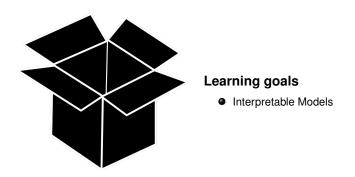
Interpretable Machine Learning Introduction and Background



HIGH-DIMENSIONAL MODEL REPRESENTATION

 A high-dimensional model representation (HDMR) decomposes the model into a sum of effect terms of increasing order:

$$\hat{f}(x) = g_{\{0\}} + g_{\{1\}}(x_1) + g_{\{2\}}(x_2) + \ldots + g_{\{1,2\}}(x_1, x_2) + \ldots + g_{\{1,\dots,\rho\}}(x_1,\dots,x_p)$$

- The features need to be independent to make the HDMR unique.
- Different techniques to estimate an additive decomposition exist, e.g., repeated expectations (partial dependence / PD) or accumulated local effects (ALE).

ADDITIVE DECOMPOSITION OF A PREDICTION FUNCTION

Consider the estimation via iterative expectations:

$$egin{aligned} g_{\{0\}} &= \mathbb{E}_{X} \left[\widehat{f}(x)
ight] \ g_{\{1\}}(x_{1}) &= \mathbb{E}_{X_{-1}} \left[\widehat{f}(x) \mid X_{1}
ight] - g_{\{0\}} \ g_{\{2\}}(x_{2}) &= \mathbb{E}_{X_{-2}} \left[\widehat{f}(x) \mid X_{2}
ight] - g_{\{0\}} \ g_{\{1,2\}}(x_{1},x_{2}) &= \mathbb{E}_{X_{-\{1,2\}}} \left[\widehat{f}(x) \mid X_{1},X_{2}
ight] - g_{\{2\}}(x_{2}) - g_{\{1\}}(x_{1}) - g_{\{0\}} \ &\vdots \ g_{\{1,\dots,
ho\}}(x) &= \widehat{f}(x) - \dots - g_{\{1,2\}}(x_{1},x_{2}) \ &- g_{\{2\}}(x_{2}) - g_{\{1\}}(x_{1}) - g_{\{0\}} \end{aligned}$$

FUNCTIONAL ANOVA

After \hat{f} has been decomposed, we can conduct a functional analysis of variance (functional ANOVA / FANOVA):

0

$$Var\left[\hat{f}(x)\right] = Var\left[g_{\{0\}} + g_{\{1\}}(x_1) + g_{\{2\}}(x_2) + \ldots + g_{\{1,2\}}(x_1, x_2) + \ldots + g_{\{1,\dots,\rho\}}(x)\right]$$

 If the features are independent, the variance can be additively decomposed without covariances:

$$Var\left[\hat{f}(x)\right] = Var\left[g_{\{0\}}\right] + Var\left[g_{\{1\}}(x_1)\right] + Var\left[g_{\{2\}}(x_2)\right] + Var\left[g_{\{1,2\}}(x_1, x_2)\right] + \dots + Var\left[g_{\{1,...,p\}}(x)\right]$$

FUNCTIONAL ANOVA

 Dividing by the prediction variance results in the fraction of variance explained by each term:

$$1 = \frac{Var\left[g_{\{0\}}\right]}{Var\left[\hat{f}(x)\right]} + \frac{Var\left[g_{\{1\}}(x_1)\right]}{Var\left[\hat{f}(x)\right]} + \frac{Var\left[g_{\{2\}}(x_2)\right]}{Var\left[\hat{f}(x)\right]} + \frac{Var\left[g_{\{1,2\}}(x_1, x_2)\right]}{Var\left[\hat{f}(x)\right]} + \dots + \frac{Var\left[g_{\{1,\dots,p\}}(x)\right]}{Var\left[\hat{f}(x)\right]}$$

 The fraction of variance explained by a term is referred to as the Sobol index:

$$S_j = rac{Var\left[g_{\{j\}}(x_j)
ight]}{Var\left[\hat{f}(x)
ight]}$$