

How to open the black box

Individual prediction explanation

Martin Jullum

Joint work with Kjersti Aas, Anders Løland and
Annabelle Redelmeier

NTNU Statistics Seminar
Trondheim, March 9th 2020



Example: Bank creates mortgage robot



Transaction history

Commercial Card - 456480111111111: 08/04/2004 - 14/04/2004

Date Processed	Description	Debit	Credit
08/04/2004	CASH ADVANCE FEE		\$5.00-
09/04/2004	SUNSHINE VILLAGE BANFF	\$86.97-	
10/04/2004	PRINCIPAL CREDIT ADJUSTMENT		\$22.00-
10/04/2004	CARD MEMBERSHIP FEE	\$19.00-	
10/04/2004	PHOTOCARD FEE	\$3.00-	
11/04/2004	PRINCIPAL DEBIT ADJUSTMENT	\$22.00-	
11/04/2004	PRINCIPAL DEBIT ADJUSTMENT	\$22.00-	
12/04/2004	PRINCIPAL CREDIT ADJUSTMENT	\$22.00-	

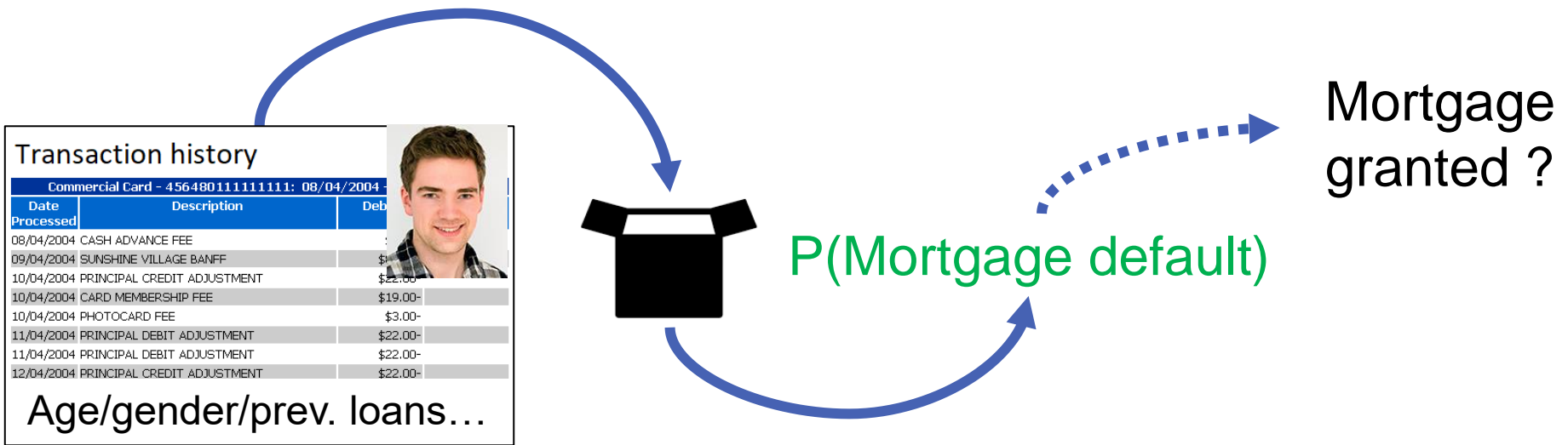
Age/gender/prev. loans...

&

Defaulted
loan?



Example: Bank creates mortgage robot



$$x \longrightarrow f(x) \longrightarrow p = 0.7$$

Why was



rejected a loan?

Why is this important?



- ▶ Customers may have a “right to an explanation”
- ▶ Builds trust to the “robot”

Individual prediction explanation

- ▶ **NOT** a general explanation of the black-box model
- ▶ $x = x^*$: Transaction history/covariates for



Explanation for $f(x^*) = 70\%$

A (mathematical) description/visualization/
characterization of how **each of the covariates**
contributed/affected the specific prediction $f(x^*) = 70\%$

Explaining a simple linear model

- ▶ Model $y = f(\mathbf{x}) = \beta_0 + \beta_1 x_1 + \beta_2 x_2$
- ▶ How would you explain predictions from this model?
- ▶ Explanation depends on:
 - β_1 and β_2
 - β_0
 - x_1^* and x_2^*
 - $E[x_1]$ and $E[x_2]$
 - $sd(x_1)$ and $sd(x_2)$
 - $corr(x_1, x_2)$
- ▶ Claim: A simple linear model is only easily interpretable if x_1 and x_2 are independent and standardized!

Prediction explanation frameworks

▶ Model-specific methods:

- **Deep Lift/Relevance propagation:** For neural networks
- **TreeSHAP:** For tree based methods

▶ Model-agnostic methods:

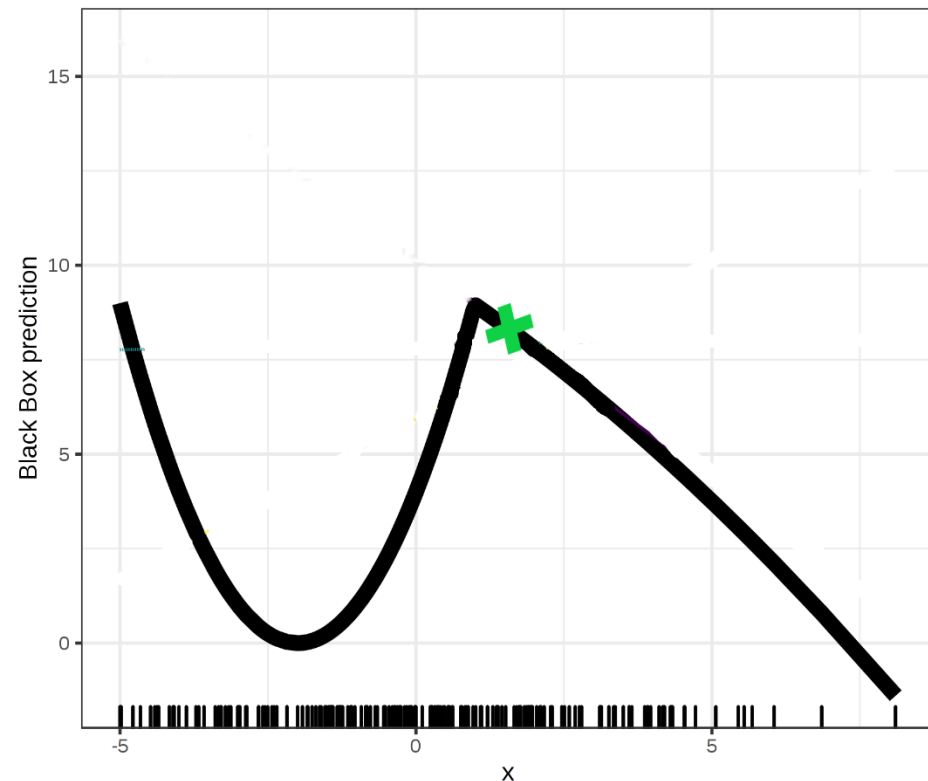
- **LIME** Local linear regression
- **Counterfactual explanations:** Which covariates should be altered to obtain a different decision?
- **Shapley values** Based on concepts from game theory



LIME

(Local Interpretable Model-agnostic Explanation)

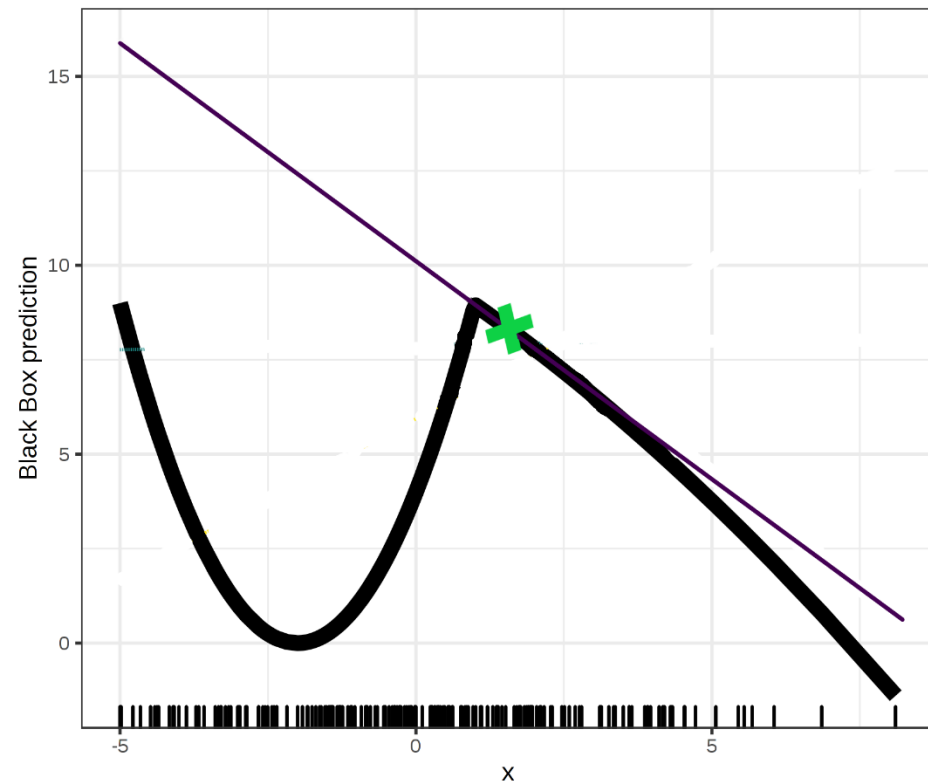
- ▶ Fits a (local) weighted linear regression model to $f(x)$ based on standardized covariates and weight determined by distance to x^*
- ▶ Importance score for each covariate: Coefficient from local model



LIME

(Local Interpretable Model-agnostic Explanation)

- ▶ Fits a (local) weighted linear regression model to $f(x)$ based on standardized covariates and weight determined by distance to x^*
- ▶ Importance score for each covariate: Coefficient from local model



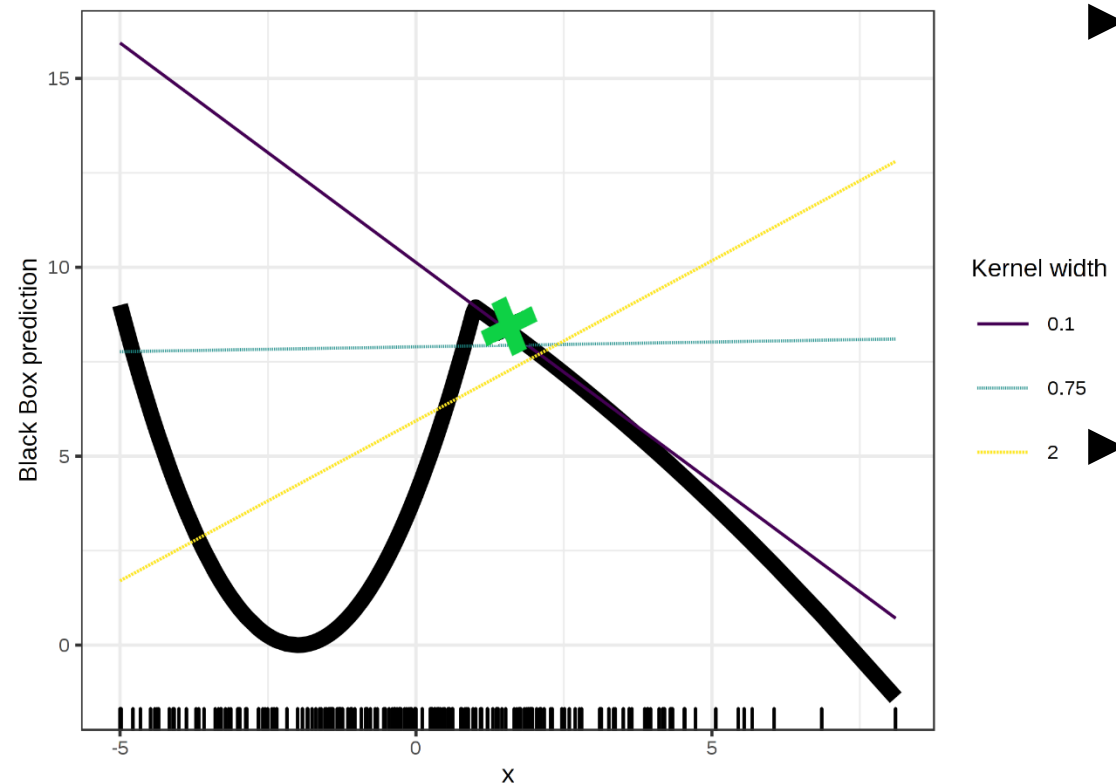
▶ Challenges

- Defining the distance and weight functions

LIME

(Local Interpretable Model-agnostic Explanation)

- ▶ Fits a (local) weighted linear regression model to $f(x)$ based on standardized covariates and weight determined by distance to x^*
- ▶ Importance score for each covariate: Coefficient from local model



▶ Challenges

- Defining the distance and weight functions
- Direct use of local model coefficients

▶ Advantages

- Simple idea
- Easy to use

Counterfactual explanations

- ▶ What is the smallest covariate change necessary to change the prediction “significantly”?
- ▶ Optimization problem:

$$(Ex) \quad \arg \min_{\mathbf{x}'} d(\mathbf{x}^*, \mathbf{x}'), \quad \text{subject to } |f(\mathbf{x}') - [f(\mathbf{x}^*) + \lambda]| \leq \varepsilon$$

- ▶ Explanation: Minimizers of (Ex)
- ▶ Challenges:
 - Choosing d, λ and ε
 - May lead to many sub-explanations
 - Need to interpret the explanations yourself
- ▶ Advantages
 - Cannot be wrong
 - Guides user on how to change prediction

Shapley values

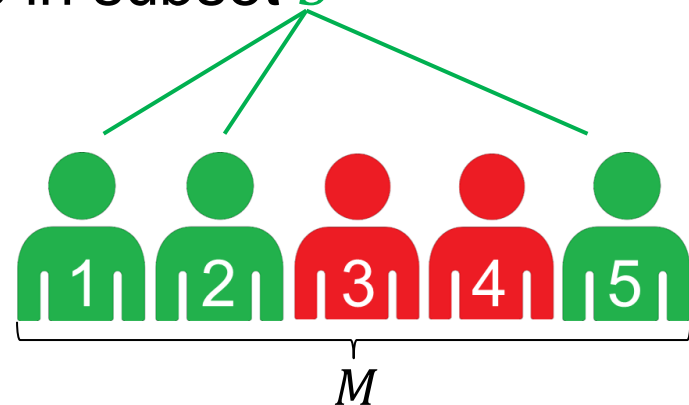


- ▶ Concept from (cooperative) game theory in the 1950s
- ▶ Used to distribute the total payoff to the players
- ▶ Explicit formula for the “fair” payment to every player j :

$$\phi_j = \sum_{S \subseteq M \setminus \{j\}} w(S) (v(S \cup \{j\}) - v(S)), \quad w(S) \text{ is a weight function}$$

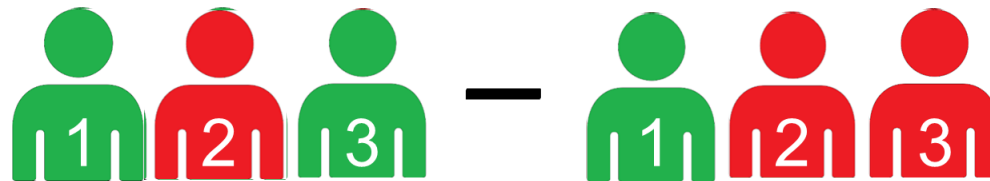
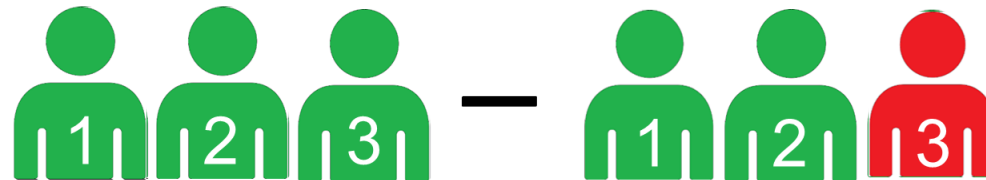
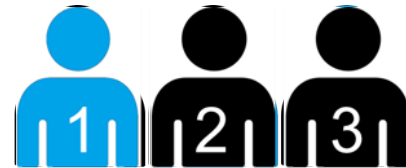
$v(S)$ is the payoff with only players in subset S

- ▶ Several mathematical optimality properties



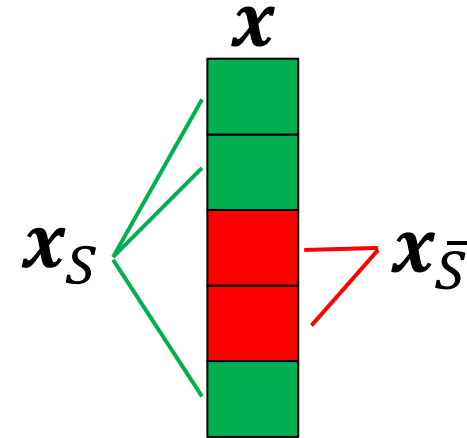
Intuition behind the Shapley formula

Game with 3 players



Shapley values for prediction explanation

- ▶ Players = covariates (x_1, \dots, x_M)
- ▶ Payoff = prediction $(f(\mathbf{x}^*))$
- ▶ Contribution function: $v(S) = E[f(\mathbf{x}) | \mathbf{x}_S = \mathbf{x}_S^*]$
- ▶ Properties



$$f(\mathbf{x}^*) = \sum_{j=0}^M \phi_j$$

$$\phi_0 = E[f(\mathbf{x})]$$

$$E[f(\mathbf{x})] = E[f(\mathbf{x}) | x_j]$$

implies $\phi_j = 0$

$$x_i, x_j \text{ same contribution}$$

implies $\phi_i = \phi_j$

- ▶ Rough interpretation of ϕ_j : **The prediction change when you don't know the value of x_j** -- averaged over all covariates

Shapley values for prediction explanation

- ▶ 2 main challenges

1. The computational complexity in the Shapley formula

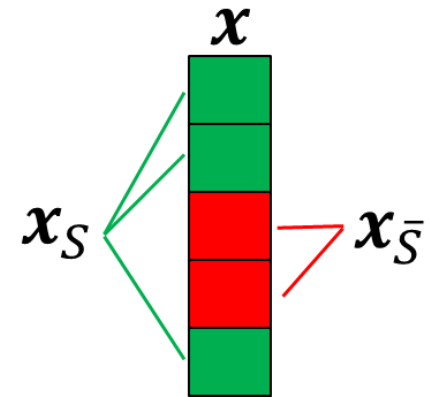
$$\phi_j = \sum_{S \subseteq M \setminus \{j\}} w(S) (v(S \cup \{j\}) - v(S))$$

- Partly solved by cleverly reducing the sum by subset sampling (KernelSHAP; Lundberg & Lee, 2017)

Shapley values for prediction explanation

- ▶ 2 main challenges

Recall



2. Estimating the contribution function

$$v(S) = E[f(\mathbf{x}) | \mathbf{x}_S = \mathbf{x}_S^*] = \int f(\mathbf{x}_{\bar{S}}, \mathbf{x}_S) p(\mathbf{x}_{\bar{S}} | \mathbf{x}_S = \mathbf{x}_S^*) d\mathbf{x}_{\bar{S}}$$

- Previous methods

- Approximates $v(S) \approx \int f(\mathbf{x}_{\bar{S}}, \mathbf{x}_S^*) p(\mathbf{x}_{\bar{S}}) d\mathbf{x}_{\bar{S}}$,
- Estimates $p(\mathbf{x}_{\bar{S}})$ using the empirical distribution of the training data
- Monte Carlo integration to solve the integral

This assumes covariates are independent!

Consequences of the independence assumption

- ▶ Requires evaluating $f(x_{\bar{S}}, x_S)$ at potentially unlikely or illegal combinations of $x_{\bar{S}}$ and x_S
- ▶ Example 1
 - Number of transactions to Switzerland: 0
 - Average transaction amount to Switzerland: 1000 NOK
- ▶ Example 2
 - Age: 17
 - Marital status: Widow
 - Profession: Professor



Shapley values for prediction explanation

- ▶ Explicit formula for a linear model $f(\mathbf{x}) = \beta_0 + \sum_{j=1}^M \beta_j x_j$ with **independent** covariates:

$$\phi_0 = \beta_0 + \sum_{j=1}^M \beta_j E[x_j], \quad \phi_j = \beta_j (x_j^* - E[x_j]), \quad j = 1, \dots, M$$

- ▶ Advantages

- Proper mathematical foundation
- Desirable set of properties

- ▶ Challenges

- Computationally heavy
- Requires good estimates of a difficult estimation problem:
 $E[f(\mathbf{x}) | \mathbf{x}_S = \mathbf{x}_S^*]$

- ▶ From our perspective the method with greatest potential – what we have work with the last two years

Our idea

Estimate $p(\mathbf{x}_{\bar{S}}|\mathbf{x}_S = \mathbf{x}_S^*)$ properly

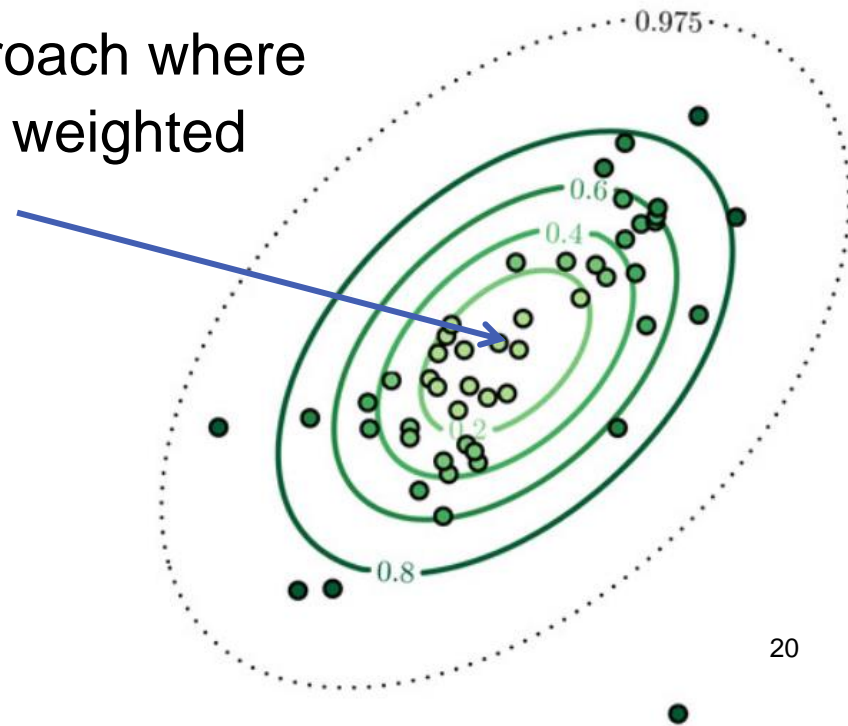
+

Monte Carlo integration to approximate

$$v(S) = E[f(\mathbf{x})|\mathbf{x}_S = \mathbf{x}_S^*] = \int f(\mathbf{x}_{\bar{S}}, \mathbf{x}_S) p(\mathbf{x}_{\bar{S}}|\mathbf{x}_S = \mathbf{x}_S^*) d\mathbf{x}_{\bar{S}}$$

Continuous covariates

- ▶ How to estimate $p(\mathbf{x}_{\bar{S}}|\mathbf{x}_S = \mathbf{x}_S^*)$ when \mathbf{x} is continuous?
- ▶ 3 approaches
 - Assume $p(\mathbf{x})$ Gaussian => analytical $p(\mathbf{x}_{\bar{S}}|\mathbf{x}_S = \mathbf{x}_S^*)$
 - Assume Gaussian copula => transformation + analytical expression
 - An empirical (conditional) approach where training observations at $\mathbf{x}_{\bar{S}}^i$ are weighted based on proximity of \mathbf{x}_S^i to \mathbf{x}_S^*



Empirical conditional approach

1. Compute the scaled Mahalanobis distance between \mathbf{x}_S^* and the columns \mathcal{S} of the training data $\mathbf{x}^1, \dots, \mathbf{x}^n$

$$D_{\mathcal{S}}(\mathbf{x}^*, \mathbf{x}^i) = \sqrt{\frac{(\mathbf{x}_S^* - \mathbf{x}_S^i)^T \Sigma_S^{-1} (\mathbf{x}_S^* - \mathbf{x}_S^i)}{|\mathcal{S}|}}$$

2. Use Gaussian kernel to get weight of each training observation:

$$w_{\mathcal{S}}(\mathbf{x}^*, \mathbf{x}^i) = \exp\left(-\frac{D_{\mathcal{S}}(\mathbf{x}^*, \mathbf{x}^i)^2}{2\sigma^2}\right)$$

3. Approximate $p(\mathbf{x}_{\bar{S}} | \mathbf{x}_S = \mathbf{x}_S^*)$ by the probability mass function where $p(\mathbf{x}_{\bar{S}} = \mathbf{x}^i | \mathbf{x}_S = \mathbf{x}_S^*) = \frac{w_{\mathcal{S}}(\mathbf{x}^*, \mathbf{x}^i)}{\sum_{k=1}^n w_{\mathcal{S}}(\mathbf{x}^*, \mathbf{x}^k)}$

Empirical conditional approach II

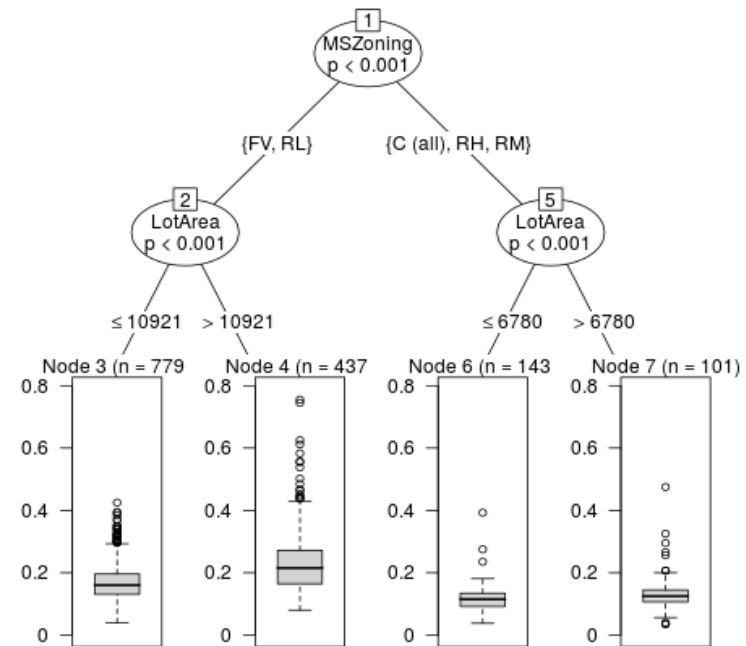
- ▶ This gives an estimator of $E[f(\mathbf{x})|\mathbf{x}_S = \mathbf{x}_S^*]$:

$$\hat{v}(S) = \frac{\sum_{k=1}^n w_S(\mathbf{x}^*, \mathbf{x}^k) f(\mathbf{x}_S^k, \mathbf{x}_S^*)}{\sum_{k=1}^n w_S(\mathbf{x}^*, \mathbf{x}^k)}$$

- ▶ It turns out that we re-invented the Nadaraya-Watson estimator (locally constant kernel estimator) aiming at estimating $E[U|V=v]$ for responses $u_i = f(\mathbf{x}_S^i, \mathbf{x}_S^*)$, and covariates $v_i = \mathbf{x}_S^i, i = 1, \dots, n$
- ▶ May then use a corrected AIC-criterion by Hurvich and Tsai (JRSS-B, 1998) to select the bandwidth parameter σ .

Categorical/mixed covariates

- ▶ How to estimate $p(\mathbf{x}_{\bar{S}} | \mathbf{x}_S = \mathbf{x}_S^*)$ when \mathbf{x} is categorical, or mixed continuous/categorical
 - Fit a **multivariate decision tree** to $U = \mathbf{x}_{\bar{S}}$ based on $V = \mathbf{x}_S$ using the training data
 - Approximate $p(\mathbf{x}_{\bar{S}} | \mathbf{x}_S = \mathbf{x}_S^*)$ by the empirical distribution of the training observations ($\mathbf{x}_{\bar{S}}$) within the terminal node of $\mathbf{x}_S = \mathbf{x}_S^*$



Multivariate decision tree

- ▶ Classical decision tree algorithms like CART work only for univariate responses
 - Multivariate generalizations exists
 - CARTs are known to be biased towards splitting on categorical covariates with many levels
- ▶ Instead, we rely on **Recursive partitioning/conditional inference trees** (Hothorn et al., 2006)
 - Decide which covariate to split on first
 - Then decides on the splitting point for that covariate

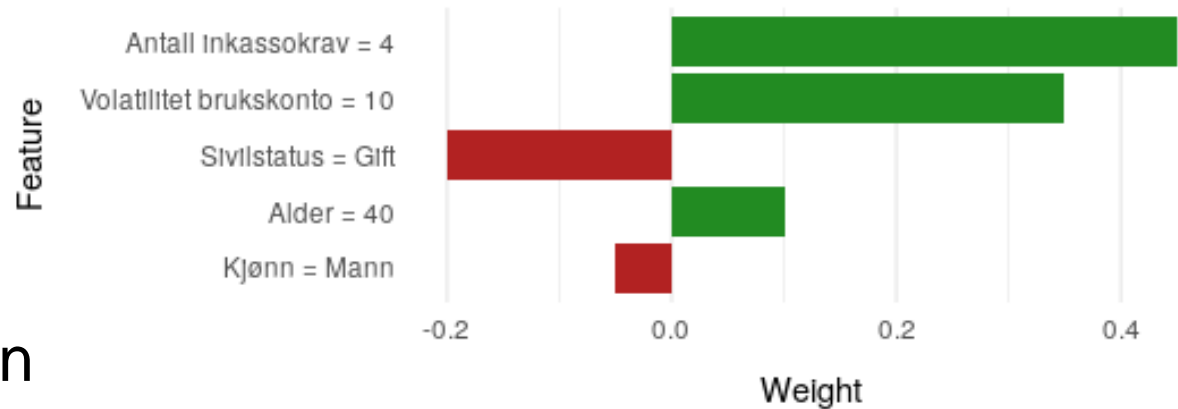
Conditional inference tree algorithm

- ▶ Multivariate response \mathbf{U} , covariates V_1, \dots, V_p
- ▶ Step 1: Decide whether or not to split by hypothesis testing:

$$H_0: p(\mathbf{U}|V_j) = p(\mathbf{U}) \quad \forall j \quad \text{vs} \quad H_A: p(\mathbf{U}|V_j) \neq p(\mathbf{U}) \text{ for some } j$$

- Hypothesis test performed by permutation test using a summary statistic for the dependence between \mathbf{U} and V_j
 - Stop tree building if not rejecting H_0 at a level α
 - If rejecting H_0 , pick the covariate with the smallest p -value.
- ▶ Step 2: Splitting criteria
 - Maximize a two-sample discrepancy statistic
 - ▶ Implemented in the R-packages *party* and *partykit*

Conclusion



- ▶ Individual prediction explanation, i.e. explaining $f(x^*)$ for specific covariate x^*
- ▶ Not straightforward to explain even a simple linear model
- ▶ Mainly three model-agnostic methods in the literature:
 - LIME, counterfactual explanations, Shapley values
- ▶ No grand truth when explaining predictions!
- ▶ Ignoring dependence between covariates can give completely wrong explanations

Want to know more?

Read our paper on arXiv
arxiv.org/abs/1903.10464



Check out our R-package
shapr on GitHub (soon CRAN) + JOSS
github.com/NorskRegnesentral/shapr