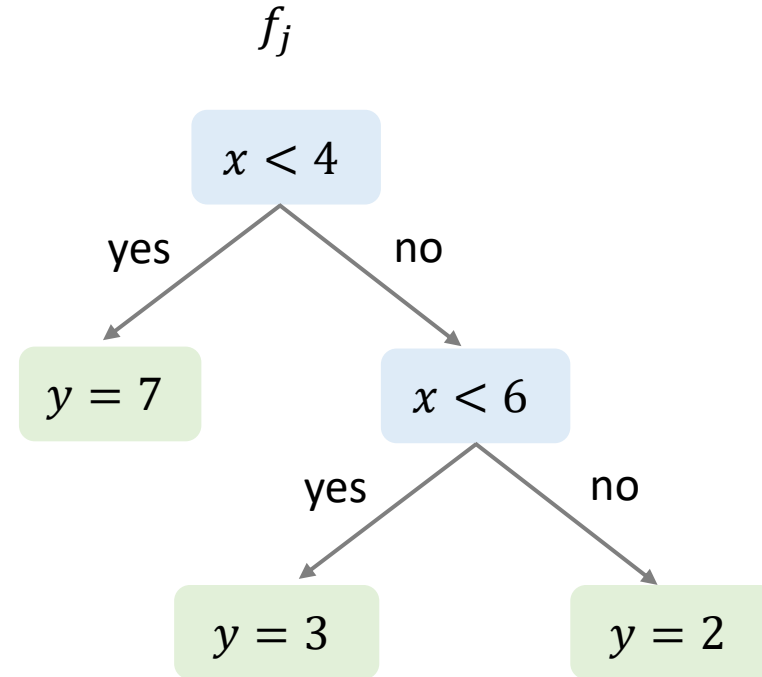| $x$ | $y$ |
|-----|-----|
| 1   | 8   |
| 5   | 5   |
| 3   | 8   |
| 9   | 7   |

Residuals

$8 - 7 = 1$

$5 - 3 = 2$

$8 - 7 = 1$

$7 - 2 = 5$

$f_j$



$x < 4$

yes — no

$y = 7$     $x < 6$

yes — no

$y = 3$     $y = 2$

# Learning GAM

- Gradient Boosting

$f_j$

**Example**

| $x$ | $y$ |
|-----|-----|
| 1 | 8 |
| 5 | 5 |
| 3 | 8 |
| 9 | 7 |

Residuals

$8 - 7 = 1$

$5 - 3 = 2$

$8 - 7 = 1$

$7 - 2 = 5$

$x < 4$

yes    no

$y = 7$    $x < 6$

yes    no

$y = 3$    $y = 2$

# Learning GAM

- Gradient Boosting

**Example**

| $x$ | $y$ |
|-----|-----|
| 1 | 8 |
| 5 | 5 |
| 3 | 8 |
| 9 | 7 |

Residuals

$8 - 7 = 1$

$5 - 3 = 2$

$8 - 7 = 1$

$7 - 2 = 5$



$S$

$x < 4$

yes     no

1

$x < 6$

yes     no

2     5

# Learning GAM

- Gradient Boosting

**Example**

Update $f_j \leftarrow f_j + S$

| $x$ | $y$ |
|:---:|:---:|
| 1 | 8 |
| 5 | 5 |
| 3 | 8 |
| 9 | 7 |

Residuals

$8 - (7 + 1) = 0$

$5 - (3 + 2) = 0$

$8 - (7 + 1) = 0$

$7 - (2 + 5) = 0$

$f_j$

```
        x < 4
      yes    no
   y = 7      x < 6
           yes    no
        y = 3      y = 2
```

$+$

$S$

```
        x < 4
      yes    no
     1         x < 6
            yes    no
          2         5
```

# Learning GAM

- Gradient Boosting

**Example**

Update $f_j \leftarrow f_j + S$

| $x$ | $y$ |
|-----|-----|
| 1 | 8 |
| 5 | 5 |
| 3 | 8 |
| 9 | 7 |

Residuals

$8 - (7 + 1) = \boxed{0}$
$5 - (3 + 2) = \boxed{0}$
$8 - (7 + 1) = \boxed{0}$
$7 - (2 + 5) = \boxed{0}$



$f_j$

$x < 4$
yes / no
$y = 7$ / $x < 6$
yes / no
$y = 3$ / $y = 2$

$+$

$S$

$x < 4$
yes / no
$1$ / $x < 6$
yes / no
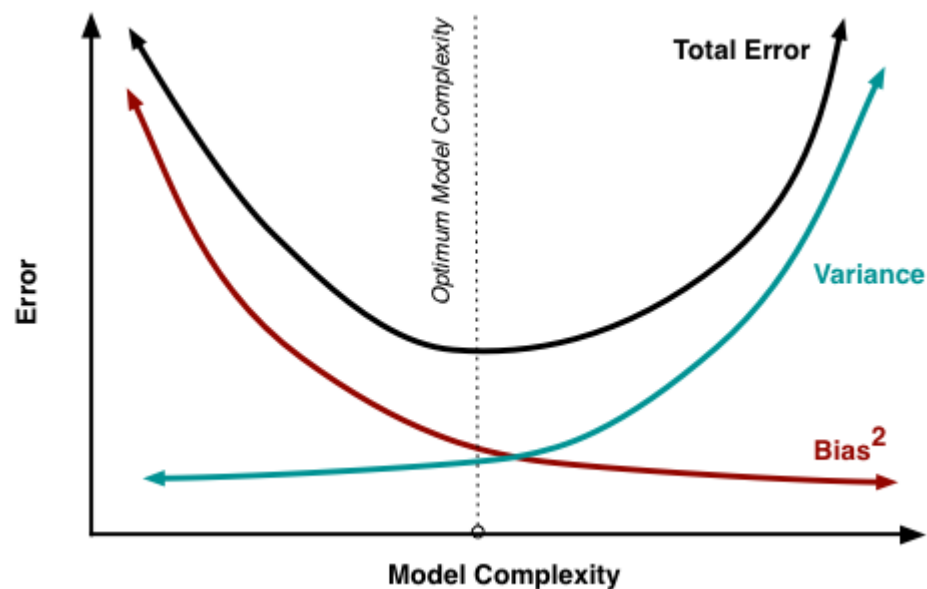$2$ / $5$

Do we learn a perfect model?

# Learning GAM

The model fits training data too well



We have low bias, but probably have high variance

Source: http://scott.fortmann-roe.com/docs/BiasVariance.html

# Learning GAM

- Gradient Boosting

Update $f_j \leftarrow f_j + \gamma \times S$

**Example**

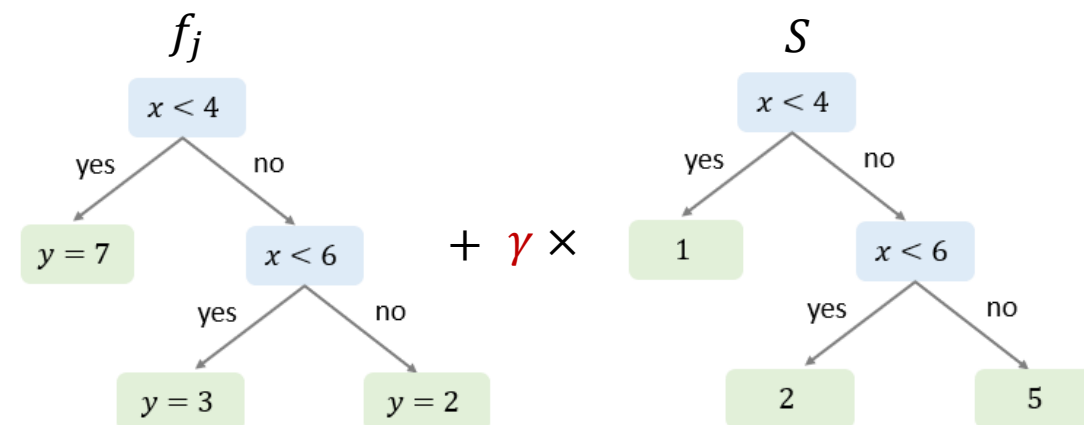| $x$ | $y$ |
|-----|-----|
| 1   | 8   |
| 5   | 5   |
| 3   | 8   |
| 9   | 7   |

Residuals

$8 - (7 + 0.1 \times 1) = 0.9$

$5 - (3 + 0.1 \times 2) = 1.8$

$8 - (7 + 0.1 \times 1) = 0.9$

$7 - (2 + 0.1 \times 5) = 4.5$



Add a learning rate to scale the contribution of the new tree

# Learning GAM

- Gradient Boosting

  - Learning tree or tree ensemble shape functions

---

**Algorithm** Gradient Boosting for GAM

---

1. $f_j \leftarrow 0, j = 1, \cdots, n$    Initialize all shape functions as zero

2. **for** $m = 1, \cdots, M$ **do**    Loop over M iterations

3.      **for** $j = 1, \cdots, n$ **do**    Loop over all features

4.      $\mathcal{R} \leftarrow \left\{ x_{ij}, y_i - \sum_k f_k \right\}_{i=1}^{N}$    Calculate residuals

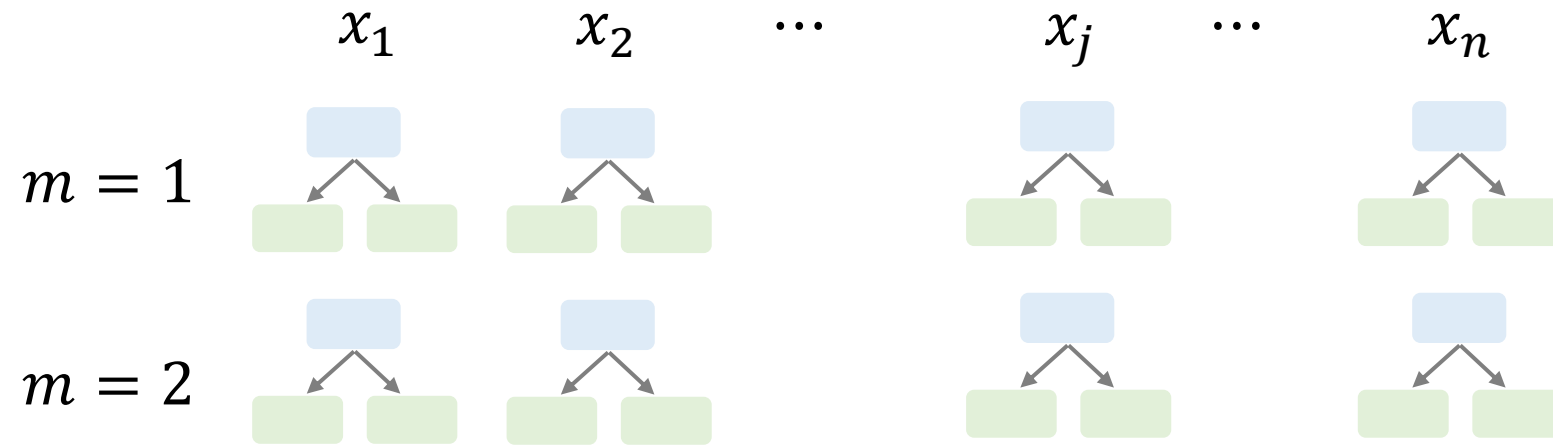5.      Learning shape function $S: x_j \rightarrow y$ using $\mathcal{R}$ as training data    Learn the one-dimensional function to predict the residuals

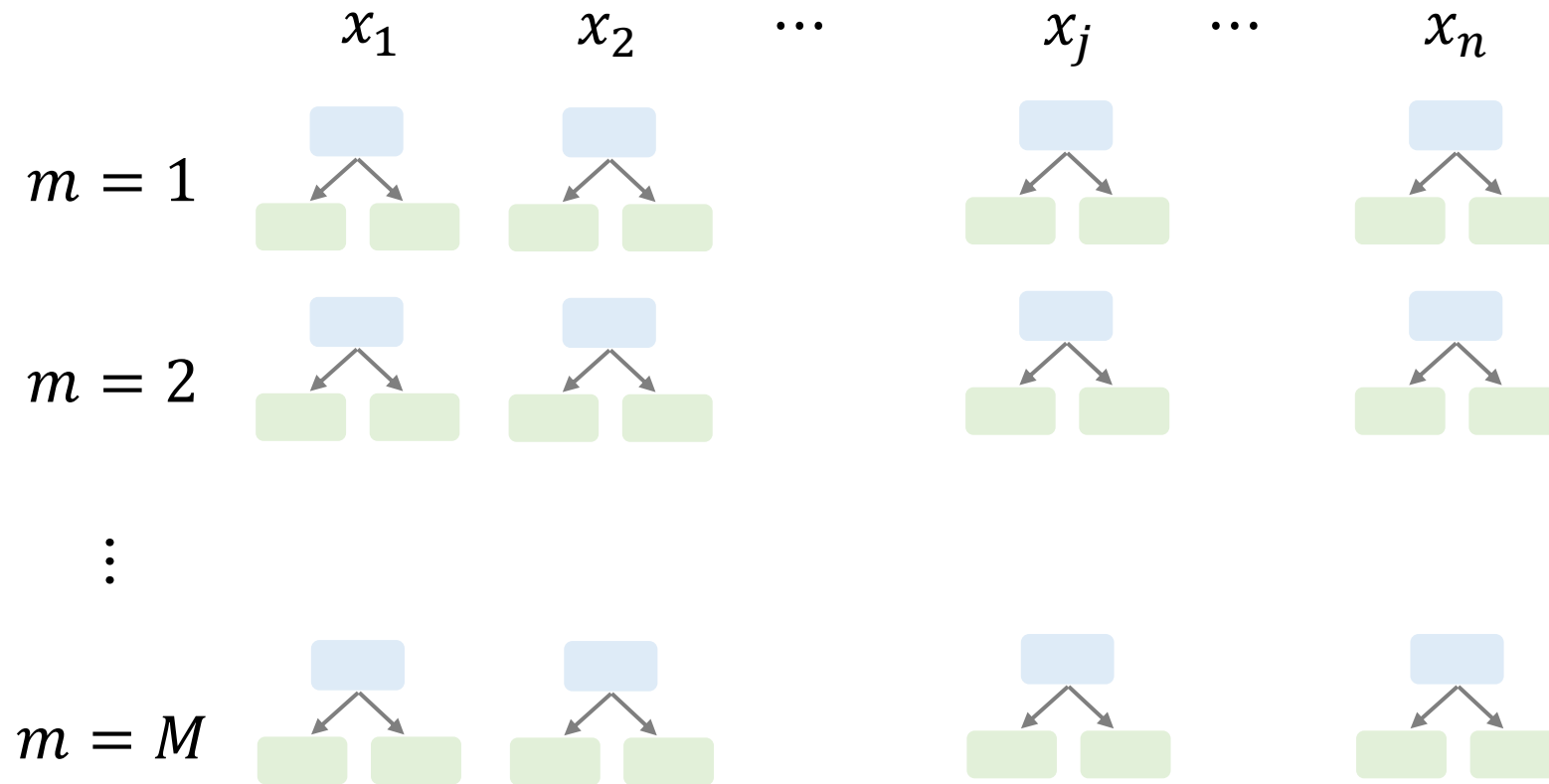6.      $f_j \leftarrow f_j + S$    Update the shape function
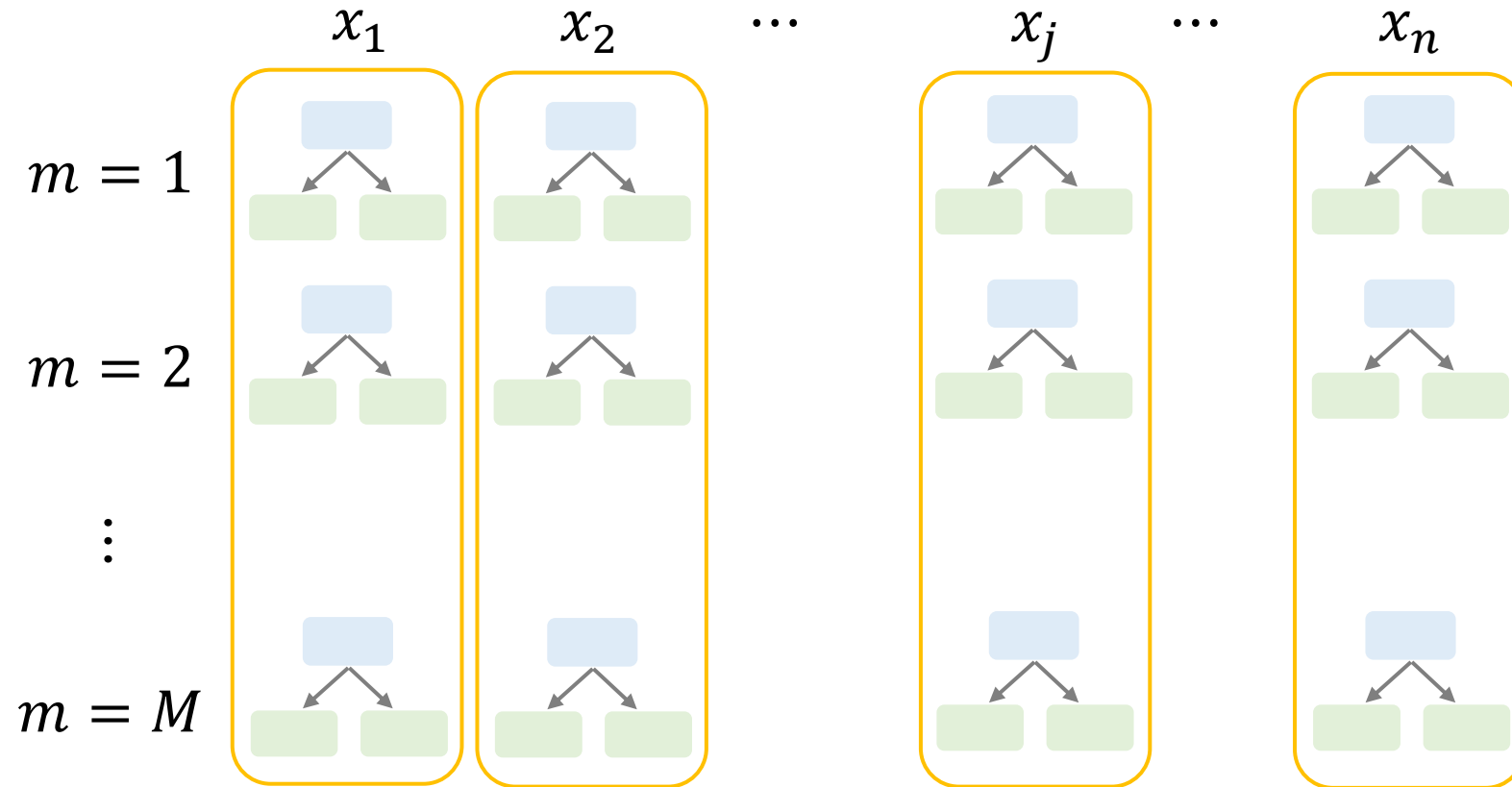
---

# Learning GAM

$$x_1 \quad\quad x_2 \quad\quad \cdots \quad\quad x_j \quad\quad \cdots \quad\quad x_n$$

$m = 1$

# Learning GAM

$x_1$  $x_2$  $\dots$  $x_j$  $\dots$  $x_n$

$m = 1$

$m = 2$

# Learning GAM

$$x_1 \quad x_2 \quad \cdots \quad x_j \quad \cdots \quad x_n$$

$m = 1$

$m = 2$

$\vdots$

$m = M$

# Learning GAM

# Question?

# Learning GAM

- Backfitting

  - Learning tree or tree ensemble shape functions

---

**Algorithm** Backfitting for GAM

---

1. $f_j \leftarrow 0, j = 1, \cdots, n$     Initialize all shape functions as zero

2. Learn $f_1$ using the training set $\{(x_{i1}, y_i)\}_{i=1}^{N}$

3. **for** $j = 2, \cdots, n$ **do**     Loop over rest features

4.     $\mathcal{R} \leftarrow \left\{ x_{ij}, y_i - \sum_{k=1}^{j-1} f_k \right\}_{i=1}^{N}$     Calculate residuals

5.     Learning shape function S: $x_j \rightarrow y$ using $\mathcal{R}$ as training data     Learn the one-dimensional function to predict the residuals

6.     $f_j \leftarrow S$     Update the shape function

7. Retrain $f_1$ based on the residuals of other $n - 1$ shape functions

---

# Learning GAM

- Least Squares

  - Learning spline shape functions
  - Reducing to fitting a linear model

$$g(y) = \beta_1 {x_1}^2 + \beta_2 \sqrt{x_2} + \cdots + \beta_n \sin x_n$$

$$\boldsymbol{y} = \boldsymbol{X\beta}$$

$$\underline{\boldsymbol{X_i}} = [{x_{i1}}^2, \sqrt{x_{i2}}, \cdots, \sin x_{in}]$$

$i^{th}$ example

$$\boldsymbol{\beta} = [\beta_1, \beta_2, \cdots, \beta_n]^T$$

**Objective**

$$\min \|\boldsymbol{y} - \boldsymbol{X\beta}\|_2$$

# Learning GAM

- Least Squares

  - Learning spline shape functions

  - Reducing to fitting a linear model

$$g(y) = \beta_1 x_1{}^2 + \beta_2 \sqrt{x_2} + \cdots + \beta_n \sin x_n$$

$$\boldsymbol{y} = \boldsymbol{X}\boldsymbol{\beta}$$

$$\underline{\boldsymbol{X_i}} = [x_{i1}{}^2, \sqrt{x_{i2}}, \cdots, \sin x_{in}]$$

$i^{th}$ example

$$\boldsymbol{\beta} = [\beta_1, \beta_2, \cdots, \beta_n]^T$$

**Objective**

$$\min \| \boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta} \|_2$$

Simple, but not flexible

# Summary

Generalized additive models (GAMs)

$$g(y) = f_1(x_1) + f_2(x_2) + \cdots + f_n(x_n)$$
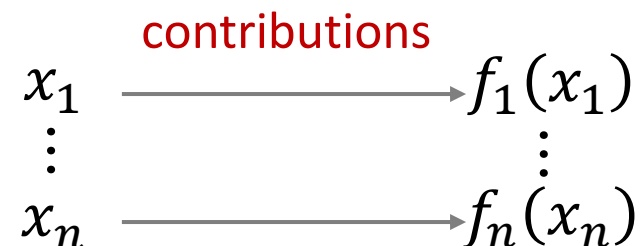
- Training

  - Shape functions for individual features: splines, trees, ensembles of trees

  - Learning methods: Least Squares, Gradient Boosting, Backfitting

# Summary

Generalized additive models (GAMs)

$$g(y) = f_1(x_1) + f_2(x_2) + \cdots + f_n(x_n)$$

- Training

  - Shape functions for individual features: splines, trees, ensembles of trees

  - Learning methods: Least Squares, Gradient Boosting, Backfitting
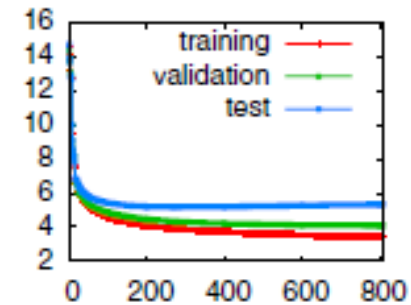
- Interpretability

# Application

- Dataset: "Concrete" (Blast Furnace Slag, Fly Ash, Superplasticizer…)

- Task: predicting the compressive strength of concrete

- Models:

| Shape Function | Least Squares | Gradient Boosting | Backfitting |
|---|---|---|---|
| Splines | P-LS/P-IRLS | BST-SP | BF-SP |
| Single Tree | N/A | BST-TR$x$ | BF-TR |
| Bagged Trees | N/A | BST-bagTR$x$ | BF-bagTR |
| Boosted Trees | N/A | BST-TR$x$ | BF-bstTR$x$ |
| Boosted Bagged Trees | N/A | BST-bagTR$x$ | BF-bbTR$x$ |

Lou, Yin, Rich Caruana, and Johannes Gehrke. "Intelligible models for classification and regression." *Proceedings of the 18th ACM SIGKDD international conference on Knowledge discovery and data mining*. 2012.
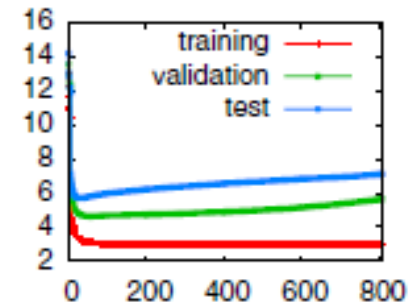
# Empirical Results

- GAMs perform better than linear or logistic regression (without feature shaping)

- Tree-based shaping methods are more accurate than spline-based methods

- Bagged-trees with 2-4 leaves as shape functions in combination with gradient boosting as learning method perform better

- Controlling the complexity of trees can avoid overfitting



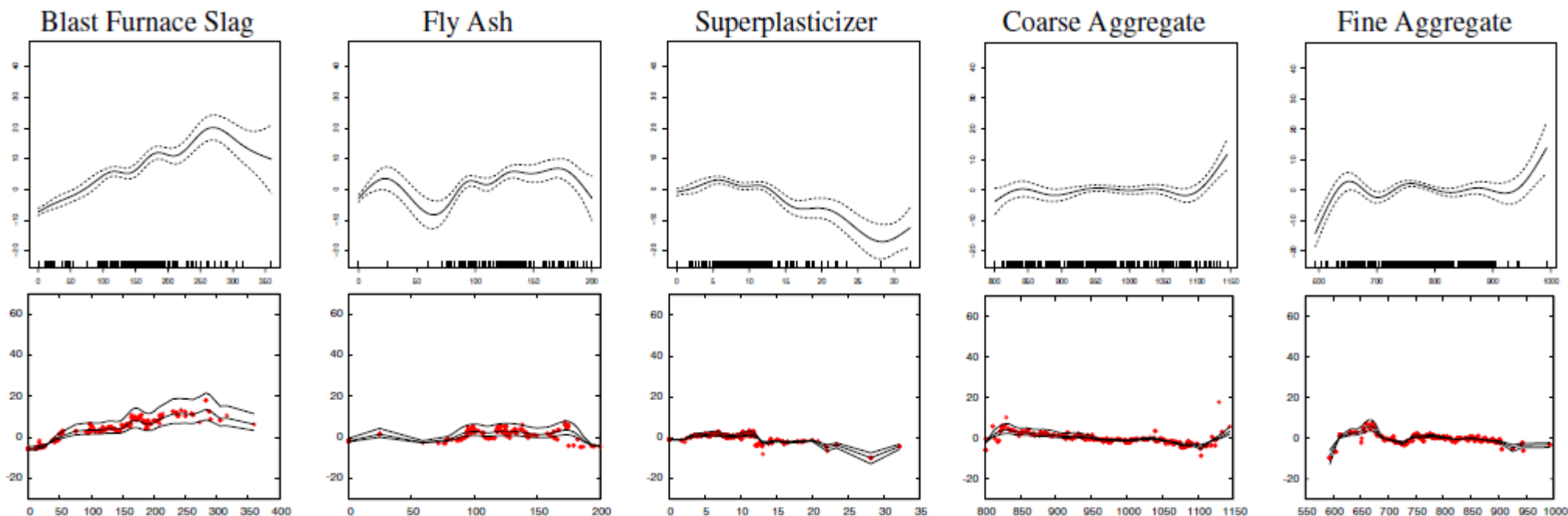(a) BST-bagTR2        (b) BST-bagTR16

(2 leaves)                    (16 leaves)

Lou, Yin, Rich Caruana, and Johannes Gehrke. "Intelligible models for classification and regression." *Proceedings of the 18th ACM SIGKDD international conference on Knowledge discovery and data mining*. 2012.

# Interpretation

Shapes of features for the "Concrete" dataset (versus the compressive strength of concrete)

(Splines)

(Bagged trees)



Lou, Yin, Rich Caruana, and Johannes Gehrke. "Intelligible models for classification and regression." *Proceedings of the 18th ACM SIGKDD international conference on Knowledge discovery and data mining*. 2012.
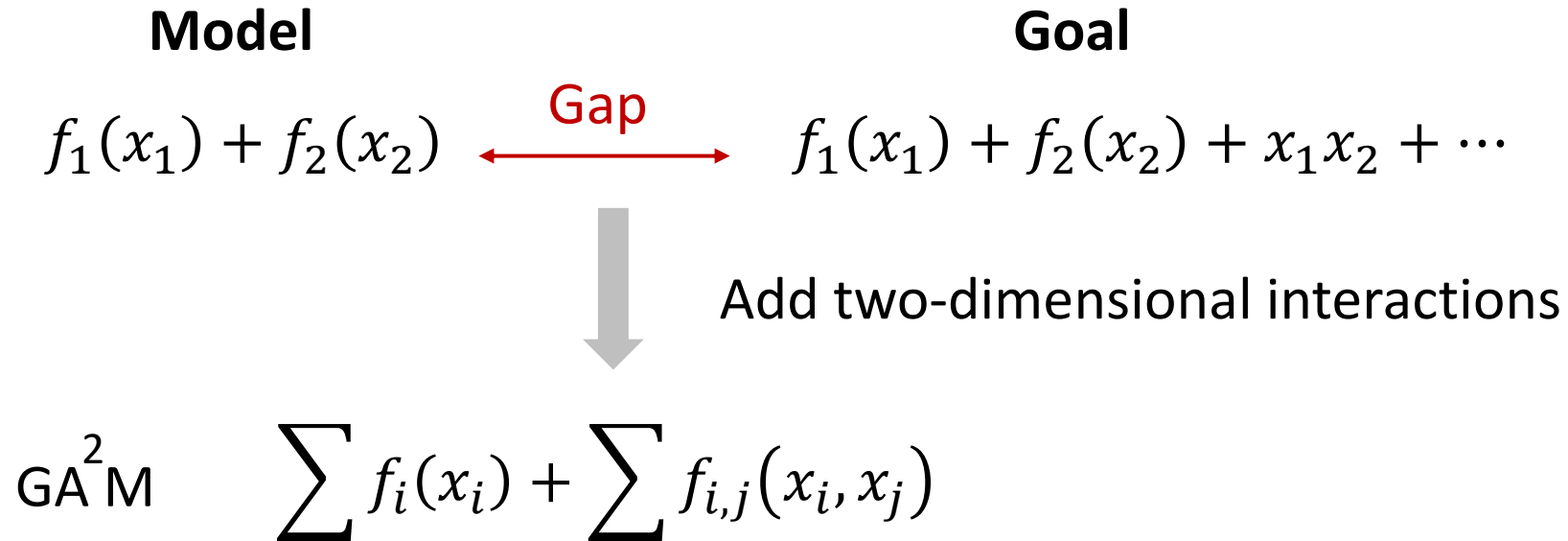
# Question?

# GA$^2$M

**Limitation**: GAMs do not consider feature dependency

<div>

**Model**                                    **Goal**

$$f_1(x_1) + f_2(x_2) \quad \xleftarrow{\text{Gap}}\rightarrow \quad f_1(x_1) + f_2(x_2) + x_1 x_2 + \cdots$$

</div>

# GA$^2$M

**Limitation**: GAMs do not consider feature dependency

<div align="center">

**Model**            **Goal**

$f_1(x_1) + f_2(x_2)$  ← Gap →  $f_1(x_1) + f_2(x_2) + x_1 x_2 + \cdots$

↓ Add two-dimensional interactions

GA$^2$M   $\sum f_i(x_i) + \sum f_{i,j}(x_i, x_j)$

</div>

# GA²M

**Definitions**

- Dataset $D = \{(\boldsymbol{x}_i, y_i)\}_{i=1}^{N}$

- $\boldsymbol{x}_i = [x_{i1}, \cdots, x_{in}]$ with $n$ features

- $y_i$ is the response

- $\boldsymbol{x} = (x_1, \cdots, x_n)$ denote the features in the dataset

- $U_1 = \{\{i\}|1 \leq i \leq n\}, U_2 = \{\{i,j\}|1 \leq i < j \leq n\}, U = U_1 \cup U_2$, i.e., $U$ contains all indices for all features and pairs of features

- For any $u \in U$, let $H_u$ denote the Hilbert space of $f_u(x_u)$

- $H = \sum_{u \in U} H_u, H_1 = \sum_{u \in U_1} H_u, H_2 = \sum_{u \in U_2} H_u$

# GA²M

GA²M

$$F(\boldsymbol{x}) = \sum_{u \in U} f_u(x_u)$$

**Objective**

$$\min_{F \in H} E[L(y, F(\boldsymbol{x}))]$$

$L$: non-negative convex loss function

regression      classification

Squared loss      Cross-entropy loss

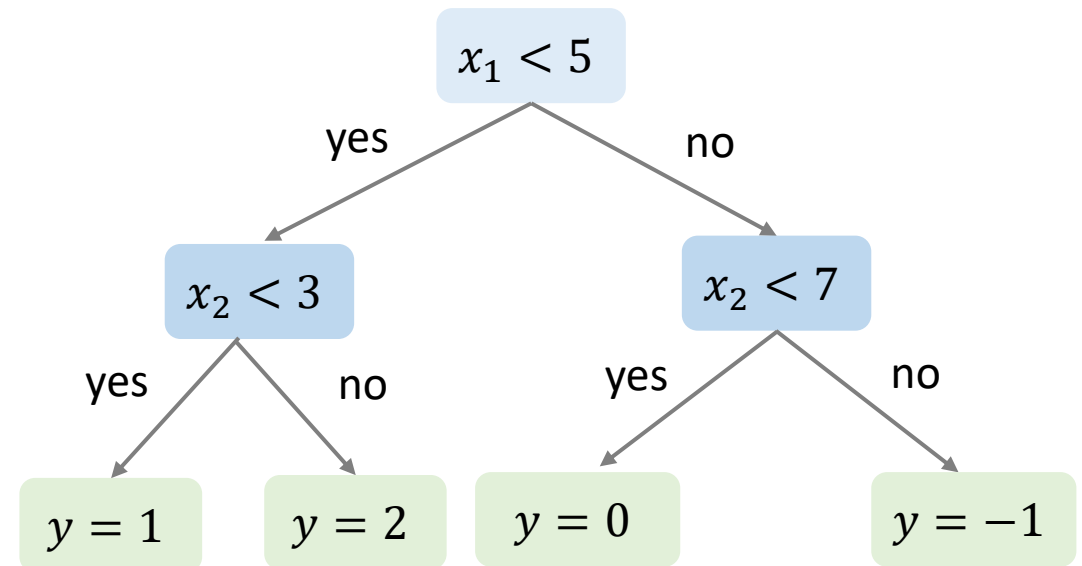$$(y - F(\boldsymbol{x}))^2 \qquad -y \log F(\boldsymbol{x}) - (1 - y) \log(1 - F(\boldsymbol{x}))$$

# GA²M

- We have known how to learn shape functions for GAMs

- Applicable to two-dimensional shape functions $f_u$, $u = \{i, j\}$

Splines

$$f_{1,2} = x_1 x_2$$



Trees



$x_1 < 5$

yes     no

$x_2 < 3$     $x_2 < 7$

yes   no    yes   no

$y = 1$   $y = 2$   $y = 0$   $y = -1$

GA$^2$M

**Challenge**

$$\sum f_i(x_i) + \sum f_{i,j}(x_i, x_j)$$

$n$ features $\longrightarrow$ $O(n^2)$ features interactions

How to find true
feature interactions?

# GA²M

**Algorithm** GA²M

1. $S \leftarrow \emptyset$      The set of the selected pairs
2. $Z \leftarrow U_2$      The set of the remaining pairs
3. **While** not converge **do**
4.      $F \leftarrow \arg \min\limits_{F \in H_1 + \sum_{u \in S} H_u} \frac{1}{2} E\left[(y - F(\boldsymbol{x}))^2\right]$
5.      $R \leftarrow y - F(\boldsymbol{x})$
6.      **for** all $u \in Z$ **do**
7.         $F_u \leftarrow E[R \mid x_u]$
8.      $u^* \leftarrow \arg \min\limits_{u \in Z} \frac{1}{2} E[(R - f_u(x_u))^2]$
9.      $S \leftarrow S \cup \{u^*\}$
10.      $Z \leftarrow Z - \{u^*\}$

# GA²M

**Algorithm** GA²M

1. $S \leftarrow \emptyset$      The set of the selected pairs
2. $Z \leftarrow U_2$      The set of the remaining pairs
3. **While** not converge **do**
4.     $F \leftarrow \arg\min\limits_{F \in H_1 + \sum_{u \in S} H_u} \dfrac{1}{2} E\left[(y - F(\boldsymbol{x}))^2\right]$    The best additive model $F$ so far in Hilbert space $H_1 + \sum_{u \in S} H_u$
5.     $R \leftarrow y - F(\boldsymbol{x})$
6.     **for** all $u \in Z$ **do**
7.        $F_u \leftarrow E[R \mid x_u]$
8.     $u^* \leftarrow \arg\min\limits_{u \in Z} \dfrac{1}{2} E[(R - f_u(x_u))^2]$
9.     $S \leftarrow S \cup \{u^*\}$
10.    $Z \leftarrow Z - \{u^*\}$

Learning shape functions for all single features ($f_i(x_i)$) and the selected feature pairs ($f_{i,j}(x_i, x_j)$). When $S = \emptyset$, $F$ is the GAM.

# GA$^2$M

---

**Algorithm  GA$^2$M**

---

1. $S \leftarrow \emptyset$       The set of the selected pairs

2. $Z \leftarrow U_2$       The set of the remaining pairs

3. **While** not converge **do**

4.     $F \leftarrow \arg\min\limits_{F \in H_1 + \sum_{u \in S} H_u} \frac{1}{2} E\left[(y - F(\boldsymbol{x}))^2\right]$    The best additive model $F$ so far in Hilbert space $H_1 + \sum_{u \in S} H_u$

5.     $R \leftarrow y - F(\boldsymbol{x})$     Calculate residual

6.     **for** all $u \in Z$ **do**

7.        $F_u \leftarrow E[R \mid x_u]$

8.     $u^* \leftarrow \arg\min\limits_{u \in Z} \frac{1}{2} E[(R - f_u(x_u))^2]$

9.     $S \leftarrow S \cup \{u^*\}$

10.    $Z \leftarrow Z - \{u^*\}$
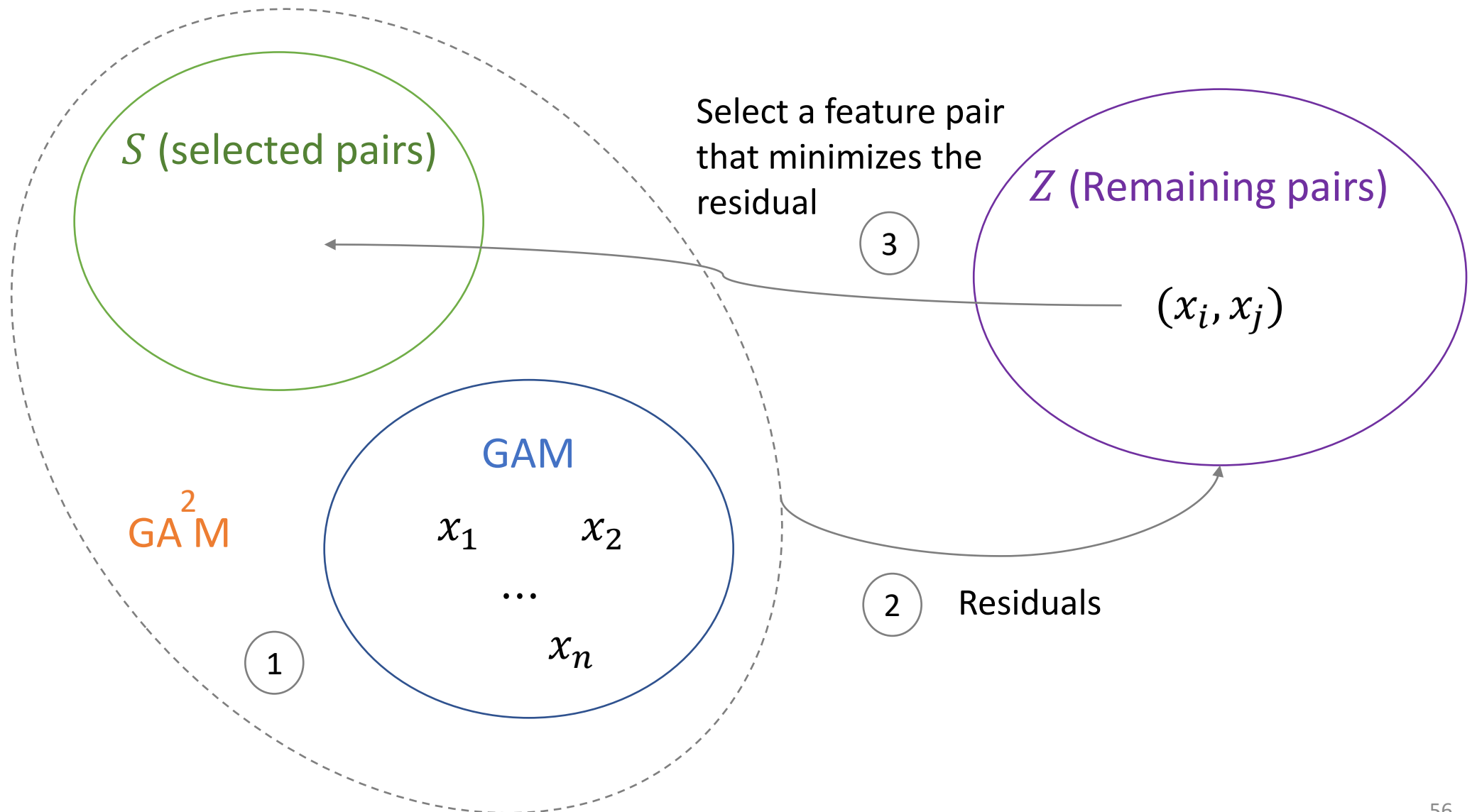
---

# GA$^2$M

---

**Algorithm**   GA$^2$M

---

1. $S \leftarrow \emptyset$          The set of the selected pairs

2. $Z \leftarrow U_2$        The set of the remaining pairs

3. **While** not converge **do**

4.     $F \leftarrow \arg\min\limits_{F \in H_1 + \sum_{u \in S} H_u} \frac{1}{2} E\left[(y - F(\boldsymbol{x}))^2\right]$     The best additive model $F$ so far in Hilbert space $H_1 + \sum_{u \in S} H_u$

5.     $R \leftarrow y - F(\boldsymbol{x})$     Calculate residual

6.     **for** all $u \in Z$ **do**       Loop over all remaining feature pairs

7.         $F_u \leftarrow E[R \mid x_u]$       Build an interaction model on the residual

8.     $u^* \leftarrow \arg\min\limits_{u \in Z} \frac{1}{2} E[(R - f_u(x_u))^2]$

    Learning a shape function for each feature pair

9.     $S \leftarrow S \cup \{u^*\}$

10.     $Z \leftarrow Z - \{u^*\}$

---

# GA$^2$M

**Algorithm** GA$^2$M

1. $S \leftarrow \emptyset$      The set of the selected pairs

2. $Z \leftarrow U_2$      The set of the remaining pairs

3. **While** not converge **do**

4. $F \leftarrow \arg \min\limits_{F \in H_1 + \sum_{u \in S} H_u} \frac{1}{2} E\left[(y - F(\boldsymbol{x}))^2\right]$     The best additive model $F$ so far in Hilbert space $H_1 + \sum_{u \in S} H_u$

5. $R \leftarrow y - F(\boldsymbol{x})$     Calculate residual

6.     **for** all $u \in Z$ **do**     Loop over all remaining feature pairs

7.         $F_u \leftarrow E[R \mid x_u]$     Build an interaction model on the residual

8. $u^* \leftarrow \arg \min\limits_{u \in Z} \frac{1}{2} E[(R - f_u(x_u))^2]$    Select the best feature pair

9. $S \leftarrow S \cup \{u^*\}$

10. $Z \leftarrow Z - \{u^*\}$

# GA$^2$M

---

**Algorithm** GA$^2$M

---

1. $S \leftarrow \emptyset$      The set of the selected pairs

2. $Z \leftarrow U_2$      The set of the remaining pairs

3. **While** not converge **do**

4. $\quad F \leftarrow \arg\min\limits_{F \in H_1 + \sum_{u \in S} H_u} \frac{1}{2} E\left[(y - F(\boldsymbol{x}))^2\right]$    The best additive model $F$ so far in Hilbert space $H_1 + \sum_{u \in S} H_u$

5. $\quad R \leftarrow y - F(\boldsymbol{x})$      Calculate residual

6. $\quad$ **for** all $u \in Z$ **do**      Loop over all remaining feature pairs

7. $\quad\quad F_u \leftarrow E[R \mid x_u]$      Build an interaction model on the residual

8. $\quad u^* \leftarrow \arg\min\limits_{u \in Z} \frac{1}{2} E[(R - f_u(x_u))^2]$    Select the best feature pair

9. $\quad S \leftarrow S \cup \{u^*\}$      Put the best feature pair in $S$

10. $\quad Z \leftarrow Z - \{u^*\}$      Remove that from $Z$

---

# GA$^2$M

---

**Algorithm**  GA$^2$M

---

1. $S \leftarrow \emptyset$ 　　　The set of the selected pairs

2. $Z \leftarrow U_2$ 　　　The set of the remaining pairs

3. **While** not converge **do**

4. 　$F \leftarrow \arg \min\limits_{F \in H_1 + \sum_{u \in S} H_u} \frac{1}{2} E\left[(y - F(\boldsymbol{x}))^2\right]$ 　The best additive model $F$ so far in Hilbert space $H_1 + \sum_{u \in S} H_u$

5. 　$R \leftarrow y - F(\boldsymbol{x})$ 　Calculate residual

6. 　**for** all $u \in Z$ **do** 　　Loop over all remaining feature pairs ⟶ $O(n^2)$

7. 　　$F_u \leftarrow E[R \mid x_u]$ 　　Build an interaction model on the residual

8. 　　$u^* \leftarrow \arg \min\limits_{u \in Z} \frac{1}{2} E[(R - f_u(x_u))^2]$ 　Select the best feature pair

9. 　　$S \leftarrow S \cup \{u^*\}$ 　Put the best feature pair in $S$

10. 　$Z \leftarrow Z - \{u^*\}$ 　Remove that from $Z$

See a fast interaction detection algorithm in [Lou et al., 2013]

---

55

$GA^2M$

# Question?

# Application

"Intelligible models for healthcare: Predicting pneumonia risk and hospital 30-day readmission"

Caruana et al., KDD 2015

# Background

- In the mid 90's, a project was funded by Cost-Effective HealthCare (CEHC) to evaluate the application of machine learning to important problems in healthcare such as predicting pneumonia risk
- **Goal**: predict the probability of death (POD) for patients with pneumonia
- High-risk: patients could be admitted to the hospital
- Low-risk: patients were treated as outpatients

# Background

- In the mid 90's, a project was funded by Cost-Effective HealthCare (CEHC) to evaluate the application of machine learning to important problems in healthcare such as predicting pneumonia risk
- **Goal**: predict the probability of death (POD) for patients with pneumonia
- High-risk: patients could be admitted to the hospital
- Low-risk: patients were treated as outpatients

**Models**

Logistic regression

Rule-based learning

k-nearest neighbor

Neural networks

...

# Background

- In the mid 90's, a project was funded by Cost-Effective HealthCare (CEHC) to evaluate the application of machine learning to important problems in healthcare such as predicting pneumonia risk
- **Goal**: predict the probability of death (POD) for patients with pneumonia
- High-risk: patients could be admitted to the hospital
- Low-risk: patients were treated as outpatients

### Models

Logistic regression

Rule-based learning

k-nearest neighbor

Neural networks        AUC=0.86 (Best performance)

…

# Background

- In the mid 90's, a project was funded by Cost-Effective HealthCare (CEHC) to evaluate the application of machine learning to important problems in healthcare such as predicting pneumonia risk
- **Goal**: predict the probability of death (POD) for patients with pneumonia
- High-risk: patients could be admitted to the hospital
- Low-risk: patients were treated as outpatients

### Models

Logistic regression      AUC=0.77 (safer to use on patients)

Rule-based learning

k-nearest neighbor

~~Neural networks~~      ~~AUC=0.86 (Best performance)~~

...

# Problem

The rule-based model learned a rule:

$$\text{HasAsthama}(x) \implies \text{LowerRisk}(x)$$

Rule-based models are interpretable
```
if (chance_of_rain > 0.75)
{ umbrella <- "yes" }
else { umbrella <- "no" }
```

# Problem

The rule-based model learned a rule:

$$HasAsthama(x) \implies LowerRisk(x)$$

counterintuitive

# Problem

The rule-based model learned a rule:

$$HasAsthama(x) \implies LowerRisk(x)$$

counterintuitive

The model captures a true pattern in the training data:

Patients:
asthma +
pneumonia

Hospital (ICU)

The aggressive care lowered the risk of dying from pneumonia

# Problem

The rule-based model learned a rule:

~~HasAsthama(x) $\implies$ LowerRisk(x)~~

Asthmatics have much higher risk!

It would be dangerous if the model predicts low risk on patients who have not been hospitalized

# Problem

How about other potential patterns?

$$\text{Pregnancy} \implies \text{Lower Risk ?}$$

# Problem

How about other potential patterns?

Pregnancy $\implies$ Lower Risk ?

MUST understand ML models in healthcare.
Otherwise, models may hurt patients
because of true patterns in data!

# Generalized Additive Models

- Better prediction performance than logistic regression
  (capture more data patterns)

- Interpretable

# Case Study: Pneumonia Risk

- There are 46 features describing each patient

- Bagged trees with gradient boosting

| Patient-history findings | | | |
|---|---|---|---|
| chronic lung disease | - | age | C |
| re-admission to hospital | - | gender | - |
| admitted through ER | - | diabetes mellitus | - |
| admitted from nursing home | - | asthma | - |
| congestive heart failure | - | cancer | - |
| ischemic heart disease | - | number of diseases | C |
| cerebrovascular disease | - | history of seizures | - |
| chronic liver disease | - | renal failure | - |
| history of chest pain | - | | |

| Physical examination findings | | | |
|---|---|---|---|
| diastolic blood pressure | C | wheezing | - |
| gastrointestinal bleeding | - | stridor | - |
| respiration rate | C | heart murmur | - |
| altered mental status | - | temperature | C |
| heart rate | C | | |

| Laboratory findings | | | |
|---|---|---|---|
| liver function tests | - | BUN level | C |
| glucose level | C | creatinine level | C |
| potassium level | C | albumin level | C |
| hematocrit | C | WBC count | C |
| percentage bands | C | pH | C |
| pO2 | C | pCO2 | C |
| sodium level | C | | |

| Chest X-ray findings | | | |
|---|---|---|---|
| positive chest x-ray | - | lung infiltrate | - |
| pleural effusion | - | pneumothorax | - |
| cavitation/empyema | - | chest mass | - |
| lobe or lung collapse | - | | |

# Prediction Performance

AUC for different learning methods

| Model | Pneumonia |
|-------|-----------|
| Logistic Regression | 0.8432 |
| GAM | 0.8542 |
| GA$^2$M | 0.8576 |
| Random Forests | 0.8460 |
| LogitBoost | 0.8493 |

**AUC: Area Under the ROC Curve**

$$TPR = \frac{TP}{TP + FN}$$

$$FPR = \frac{FP}{FP + TN}$$

# Interpretation



age

asthma

BUN level

(Blood Urea Nitrogen)

cancer

Older people have
higher risk

# Interpretation



age



asthma

GAMs also found the
pattern: asthma lowers
the risk



BUN level
(Blood Urea Nitrogen)



cancer

# Interpretation



age



asthma



BUN level

(Blood Urea Nitrogen)



cancer

GAMs also found the
pattern: asthma lowers
the risk

↓

Repair: eliminate this
term

# Interpretation



age     asthma     BUN level     cancer

(Blood Urea Nitrogen)

Most patients have BUN=0 because, as in many medical datasets, if the variable is not measured or assumed normal it is coded as 0

# Interpretation



age

asthma
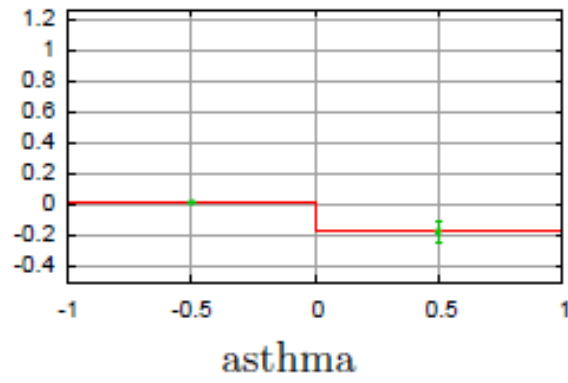
BUN level

(Blood Urea Nitrogen)

cancer

BUN levels below 30 appear to be low risk, while levels from 50-200 indicate higher risk
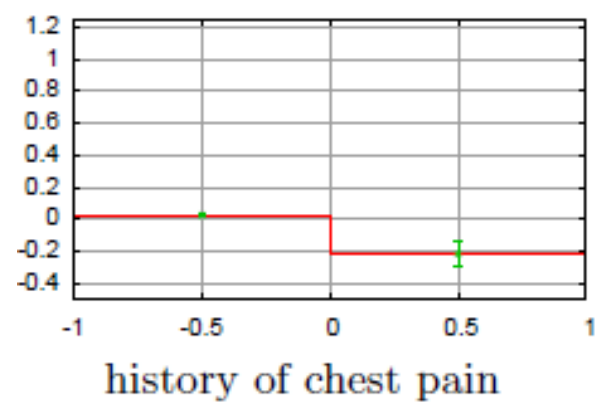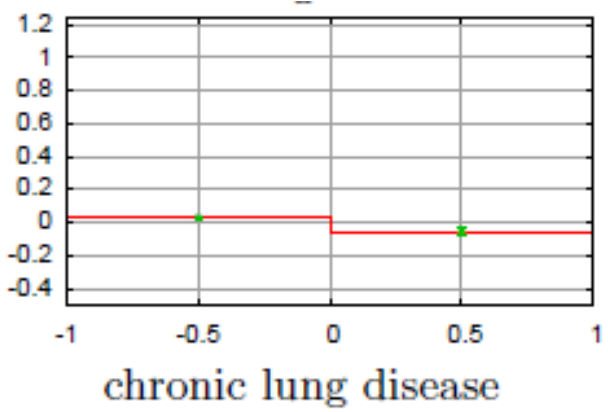
# Interpretation



age



asthma



BUN level
(Blood Urea Nitrogen)



cancer

Having cancer significantly increases the risk of dying from pneumonia

# Interpretation



age     asthma     BUN level     cancer
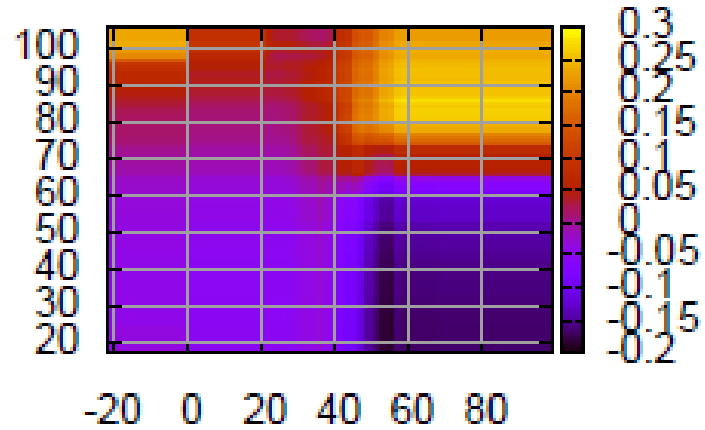
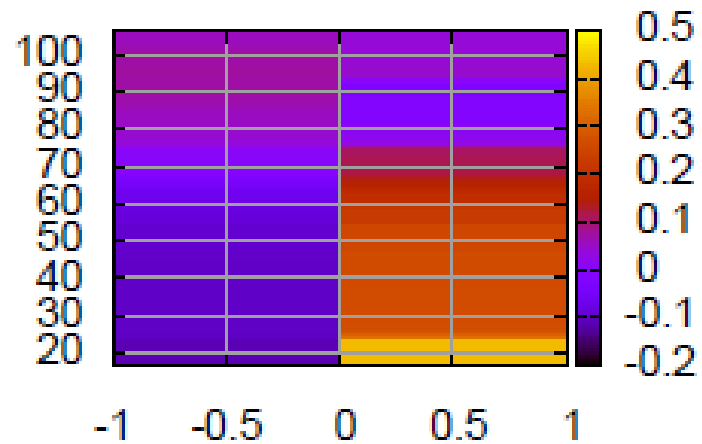(Blood Urea Nitrogen)

chronic lung disease     history of chest pain

Chronic lung disease and a history of chest pain both lower risk (similar problem as asthma)

# Interpretation



age vs. respiration rate

Old people with high respiration rate have the highest risk



age vs. cancer

- Risk is highest for the youngest patients
- It declines for patients who acquire cancer later in life
- For patients without cancer, risk rises as expected with age

# Takeaway

- If a model contains a modest number of terms (e.g., less than 50), it is best to show terms in the model to experts in the order they are most familiar with

- When the number of terms grows large, it is best to provide a well-defined ordering of the terms for a patient (from terms that increase risk most to terms that decrease risk most)

# Question?

# Reference

- Lou, Yin, Rich Caruana, and Johannes Gehrke. "Intelligible models for classification and regression." *Proceedings of the 18th ACM SIGKDD international conference on Knowledge discovery and data mining*. 2012.
- Lou, Yin, et al. "Accurate intelligible models with pairwise interactions." *Proceedings of the 19th ACM SIGKDD international conference on Knowledge discovery and data mining*. 2013.
- Caruana, Rich, et al. "Intelligible models for healthcare: Predicting pneumonia risk and hospital 30-day readmission." *Proceedings of the 21th ACM SIGKDD international conference on knowledge discovery and data mining*. 2015.