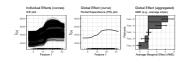
Interpretable Machine Learning

Individual Conditional Expectation (ICE) Plot





- ICE curves as local effect method
- How to sample grid points for ICE curves



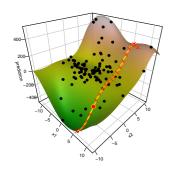
MOTIVATION

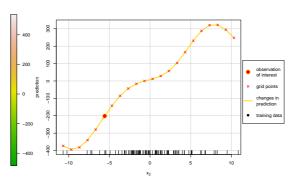
Question: How does varying a single feature of an obs. affect its predicted outcome?

Idea: For a given observation, change the value of the feature of interest, and visualize how prediction changes

Example: On model prediction surface (left), select observation and visualize changes in prediction for different values of x_2 , while keeping x_1 fixed

 \Rightarrow local interpretation







INDIVIDUAL CONDITIONAL EXPECTATION (ICE)

► Goldstein et. al (2013)

Partition each observation ${\bf x}$ into ${\bf x}_S$ (feature(s) of interest) and ${\bf x}_{-S}$ (remaining features)

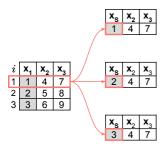


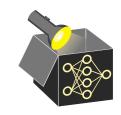
	x _s	x	-s
i	X ₁	X ₂	X ₃
1	1	4	7
2	2	5	8
3	3	6	9

 \leadsto In practice, \mathbf{x}_S consists of one or two features (i.e., $|S| \leq 2$ and $-S = S^{\complement}$).

Formal definition of ICE curves:

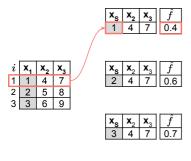
- Define grid points $\mathbf{x}_{S}^{*} = \mathbf{x}_{S}^{*(1)}, \dots, \mathbf{x}_{S}^{*(g)}$ to vary \mathbf{x}_{S}
- $$\begin{split} \bullet & \text{ Plot point pairs } \left\{ \left(\mathbf{x}_{S}^{*^{(k)}}, \hat{f}_{S,ICE}^{(i)}(\mathbf{x}_{S}^{*^{(k)}})\right) \right\}_{k=1}^{g} \\ & \text{ where } \hat{f}_{S,ICE}^{(i)}(\mathbf{x}_{S}^{*}) = \hat{f}(\mathbf{x}_{S}^{*}, \mathbf{x}_{-S}^{(i)}) \end{split}$$
- For each k connect point pairs to obtain ICE curve
- ightharpoonup ICE curves visualize how prediction of *i*-th observation changes after varying its feature values indexed by S using grid points \mathbf{x}_S^* while keeping all values in -S fixed

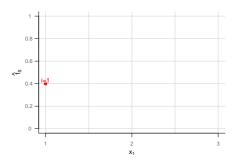




1. Step - Grid points:

- ullet Sample grid values $\mathbf{x}_{S}^{*^{(1)}},\ldots,\mathbf{x}_{S}^{*^{(g)}}$ along possible values of feature $S\left(|S|=1\right)$
- For $\mathbf{x}^{(i)} = (\mathbf{x}_{S}, \mathbf{x}_{-S})$, replace \mathbf{x}_{S} with those grid values
- \Rightarrow Creates new artificial points for *i*-th observation (here: $\mathbf{x}_S^* = x_1^* \in \{1, 2, 3\}$ scalar)



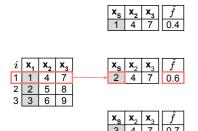


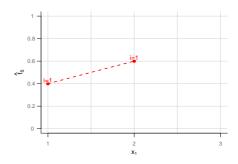


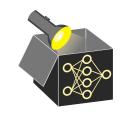
2. Step - Predict and visualize:

For each artificially created data point of *i*-th observation, plot prediction $\hat{f}_{S,ICE}^{(i)}(\mathbf{x}_S^*)$ vs. grid values \mathbf{x}_S^* :

$$\hat{f}_{1,ICE}^{(i)}(x_1^*) = \hat{f}(x_1^*, \mathbf{x}_{2,3}^{(i)}) ext{ vs. } x_1^* \in \{1, 2, 3\}$$



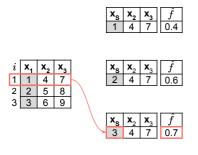


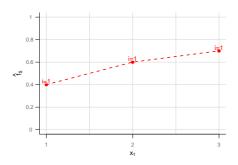


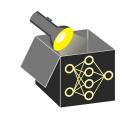
2. Step - Predict and visualize:

For each artificially created data point of *i*-th observation, plot prediction $\hat{f}_{S,ICE}^{(i)}(\mathbf{x}_S^*)$ vs. grid values \mathbf{x}_S^* :

$$\hat{f}_{1,ICE}^{(i)}(x_1^*) = \hat{f}(x_1^*, \mathbf{x}_{2,3}^{(i)}) ext{ vs. } x_1^* \in \{1, 2, 3\}$$



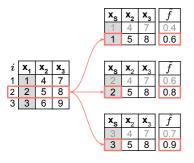


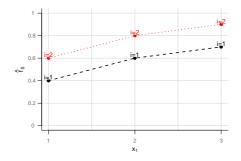


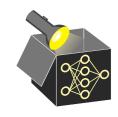
2. Step - Predict and visualize:

For each artificially created data point of *i*-th observation, plot prediction $\hat{f}_{S,ICE}^{(i)}(\mathbf{x}_S^*)$ vs. grid values \mathbf{x}_S^* :

$$\hat{f}_{1,ICE}^{(i)}(x_1^*) = \hat{f}(x_1^*, \mathbf{x}_{2,3}^{(i)}) ext{ vs. } x_1^* \in \{1, 2, 3\}$$

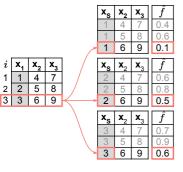


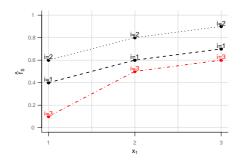




3. Step - Repeat for other observations:

ICE curve for i = 2 connects all predictions at grid values associated to i-th obs.







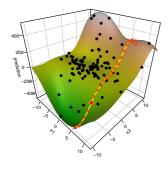
3. Step - Repeat for other observations:

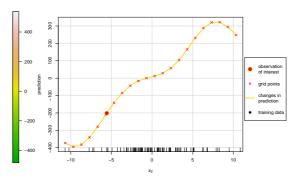
ICE curve for i = 3 connects all predictions at grid values associated to i-th obs.

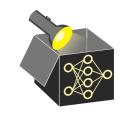
ICE CURVES - INTERPRETATION

Example: Prediction surface of a model (left), select observation and visualize changes in prediction for different values of x_2 while keeping x_1 fixed

 $\Rightarrow \text{local interpretation}$

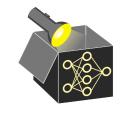




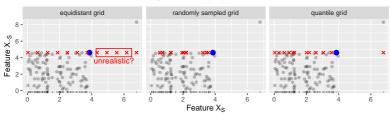


COMMENTS ON GRID VALUES

- Plotting ICE curves involves generating grid values x_S*; visualized on x-axis
- Three common strategies for grid definition:
 - Equidistant grid values within feature range
 - Random samples from observed feature values
 - Quantiles of observed feature values
- Marginal realism: Random and quantile grids better reflect the marginal distribution of $x_S \Rightarrow$ reduce unrealistic values along x_S

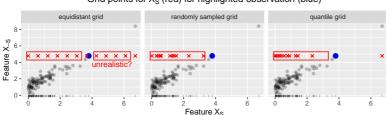


Grid points for X_S (red) for highlighted observation (blue)



COMMENTS ON GRID VALUES

- Plotting ICE curves involves generating grid values x_S^{*}; visualized on x-axis
- Three common strategies for grid definition:
 - Equidistant grid values within feature range
 - Random samples from observed feature values
 - Quantiles of observed feature values
- Marginal realism: Random and quantile grids better reflect the marginal distribution of $x_S \Rightarrow$ reduce unrealistic values along x_S
- However: For correlated features, extrapolation remains:

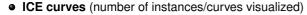






PRACTICAL CONSIDERATIONS

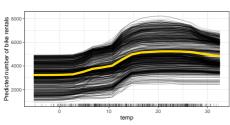
- Grid resolution (instances × grid over feature of interest)
 - Too coarse ⇒ may miss sharp nonlinearities or discontinuities
 - Too fine ⇒ high runtime (without gaining much)
 - Fix: cap at $\approx 50-100$ grid points; vectorize predictions by feeding the model a single data frame containing all grid-modified instances



- Too few ⇒ hides variability across instances, misses subgroup differences
- \bullet Too many \Rightarrow visual overload (many overlapping curves), time intensive
- Fix: Stratified or cluster-based subsample (e.g., 100); facet by subgroup

Default values for popular libraries:

Library	Grid	ICE curves
sklearn (Py)	100	1 000 (random
PDPbox (Py)	10	num. rows
iml (R)	20	num. rows
pdp (R)	51	num. rows



ICE curves (black lines) and their point-wise average across the grid (yellow line)

