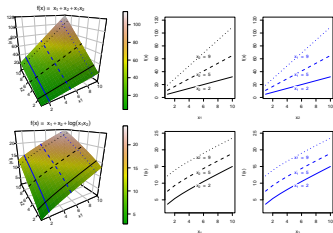
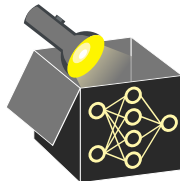


Interpretable Machine Learning

Feature Interactions

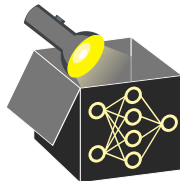


Learning goals

- Feature interactions
- Difference to feature dependencies

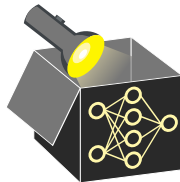
FEATURE INTERACTIONS

- Feature dependencies concern data distribution
- Feature interactions may occur in structure of **both** model or DGP (e.g., functional relationship between X and (X) or X and $Y = f(X)$)
 - ~> Feature dependencies may lead to feature interactions in a model



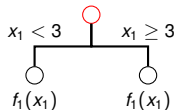
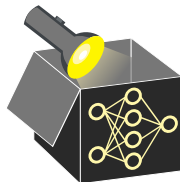
FEATURE INTERACTIONS

- Feature dependencies concern data distribution
- Feature interactions may occur in structure of **both** model or DGP (e.g., functional relationship between X and (X) or X and $Y = f(X)$)
 - ~> Feature dependencies may lead to feature interactions in a model
- No. of potential interactions increases exponentially with no. of features
 - ~> Difficult to identify interactions, especially when features are dep.

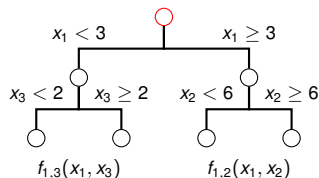


FEATURE INTERACTIONS

- Feature dependencies concern data distribution
- Feature interactions may occur in structure of **both** model or DGP (e.g., functional relationship between X and (X) or X and $Y = f(X)$)
 - ~> Feature dependencies may lead to feature interactions in a model
- No. of potential interactions increases exponentially with no. of features
 - ~> Difficult to identify interactions, especially when features are dep.
- Interactions: Feature's effect on the prediction depends on other features
 - ~> Example: $() = x_1 x_2 \Rightarrow$ Effect of x_1 on depends on x_2 and vice versa



No interaction



Interactions: x_1 and x_3 ,
 x_1 and x_2

No interactions: x_2 and x_3

FEATURE INTERACTIONS

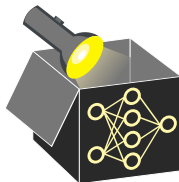
► FRIEDMAN_POPESCU

Definition: A function $f()$ contains an interaction between x_j and x_k if a difference in $f()$ -values due to changes in x_j will also depend on x_k , i.e.:

$$\mathbb{E} \left[\frac{\partial^2 f()}{\partial x_j \partial x_k} \right]^2 > 0$$

\Rightarrow If x_j and x_k don't interact, $f()$ is sum of 2 functions, each indep. of x_j, x_k :

$$f() = f_{-j}(x_1, \dots, x_{j-1}, x_{j+1}, \dots, x_p) + f_{-k}(x_1, \dots, x_{k-1}, x_{k+1}, \dots, x_p)$$

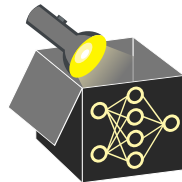


FEATURE INTERACTIONS

Example: $f() = x_1 + x_2 + x_1 \cdot x_2$ (not separable)

$$\mathbb{E} \left[\frac{\partial^2 (x_1 + x_2 + x_1 \cdot x_2)}{\partial x_1 \partial x_2} \right]^2 = \mathbb{E} \left[\frac{\partial (1 + x_2)}{\partial x_2} \right]^2 = 1 > 0$$

\Rightarrow interaction between x_1 and x_2

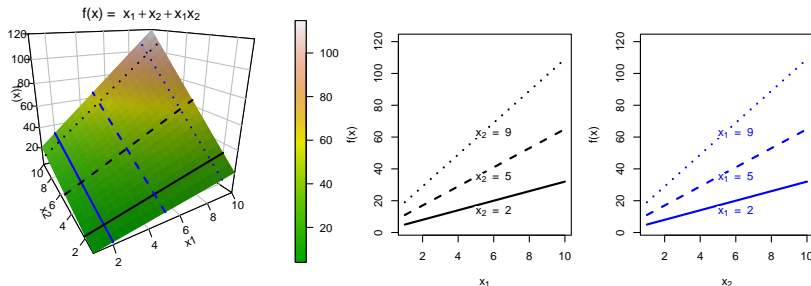
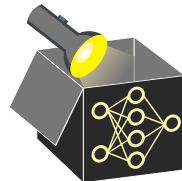


FEATURE INTERACTIONS

Example: $f() = x_1 + x_2 + x_1 \cdot x_2$ (not separable)

$$\mathbb{E} \left[\frac{\partial^2 (x_1 + x_2 + x_1 \cdot x_2)}{\partial x_1 \partial x_2} \right]^2 = \mathbb{E} \left[\frac{\partial (1 + x_2)}{\partial x_2} \right]^2 = 1 > 0$$

⇒ interaction between x_1 and x_2



- Effect of x_1 on $f()$ varies with x_2 (and vice versa)

⇒ Different slopes

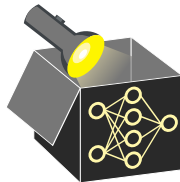
FEATURE INTERACTIONS

Example of separable function:

$$f() = x_1 + x_2 + \log(x_1 \cdot x_2) = x_1 + x_2 + \log(x_1) + \log(x_2)$$

$$\Rightarrow f() = f_1(x_1) + f_2(x_2) \text{ with } f_1(x_1) = x_1 + \log(x_1) \text{ and } f_2(x_2) = x_2 + \log(x_2)$$

$$\Rightarrow \text{no interactions due to separability, also } \mathbb{E} \left[\frac{\partial^2 f()}{\partial x_1 \partial x_2} \right]^2 = 0$$



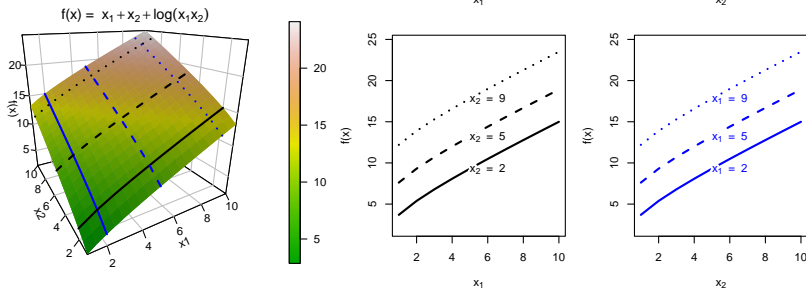
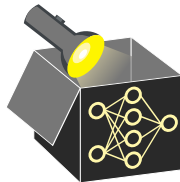
FEATURE INTERACTIONS

Example of separable function:

$$f() = x_1 + x_2 + \log(x_1 \cdot x_2) = x_1 + x_2 + \log(x_1) + \log(x_2)$$

$$\Rightarrow f() = f_1(x_1) + f_2(x_2) \text{ with } f_1(x_1) = x_1 + \log(x_1) \text{ and } f_2(x_2) = x_2 + \log(x_2)$$

$$\Rightarrow \text{no interactions due to separability, also } \mathbb{E} \left[\frac{\partial^2 f()}{\partial x_1 \partial x_2} \right]^2 = 0$$



- Effect of x_1 on $f()$ stays the same for different x_2 values (and vice versa)

\Rightarrow Parallel lines at different horizontal (blue) or vertical (black) slices