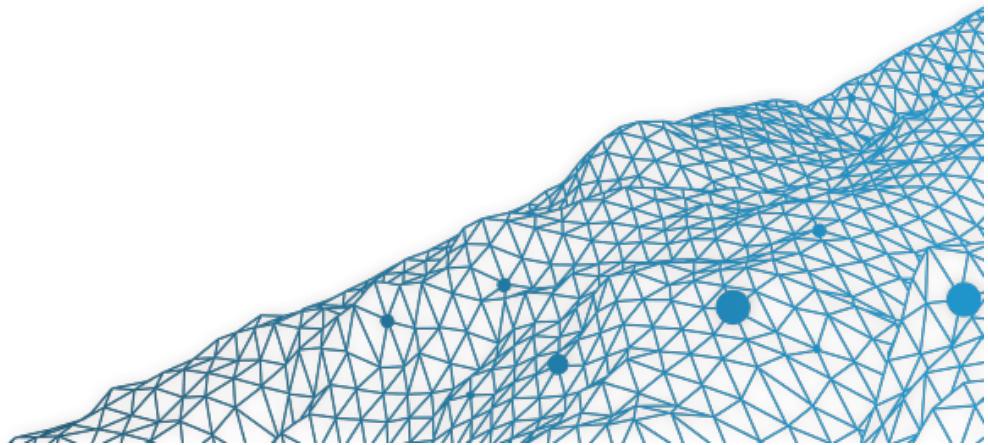


12 Interpretability & Explainability

Arthur Charpentier (Université du Québec à Montréal)

Machine Learning & Econometrics

SIDE Summer School - July 2019



Interpretability of Machine Learning

Consider a regression model

$$y_i = \beta_0 + \beta_1 x_{1,i} + \beta_2 x_{2,i} + \varepsilon_i$$

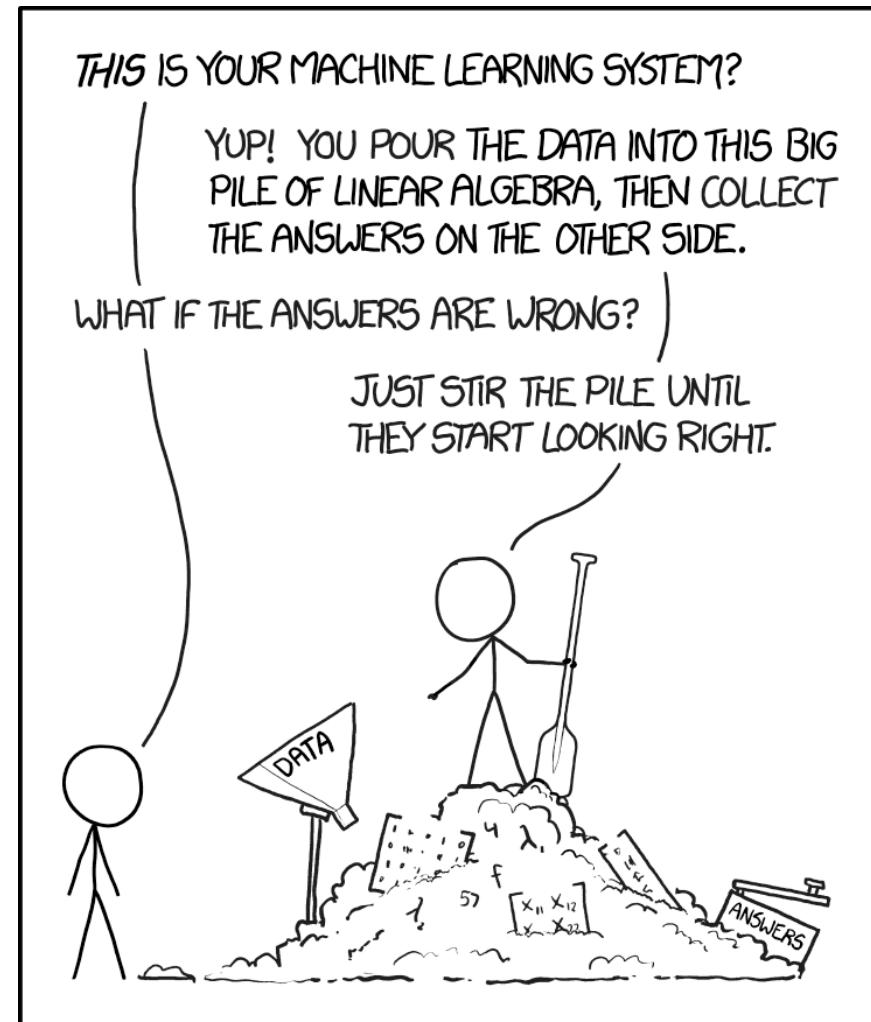
With Woodridge (2009, [Introductory Econometrics](#)) notations, holding x_2 fixed (*Ceteris Paribus* interpretation),

$$\beta_1 = \frac{\partial y}{\partial x_1}$$

assuming $\mathbb{E}[\varepsilon|x_1, x_2] = 0$, or with notions used so far

$$\beta_1 = \frac{\partial m(\mathbf{x})}{\partial x_1} (= \text{constant})$$

(source [Randall Munroe \(xkcd, 2016\)](#))



Linear Models

Ceteris paribus can be translated into “all other things being equal” or “holding other factors constant.”

In a loglinear regression, $\log y_i = \beta_0 + \beta_1 x_{1,i} + \beta_2 x_{2,i} + \varepsilon_i$

$$\beta_1 = \frac{\partial \Delta y}{y \partial x_1}$$

In a logistic linear regression, $\text{logit}(p_i) = \beta_0 + \beta_1 x_{1,i} + \beta_2 x_{2,i}$

$$\exp[\beta_1] = \frac{\partial y}{y \partial x_1}$$

Mutatis mutandis approximately translates as “allowing other things to change accordingly” or “the necessary changes having been made.”

Interpretability of Machine Learning

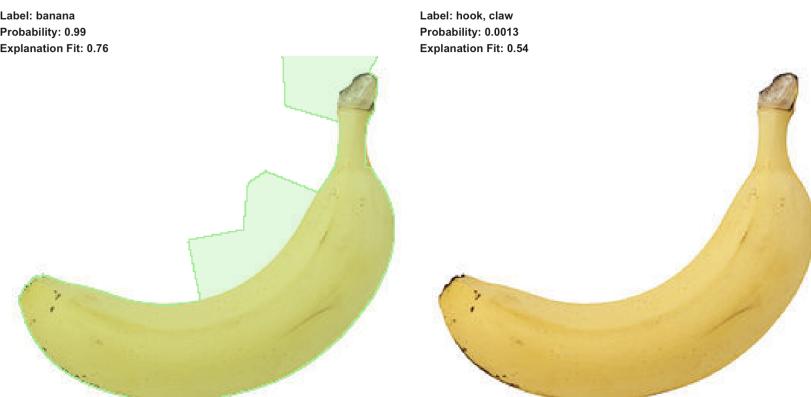
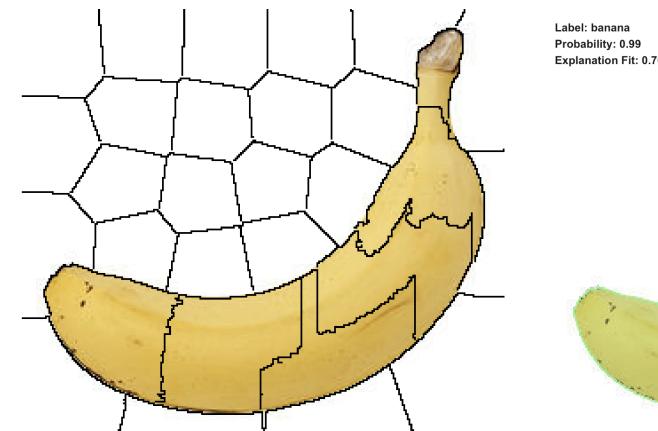
What do we mean by interpreting a machine learning model, and why do we need it? Is it to trust the model? Or try to find causal relationships in the analyzed phenomenon? Or to visualize it? see Lipton (2017, [Mythos of Model Interpretability](#)), Lakkaraju *et al.* (2019, [Faithful and Customizable Explanations of Black Box Models](#)) Molnar (2019. [Interpretable Machine Learning: A Guide for Making Black Box Models Explainable](#)), Guidotti *et al.* (2018, [A Survey of Methods for Explaining Black Box Models](#)) Gilpin *et al.* (2019, [Explaining Explanations: An Overview of Interpretability of Machine Learning](#)) Lundberg & Lee (2017, [A Unified Approach to Interpreting Model Predictions](#)), etc.

Interpretability of Predictive Models

E.g. classifier on labeled pictures $\{(\mathbf{x}_i, y_i)\}$



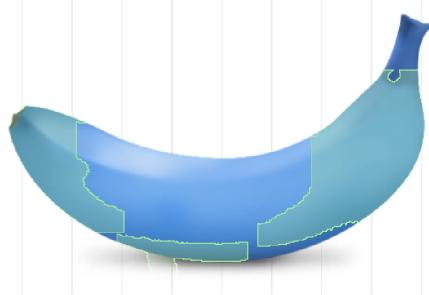
Given a new picture (\mathbf{x}) we want a label $\hat{y} = \text{label}(\mathbf{x})$



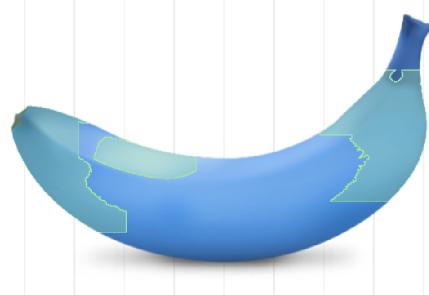
Interpretability of Predictive Models

There were neither yellow zucchini nor blue banana in the training dataset

Label: hook, claw
Probability: 0.56
Explanation Fit: 0.29



Label: nipple
Probability: 0.084
Explanation Fit: 0.41



Label: screwdriver
Probability: 0.53
Explanation Fit: 0.68

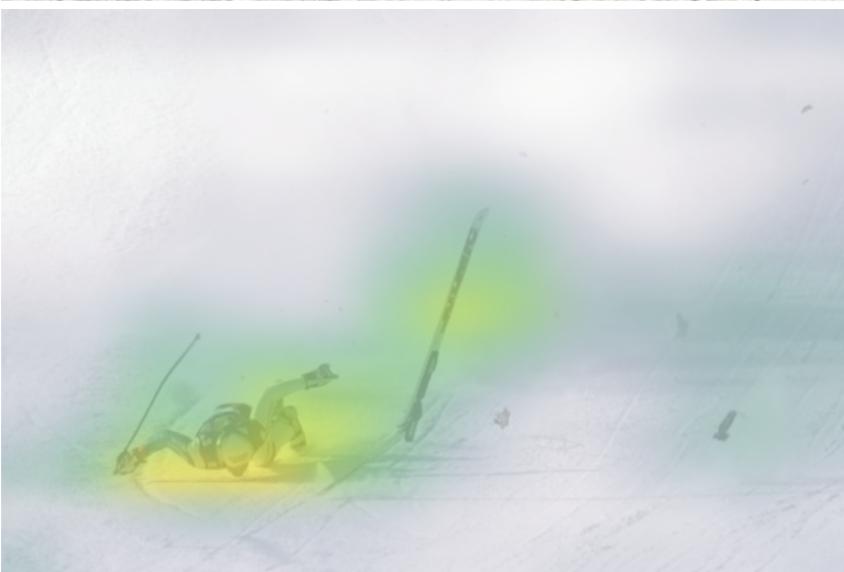


Label: zucchini, courgette
Probability: 0.089
Explanation Fit: 0.38



why did $\text{label}(x)$ recognize a claw and a screwdriver (56% and 53% chances)

Interpretability of Predictive Models \hat{y} = “ski” and “car”



Interpretability of Predictive Models

For works that describe machine learning models as black boxes, transparency and interpretability are closely related, if not the same concept.

We can open the black box either

- by explaining the model,
- by explaining the outcome
- by inspecting the black box internally
- by providing a transparent solution.

Interpretability of Predictive Models

Neural nets and random forests are considered as black boxes, Ribeiro *et al.* (2016, “Why Should I Trust You?”: Explaining the Predictions of Any Classifier).

Following Guidotti *et al.* (2018, A Survey of Methods for Explaining Black Box Models) In this general setting, a (black box) model is $m : \mathcal{X} \mapsto \mathcal{Y}$ (neural nets, SVM, etc), explanators ϵ will be described after (Features Importance, Sensitivity Analysis, Partial Dependence Plot, etc).

- Black-box model explanation

Given a black box predictor m and a dataset $\mathcal{D}_n = \{(\mathbf{x}_i, y_i)\} \in (\mathcal{X} \times \mathcal{Y})^n$, the black box model explanation problem consists in finding a function $f : (\mathcal{X} \mapsto \mathcal{Y}) \times (\mathcal{X} \times \mathcal{Y})^n \rightarrow (\mathcal{X} \mapsto \mathcal{Y})$ that returns a comprehensible global predictor c_g , i.e., $f(m, \mathcal{D}_n) = c_g$, such that c_g is able to mimic the behavior of m , and exists a global explanator function $\epsilon_g : (\mathcal{X} \mapsto \mathcal{Y}) \rightarrow \mathcal{E}$ that can derive from c_g a set of explanations $E \in \mathcal{E}$ modeling in a human understandable way the logic behind c_g , i.e., $\epsilon(c_g) = E$.

Interpretability of Predictive Models

- Black-box outcome explanation

Given a black box predictor b and a dataset $\mathcal{D}_n = \{(\mathbf{x}_i, y_i)\}$, the black box outcome explanation problem consists in finding a function

$f : (\mathcal{X} \mapsto \mathcal{Y}) \times (\mathcal{X} \times \mathcal{Y})^n \rightarrow (X \mapsto Y)$ that returns a comprehensible local predictor c_ℓ , i.e., $f(m, \mathcal{D}_n) = c_\ell$, such that c_ℓ is able to mimic the behavior of m , and exists a local explanator function $\epsilon_\ell : (\mathcal{X} \mapsto \mathcal{Y}) \times (\mathcal{X} \mapsto \mathcal{Y}) \times \mathcal{X} \rightarrow \mathcal{E}$ that can derive from the black box model m , the comprehensible local predictor c_ℓ , and a data record \mathbf{x} , a human understandable explanation $e \in E$ for the data record \mathbf{x} , i.e., $\epsilon_\ell(m, c_\ell, \mathbf{x}) = e$

Interpretability of Predictive Models

- Black box inspection explanation

Given a black box predictor b and a dataset $\mathcal{D}_n = \{(x_i, y_i)\}$, the black box inspection problem consists in finding a function $f : (\mathcal{X} \mapsto \mathcal{Y}) \times (\mathcal{X} \times \mathcal{Y})^n \rightarrow V$ that returns a visual representation of the behavior of the black box, $f(m, \mathcal{D}_n) = v$ with V being the set of all possible representations.

Interpretability of Predictive Models

The transparent box design problem consists in providing a model which is locally or globally interpretable on its own,

- Transparent box design problem

Given a dataset $\mathcal{D}_n = \{(\mathbf{x}_i, y_i)\}$, the transparent box design problem consists in finding a learning function $L : (\mathcal{X} \times \mathcal{Y})^n \rightarrow (\mathcal{X} \mapsto \mathcal{Y})$ that returns a (locally or globally) comprehensible predictor c , i.e., $L(\mathcal{D}_n) = c$. This implies that there exists an explanator function, local ϵ_ℓ or global ϵ_g , that takes as input the comprehensible predictor c and returns a human understandable explanation $e \in E$, or a set of explanations E .

Interpretability of Predictive Models

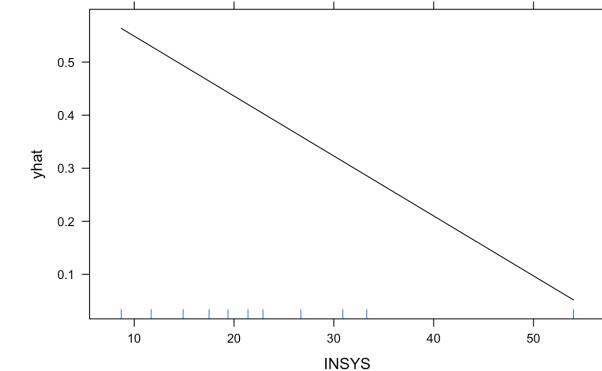
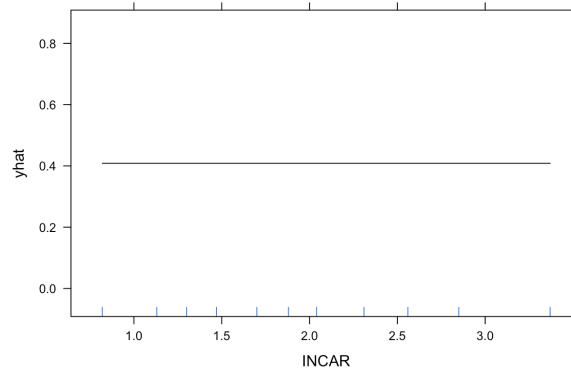
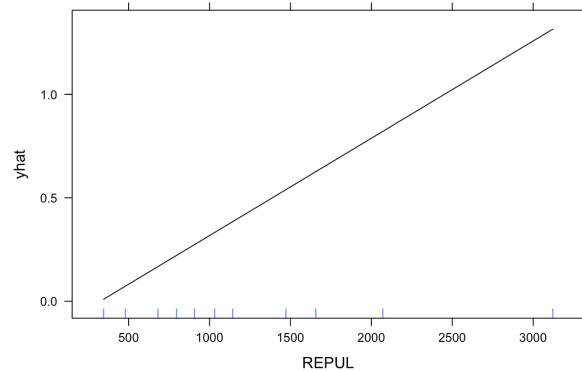
- Partial Dependence Plot

Introduced in Friedman (2001, [Greedy function approximation: A gradient boosting machine](#)). Let \mathbf{x} be split in two parts : \mathbf{x}_s (variable(s) of interest) and \mathbf{x}_c the complementary, $\mathbf{x} = (\mathbf{x}_s, \mathbf{x}_c)$. Partial dependence of \mathbf{x}_s is

$$p(\mathbf{x}_s) = \mathbb{E}[m(\mathbf{x}_s, \mathbf{x}_c)] \text{ and } \hat{p}(\mathbf{x}_s) = \frac{1}{n} \sum_{i=1}^n \hat{m}(\mathbf{x}_s, \mathbf{x}_{i,c})$$

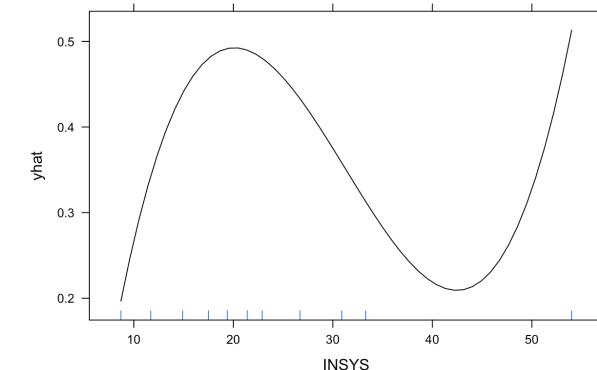
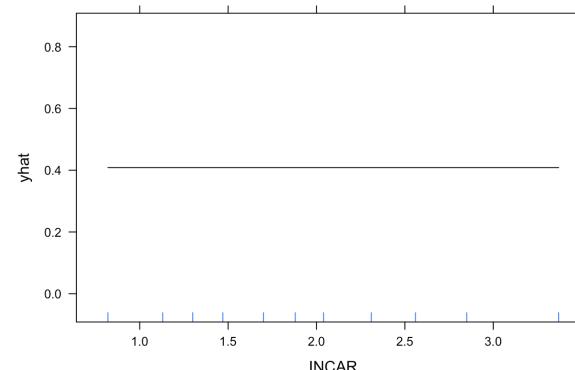
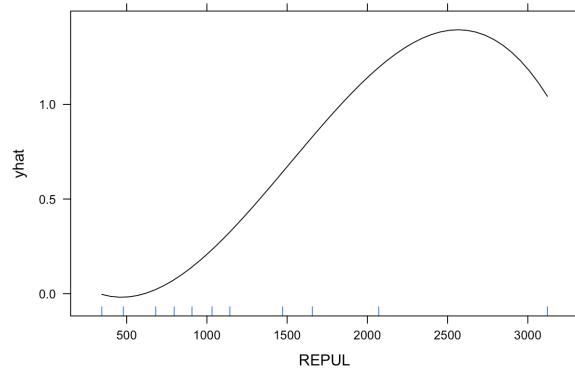
See `pdp` R package and `pdp::partial(model, pred.var = "REPUL", plot = TRUE)`

Consider a (standard) linear model on the `myocarde` dataset

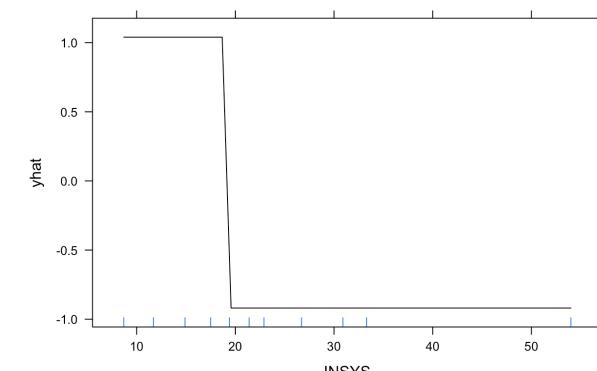
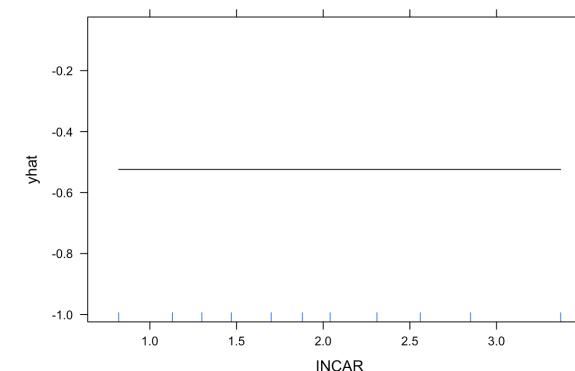
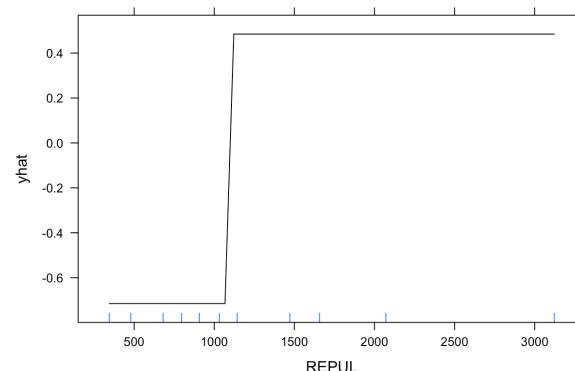


Interpretability of Predictive Models

Consider an additive model (GAM) on the `myocarde` dataset

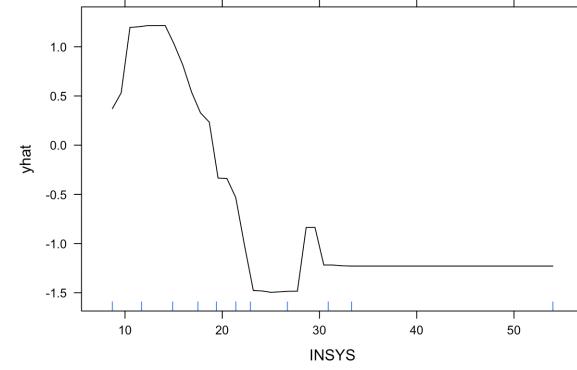
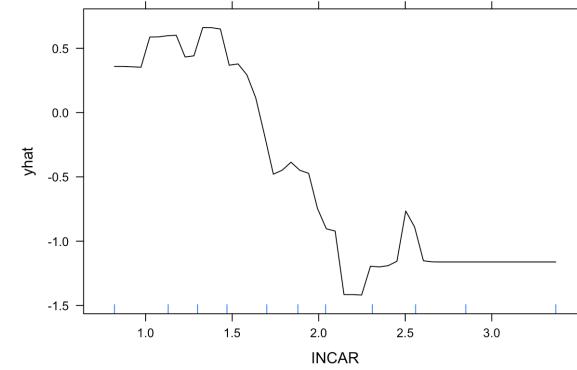
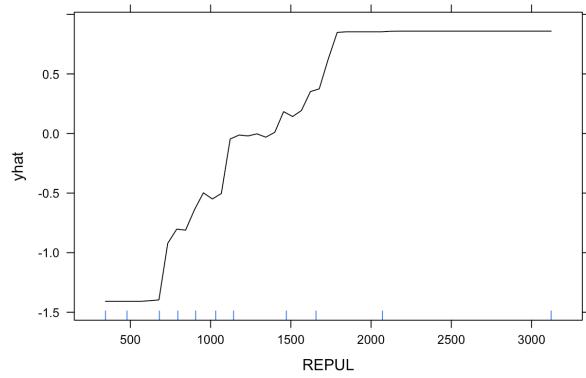


or a classification tree

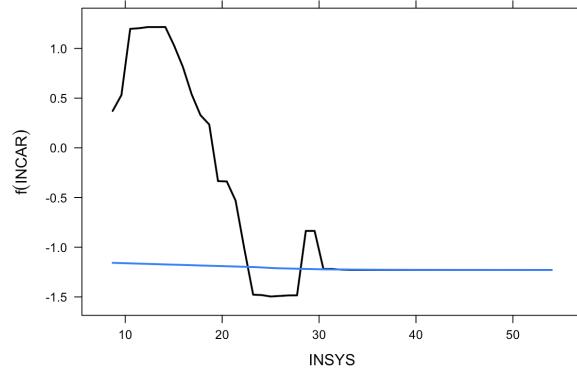
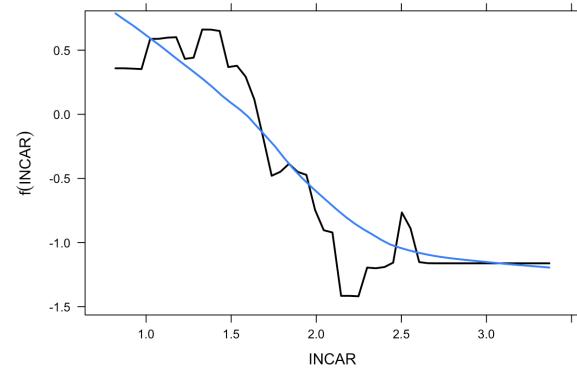
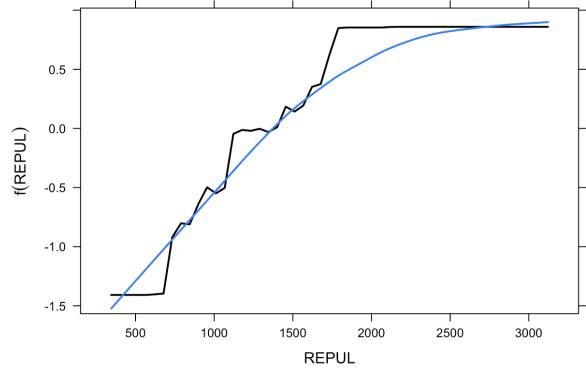


Interpretability of Predictive Models

Consider a random forest model on the `myocarde` dataset



or the smooth version with `plotPartial(smooth = TRUE, ...)`

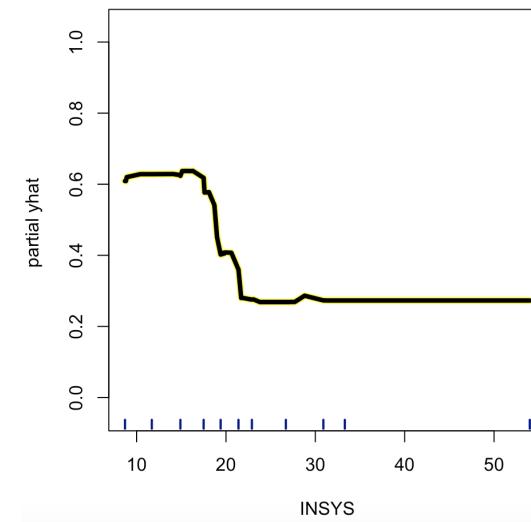
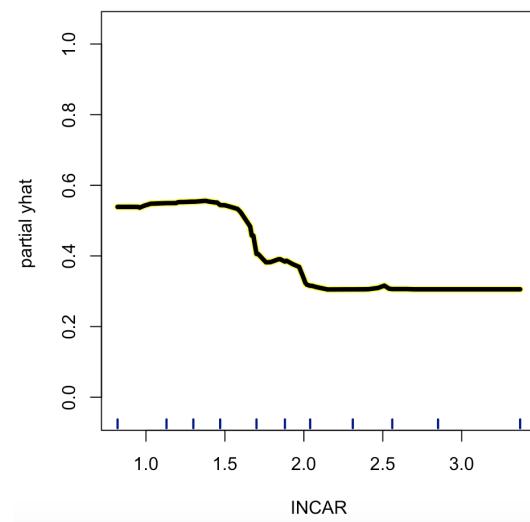
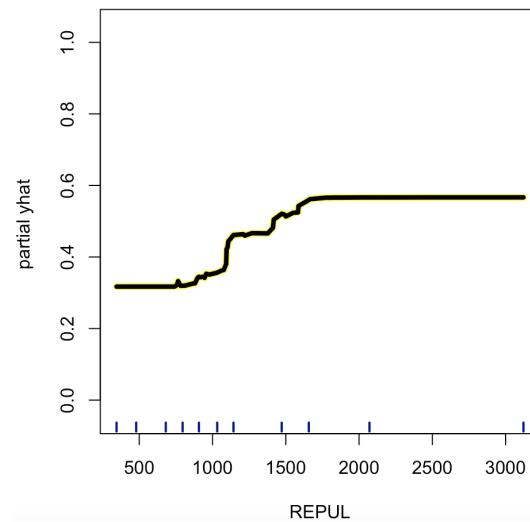


Interpretability of Predictive Models

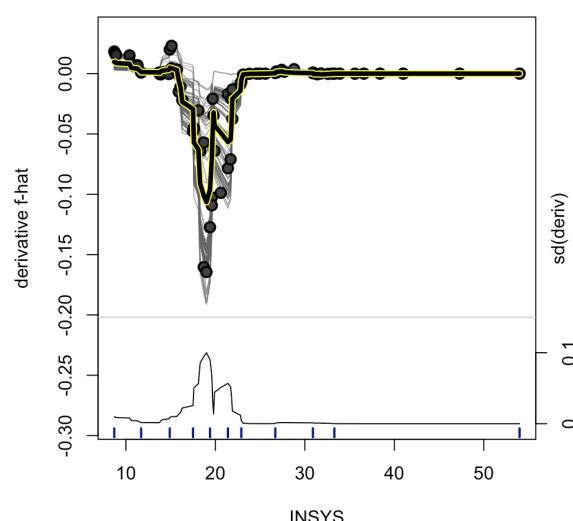
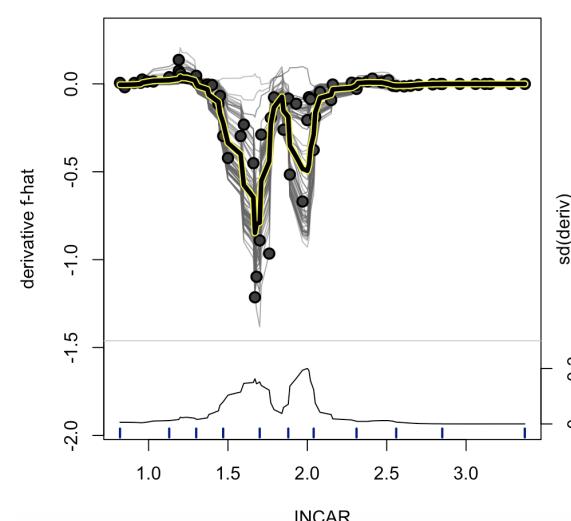
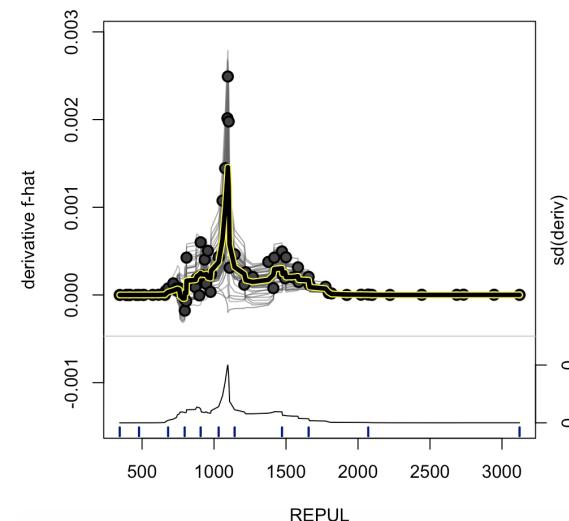
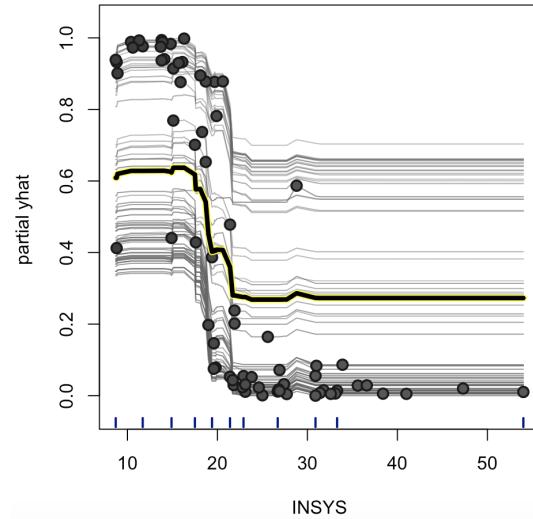
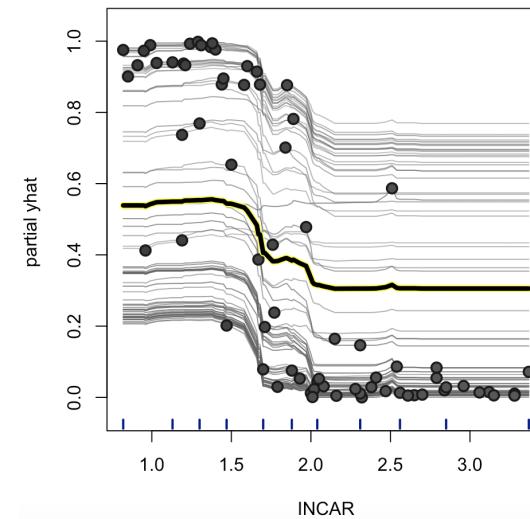
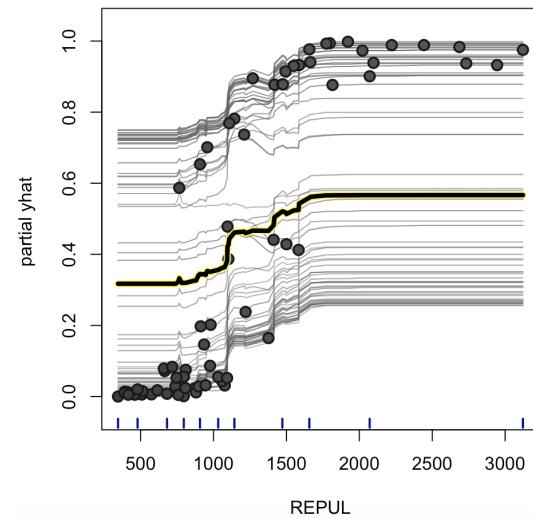
- Individual Conditional Expectation

Extention of Partial Dependence Plots, introduced in Goldstein *et al.* (2013, [Peeking inside the black box: Visualizing statistical learning with plots of individual conditional expectation](#)), “Visually, ICE plots disaggregate the output of classical PDPs. Rather than plot the target covariates’average partial effect on the predicted response, we instead plot the n estimated conditional expectation curves”

See `ICEbox` package and `ICEbox::ice`. Here m is a random forest,



Interpretability of Predictive Models



Interpretability of Predictive Models

- Accumulated Local Effects

Introduced in Apley (2016, [Visualizing the effects of predictor variables in black box supervised learning models](#)). Partial dependence of \mathbf{x}_s is

$$p(\mathbf{x}_s) = \mathbb{E}[m(\mathbf{x}_s, \mathbf{S}_c)] \text{ and } \hat{p}(\mathbf{x}_s) = \frac{1}{n} \sum_{i=1}^n \hat{m}(\mathbf{x}_s, \mathbf{x}_{i,c})$$

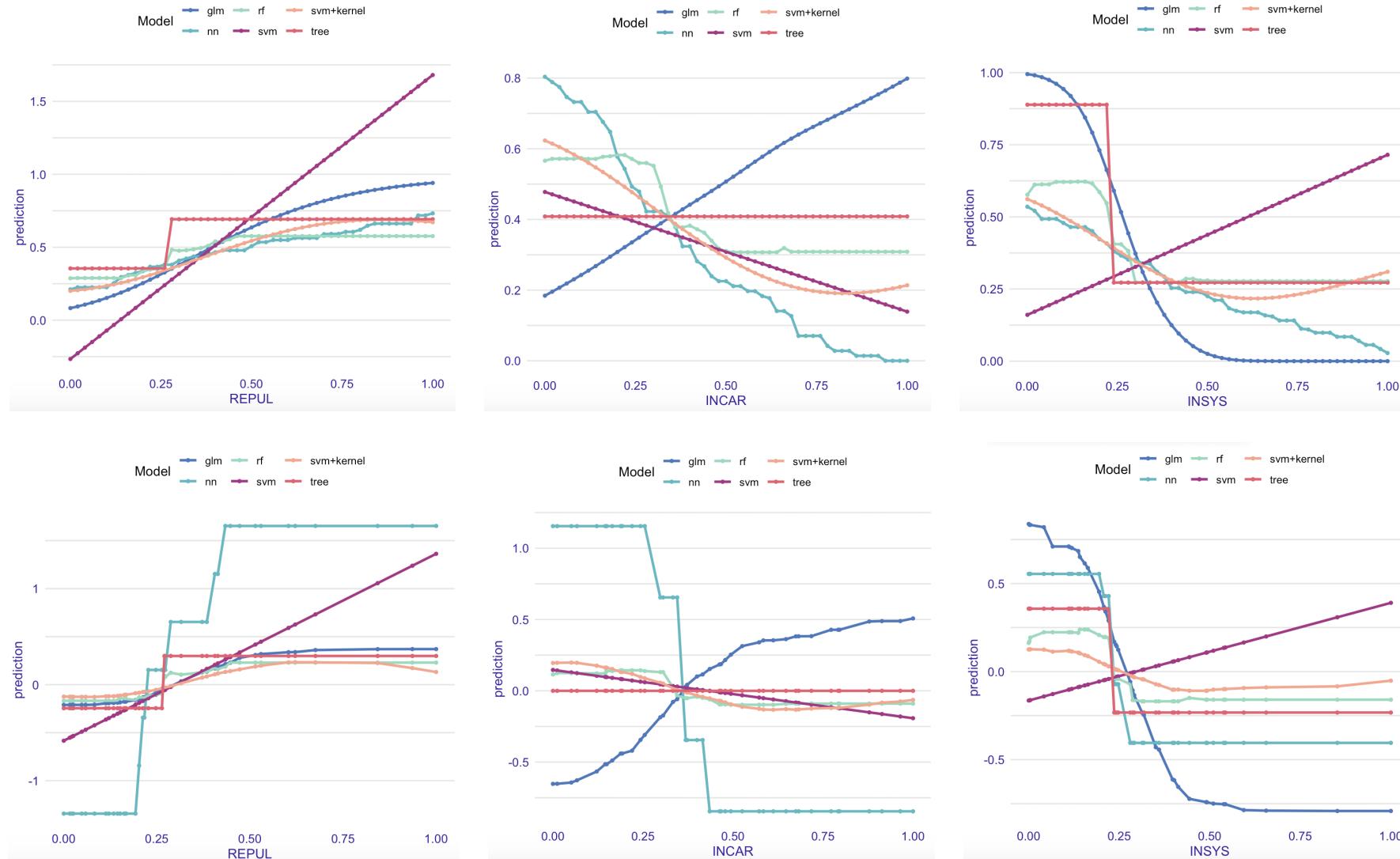
Here, we focus on

$$a(\mathbf{x}_s) = \int_{-\infty}^{\mathbf{x}_s} \mathbb{E} \left[\frac{\partial m(\mathbf{z}_s, \mathbf{S}_c)}{\partial \mathbf{x}_s} \right] d\mathbf{z}_s$$

See `DALEX` package, see <https://pbiecek.github.io> and function

```
DALEX::single_variable(explain(), variable = x, type = "ale") - as there is a  
DALEX::single_variable(explain(), variable = x, type = "pdp")
```

Partial Dependence Plot (on top) and Accumulated Local Effects (below)

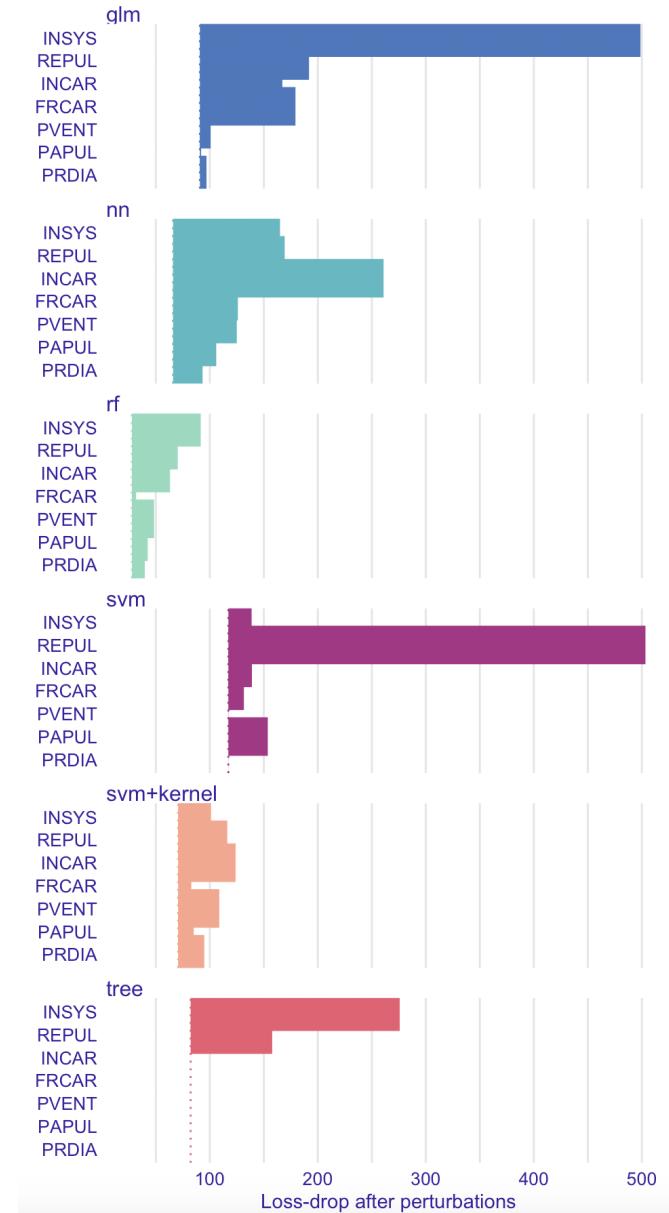


Interpretability of Predictive Models

- Feature Interaction

Friedman & Popescu (2008, ([Predictive learning via rule ensembles](#)), see Greenwell *et al.* ([A simple and effective model-based variable importance measure](#))

`iml::Interaction`



Interpretability of Predictive Models

- Feature (Variable) Importance

Breiman (2001, Random Forests)

`iml::FeatureImp` Or `DALEX::variable_importance`

Interpretability of Predictive Models

- Local Surrogate (LIME) - Local Interpretable Model-Agnostic Explanations

Alvarez-Melis & Jaakkola (2018, [On the robustness of interpretability methods](#))

`lime::explain` (see the [vignette](#) or the [ceterisParibus package](#) (and the *what-if* plot))

Interpretability of Predictive Models

- Shapley Value

Owen & Prieur (2016, [On Shapley value for measuring importance of dependent inputs](#))

The value of covariates $\{\mathbf{x}_u\}$, with $u \subset \{1, 2, \dots, p\}$ in a model m is

$$\text{val}(u) = \text{Var} [\mathbb{E}(m(\mathbf{x})|\mathbf{x}_u)]$$

going from 0 when $u = \emptyset$, σ^2 when $u = \{1, 2, \dots, p\}$.

We wish to define contributions φ_k satisfying

- efficiency: $\sum_{k=1}^p \varphi_k = \text{val}(\{1, 2, \dots, p\})$
- symmetry: if $\text{val}(u \cup \{k_1\}) = \text{val}(u \cup \{k_2\})$ where $k_1, k_2 \notin u$, then $\varphi_{k_1} = \varphi_{k_2}$
- dummy: if $\text{val}(u \cup \{k\}) = \text{val}(u)$ then $\varphi_k = 0$
- additivity: if we consider two value functions val_1 and val_2 , with contributions φ_1 and φ_2 , then value $\text{val} = \text{val}_1 + \text{val}_2$ has contributions $\varphi = \varphi_1 + \varphi_2$

Interpretability of Predictive Models

Shapley Value

Shapley (1953, [A Value for \$n\$ -person Game](#)) proved that the only solution is

$$\varphi_k = \sum_{u \in \{1, \dots, p\} \setminus \{k\}} \binom{p-1}{|u|}^{-1} [\text{val}(u \cup \{k\}) - \text{val}(u)]$$

For a linear model m ,

$$\varphi_k = \sum_{u \in \{1, \dots, p\} \setminus \{k\}} \binom{p-1}{|u|}^{-1} [R_{u \cup \{k\}}^2 - R_u^2]$$

see `iml::Shapley`

Calculating the Shapley Value is computationally expensive.

Molnar (2019, [Interpretable Machine Learning](#)) about *what-if* tools