

Generalized Additive Models via Direct Optimization of Regularized Decision Stump Forests

Magzhan Gabidolla and Miguel Á. Carreira-Perpiñán, EECS, UC Merced

1 Introduction

- Generalized Additive Models (GAMs) are an important model class in interpretable machine learning.
- One effective approach in learning these models is through an ensemble of decision stumps.
- Since each stump depends on a single feature, stumps using the same feature x_d can be grouped to define a *shape function* $f_d(x_d)$ for this feature.
- Let $s(\mathbf{x}; \theta): \mathbb{R}^D \rightarrow \mathbb{R}$ be a decision stump with 4 learnable parameters: $\theta = \{\phi, \tau, \mu^l, \mu^r\}$. $\phi \in \{1, \dots, D\}$ is a feature index to split, $\tau \in \mathbb{R}$ is a threshold value, $\mu^l, \mu^r \in \mathbb{R}$ are the left and right leaf prediction values:

$$s(\mathbf{x}; \theta) = \begin{cases} \mu^l, & \text{if } x_\phi < \tau \\ \mu^r, & \text{if } x_\phi \geq \tau. \end{cases}$$

- A stump forest $F(\mathbf{x}; \Theta)$ is defined as a sum of T stumps plus a bias term μ :

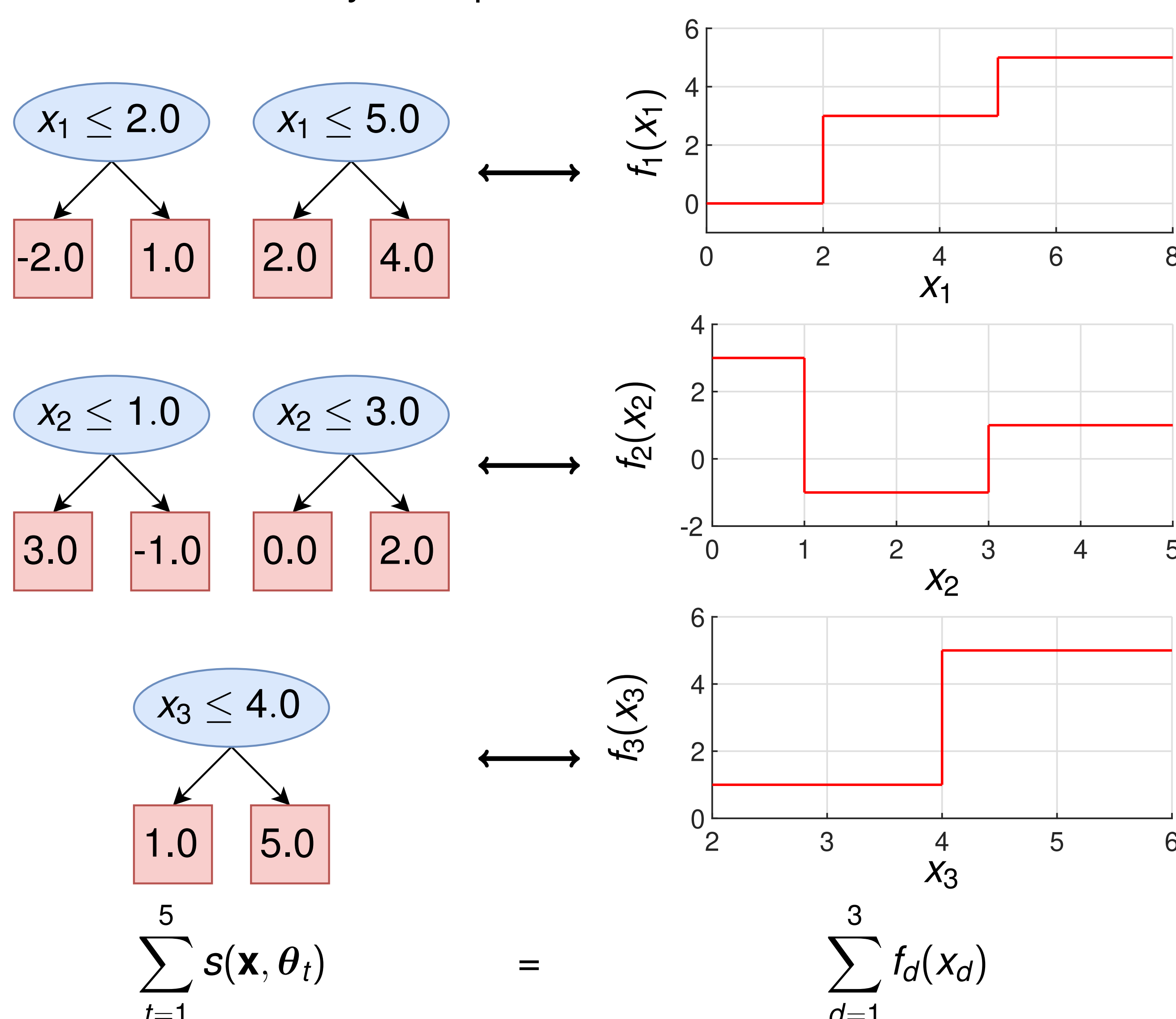
$$F(\mathbf{x}; \Theta) = \mu + \sum_{t=1}^T s(\mathbf{x}; \theta_t).$$

- Regrouping the stumps by feature indices, we obtain an additive model:

$$F(\mathbf{x}; \Theta) = \mu + \sum_{t=1}^T s(\mathbf{x}; \theta_t) = \mu + \sum_{d=1}^D \sum_{t: \phi_t=d} s(\mathbf{x}; \theta_t) = \mu + \sum_{d=1}^D f_d(x_d),$$

where $f_d(x_d) = \sum_{t: \phi_t=d} s(\mathbf{x}; \theta_t)$.

- Illustration on a toy example:



2 Alternating Optimization

- Consider a regression problem with a squared error loss:

$$\min_{\Theta} \frac{1}{2} \sum_{n=1}^N [y_n - F(\mathbf{x}_n; \Theta)]^2.$$

- We apply alternating optimization with the following steps:
- Individual stumps.** The optimization problem over a given stump $s(\cdot; \theta_t)$ when others are fixed is:

$$\min_{\theta_t} \frac{1}{2} \sum_{n=1}^N [y_n - \sum_{u \neq t} s(\mathbf{x}_n; \theta_u) - s(\mathbf{x}_n; \theta_t)]^2$$

This is a standard regression problem over a stump but with targets corresponding to the residuals. It can be solved exactly through enumeration over each (feature, threshold) pair as in traditional decision tree algorithms.

- All leaf parameters.** Once the splits $\{\phi_t, \tau_t\}_{t=1}^T$ are fixed, the problem simplifies to a linear regression on features corresponding to stump partitions. To see this, we can rewrite the predictive function of a stump using an indicator function: $s(\mathbf{x}; \theta) = \mu^l I(x_\phi < \tau) + \mu^r I(x_\phi \geq \tau)$. With this notation, the objective function over all leaves is:

$$\min_{\mu, \{\mu_t^l, \mu_t^r\}_{t=1}^T} \frac{1}{2} \sum_{n=1}^N \left[y_n - \sum_{t=1}^T \mu_t^l I(x_{\phi_t} < \tau_t) + \mu_t^r I(x_{\phi_t} \geq \tau_t) \right]^2.$$

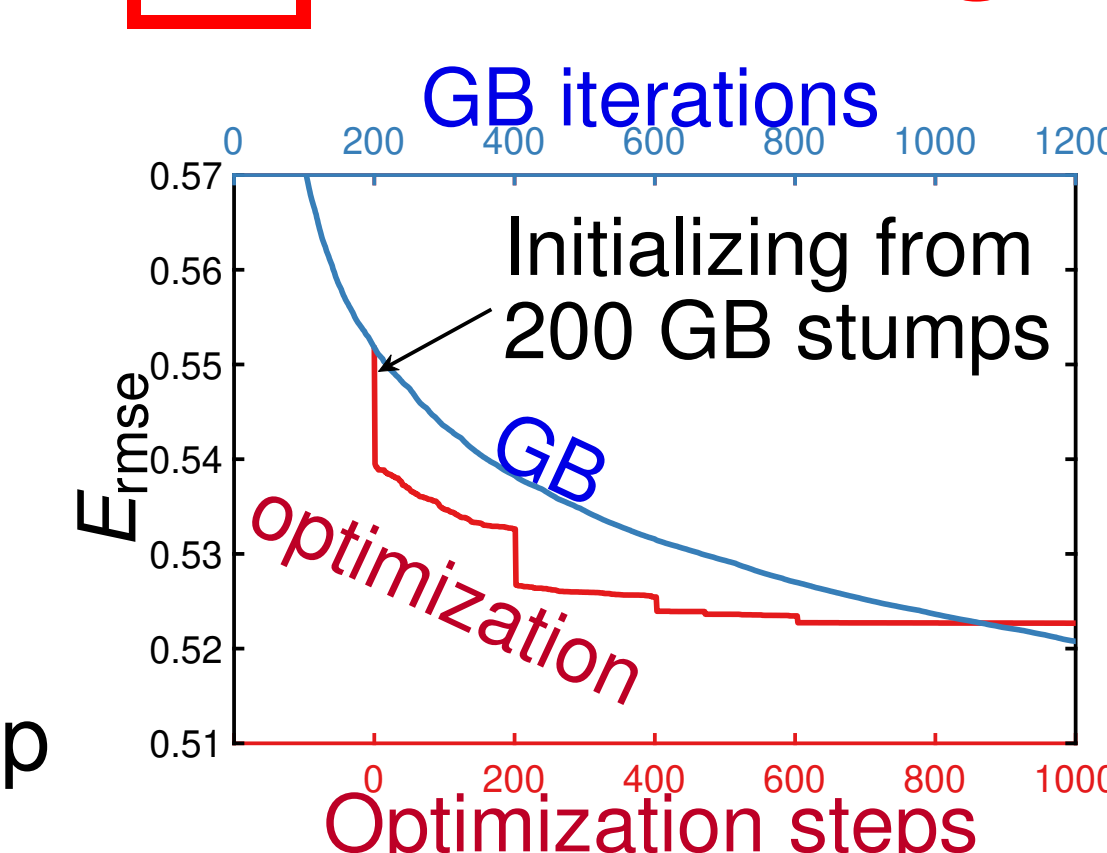
Since in this step all the splits $\{\phi_t, \tau_t\}_{t=1}^T$ are fixed, the indicator functions $I(\cdot)$ are just constants multiplying the unknowns $\{\mu_t^l, \mu_t^r\}_{t=1}^T$ that appear linearly inside the squared error. And so, eq. 1 is a simple linear regression problem on features induced by the stump splits, and can also be solved exactly and efficiently.

4 Experiments, comparison

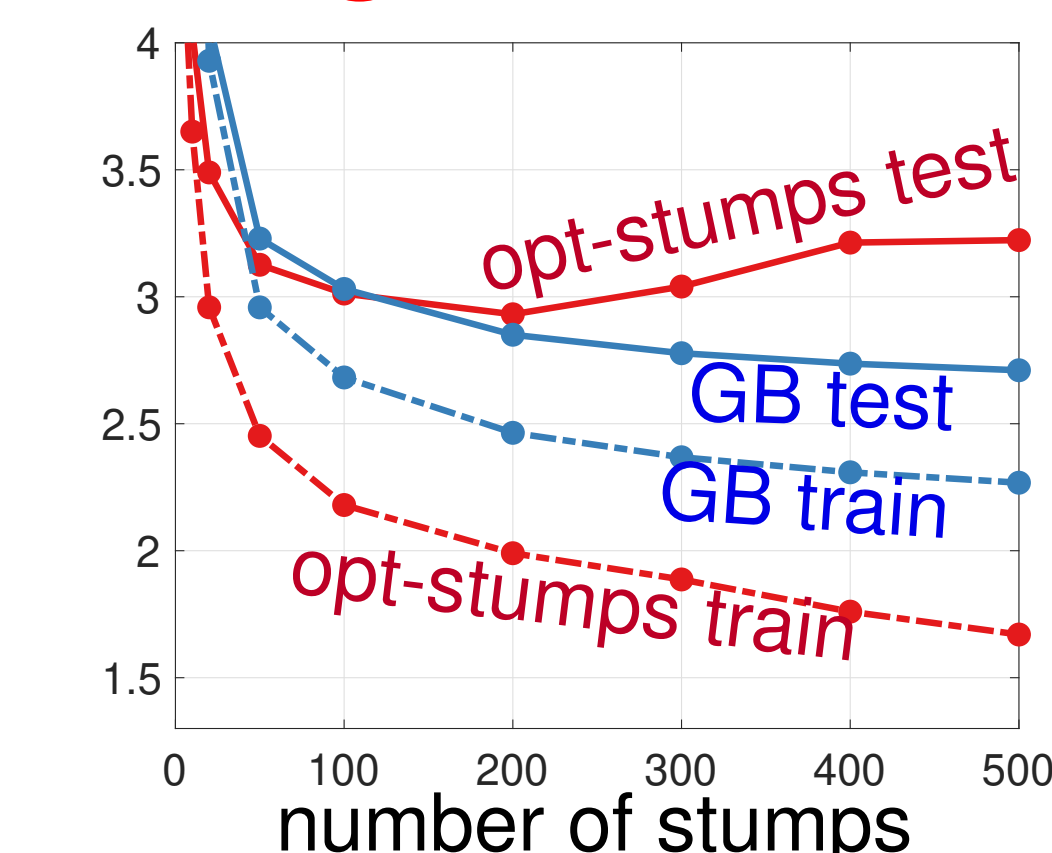
Table: RMSE error. Green color is the best test error, and blue is the second best.

Dataset	Ours	GB	EBM	Splines	FastSparse
Cpuact train $\times 10^{-2}$ N=8.2k test $\times 10^{-2}$ size D=21	2.12±0.01 2.37±0.03 642±0 9.4±0.3	2.20±0.04 2.43±0.06 3.4k±133 46±17	2.19±0.02 2.50±0.05 16.6k±36 39±2	2.53±0.02 2.69±0.06 271±3 37±0.03	2.76±0.03 2.91±0.17 119±4 3.8±0.5
Wine train $\times 10^{-2}$ N=6.5k test $\times 10^{-2}$ size D=11	65.70±0.15 70.02±0.66 724±12 6.0±0.3	68.13±0.27 70.92±0.51 770±32 2.87±0.58	66.73±0.27 70.12±0.39 3.9k±11 4.44±1.33	67.99±0.29 71.79±1.40 197±7 56±16	68.01±0.25 71.77±0.63 182±4 0.57±0.07
Housing train $\times 10^{-2}$ N=21k test $\times 10^{-2}$ size D=8	51.84±0.16 54.80±0.65 1.4k±20 13.6±0.4	54.24±0.27 56.15±0.58 2.4k±31 42±8	52.70±0.04 55.23±0.68 7.2k±8 36±2	53.37±0.21 55.49±0.61 528±2 37±2	54.62±0.20 56.29±0.65 579±9 3.94±0.73
Diamond train $\times 10^2$ N=54k test $\times 10^2$ size D=26	9.95±0.02 10.15±0.08 934±16 25.1±0.9	10.07±0.05 10.19±0.08 1182±81 140±58	10.11±0.03 10.23±0.06 3.4k±7 20±2	10.02±0.02 10.96±1.45 273±24 42±0.4	10.01±0.02 10.17±0.09 516±11 45±10

3 Overfitting and Regularization



Optimization ability. Comparison of FAO against GB for stump forests in train error for the California Housing dataset.



Overfitting problem of optimized stump forests on the cpuact dataset. Here GB avoids overfitting by using a learning rate 0.3.

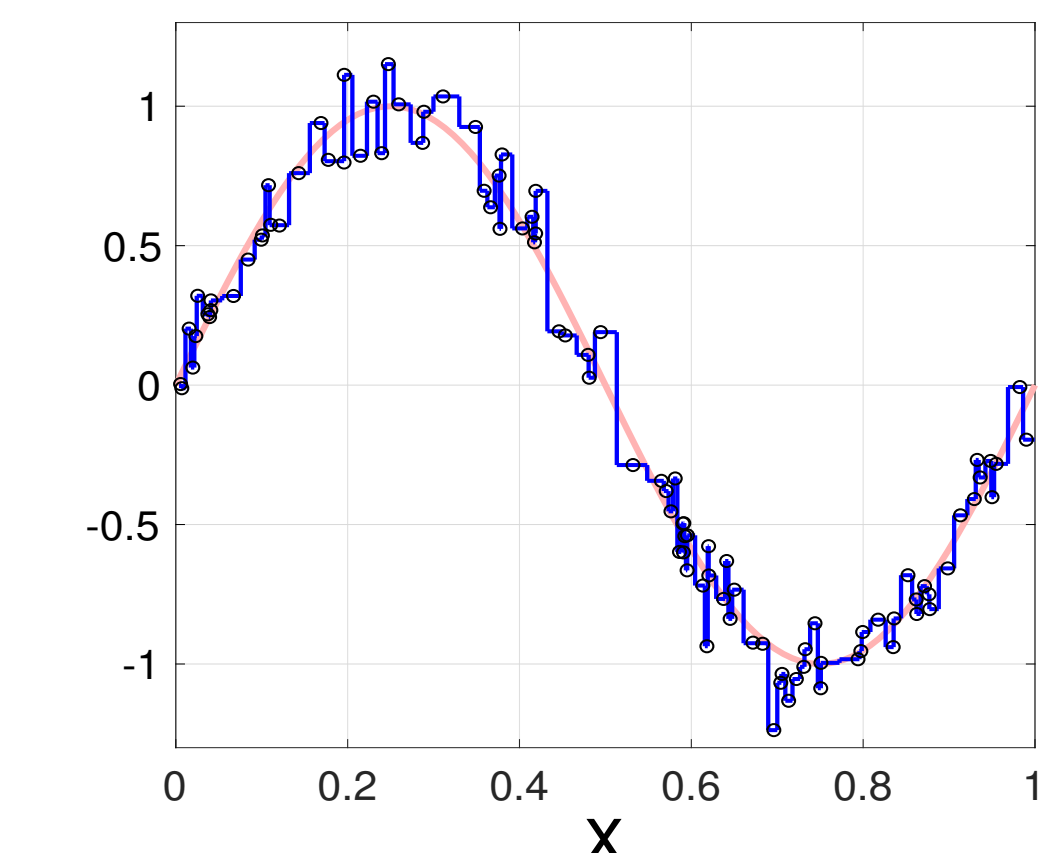


Illustration on a simple 1D regression problem with 100 stumps. Ground truth targets are from a sine function plotted in light red.

- We apply an ℓ_1 discontinuity penalty to shape functions:

$$r(\{\tau_t, \mu_t^l, \mu_t^r\}_{t: \phi_t=d}) = \sum_{t=0}^{T_d-1} |\beta_{t+1}^d - \beta_t^d|$$

where β_t is the value of the t -th constant piece from the left, and T_d is the number of constant pieces for feature d . This penalizes the difference between adjacent constant piece values.

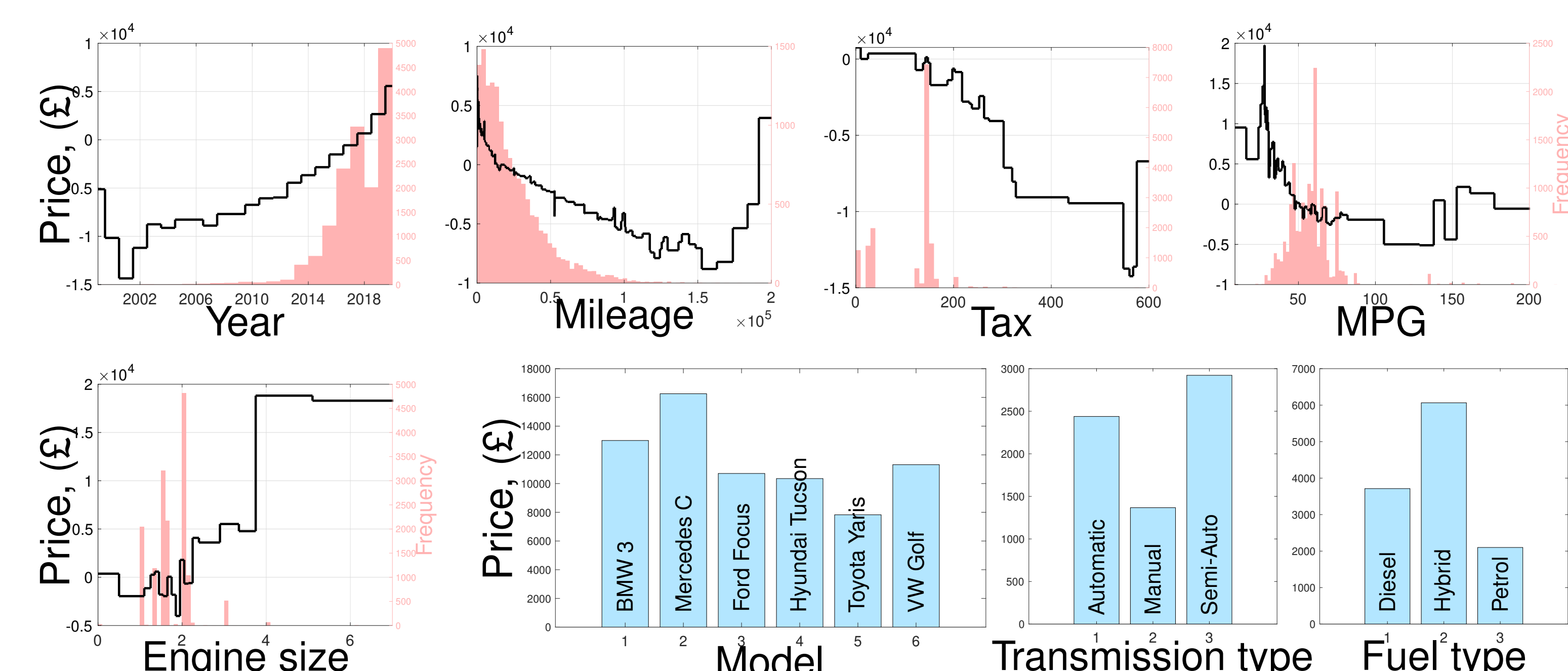
- We also propose is to penalize the deviation from the bias for each leaf value:

$$\sum_{t=1}^T (\mu_t^l - \mu)^2 + (\mu_t^r - \mu)^2.$$

- With these two types of regularization, our final objective function is:

$$\min_{\Theta} \frac{1}{2} \sum_{n=1}^N [y_n - F(\mathbf{x}_n; \Theta)]^2 + \lambda \sum_{d=1}^D r(\{\tau_t, \mu_t^l, \mu_t^r\}_{t: \phi_t=d}) + \alpha \sum_{t=1}^T ((\mu_t^l - \mu)^2 + (\mu_t^r - \mu)^2)$$

5 Experiments, interpretability



Visualization of the resulting GAM on the UK used car price prediction dataset. For the numerical features, the light red bars show the histogram of the training points with the frequency values given on the right y-axis.