

Proper prediction explanation with *shapr*

Martin Jullum



NAV
Oslo, June 26th 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-	

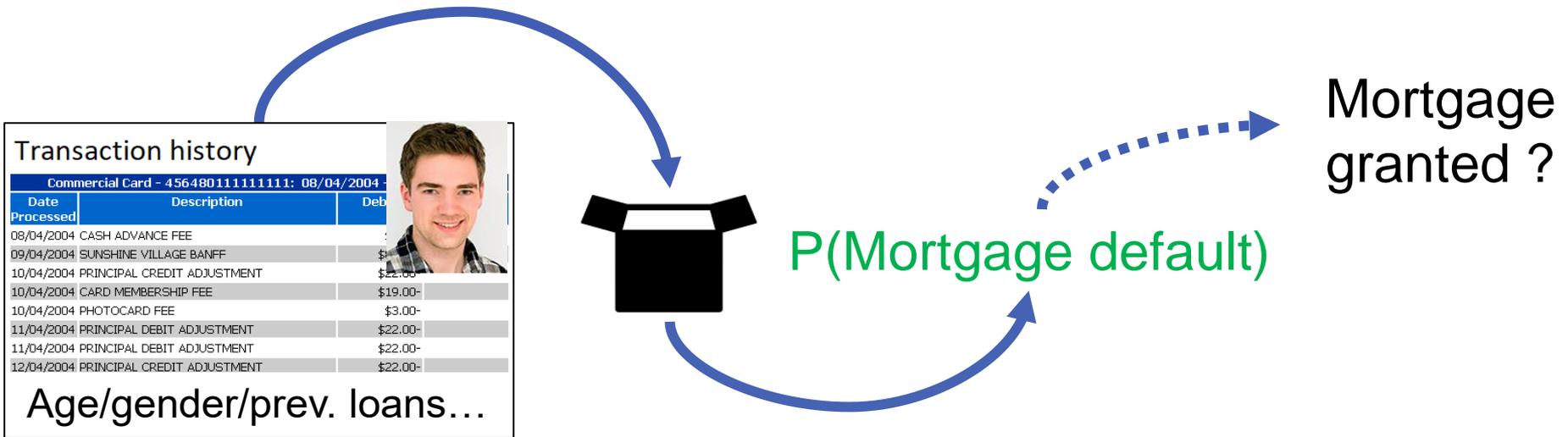
Age/gender/prev. loans...

&

Defaulted
loan?



Example: Bank creates mortgage robot



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

Why was



rejected a loan?

Individual prediction explanation

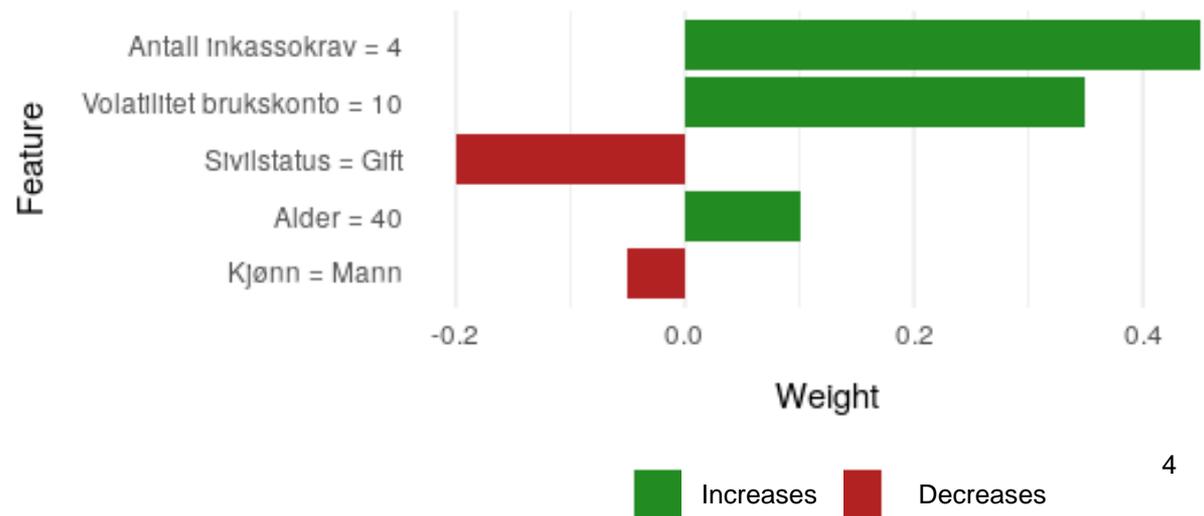
- ▶ **NOT** a general explanation of the black-box model

- ▶ $x = x^*$: Transaction history/features for



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

- ▶ How did each feature contribute to increase/decrease the predicted value $f(x^*) = 70\%$?



Why is this important?



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

Shapley values

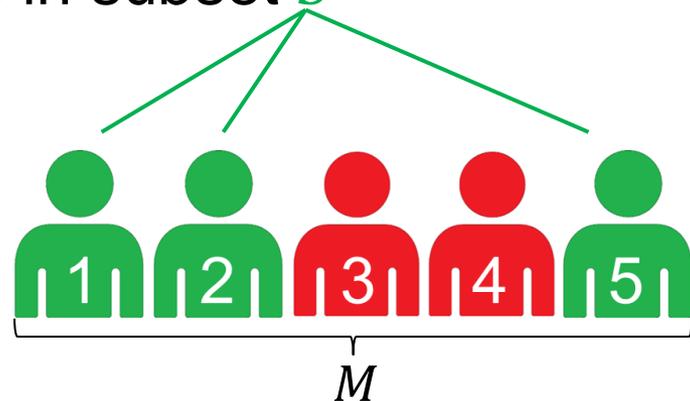


- ▶ Concept from (cooperative) game theory in the 1950s
- ▶ Used to distribute the total payoff to the players
- ▶ Explicit formula for the “fair” payment to every player j :

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

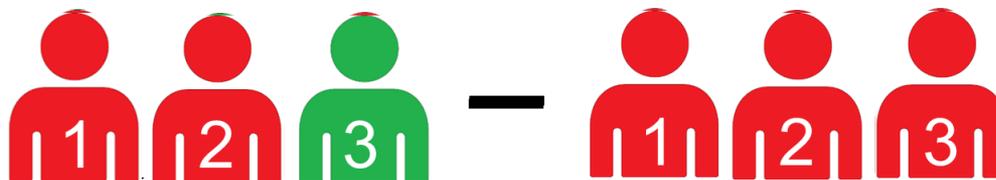
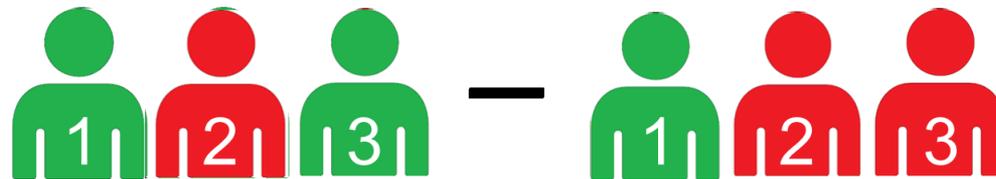
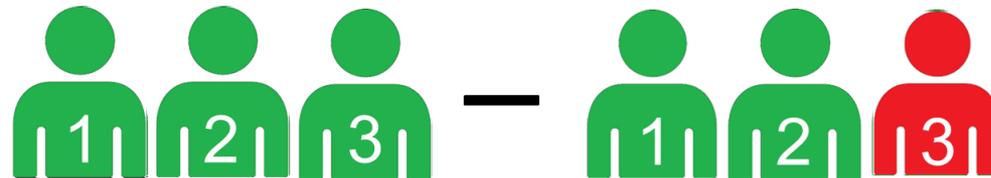
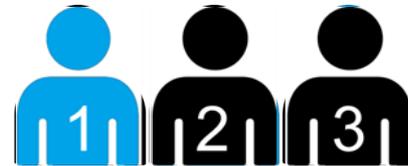
$v(S)$ is the payoff with only players in subset S

- ▶ Several mathematical optimality properties



Intuition behind the Shapley formula

Game with 3 players



Shapley values for prediction explanation

- ▶ Players = features (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

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

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

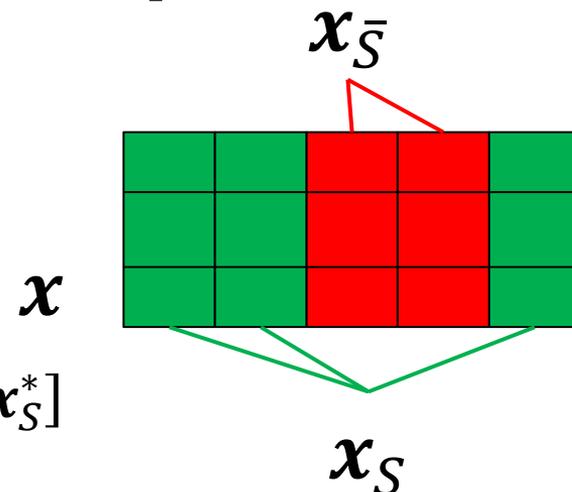
$$E[f(\mathbf{x})] = E[f(\mathbf{x}) | x_j]$$

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 features**



Shapley values for prediction explanation

► 2 main challenges

1. The computational complexity in the Shapley formula

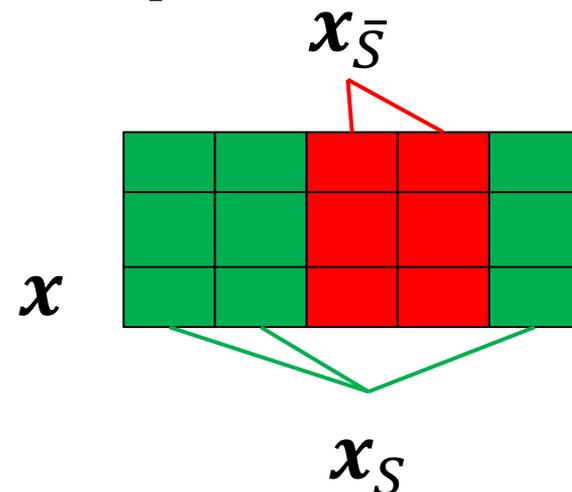
$$\phi_j = \sum_{S \subseteq M \setminus \{j\}} w(S) (v(S \cup \{j\}) - v(S))$$

- 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



Our idea

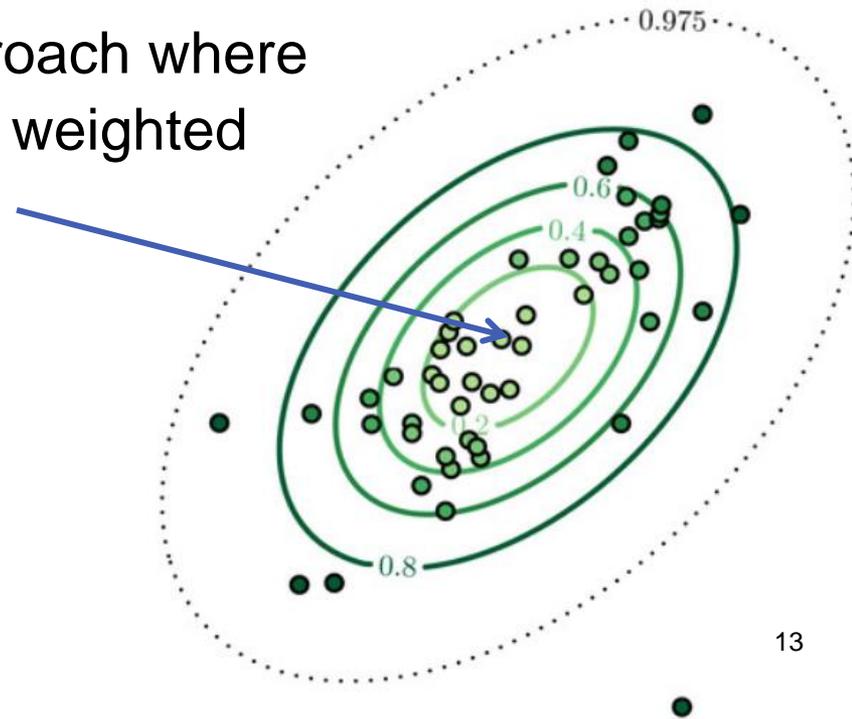
Estimate $p(x_{\bar{S}} | x_S = x_S^*)$ properly

+

Monte Carlo integration

Continuous features

- ▶ 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^*



Explaining sick leave predictions



- ▶ NAV is modelling how long individuals are on sick leave
 - Used by case workers to schedule follow-up meetings
- ▶ Case workers need to understand the “reasoning” of the individual predictions
- ▶ Modelling based on
 - age, gender, sick leave history, type of business etc.
 - Several **categorical** features with **many levels**
- ▶ Need methodology for prediction explanation which can handle categorical features

Our idea

Estimate $p(x_{\bar{S}} | x_S = x_S^*)$ properly

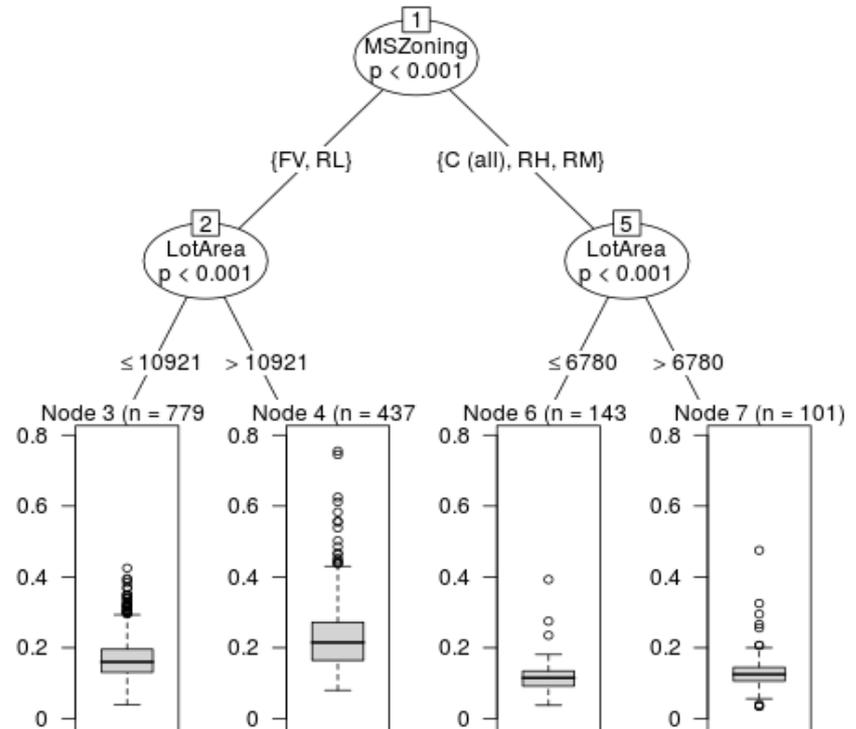
+

Monte Carlo integration

Our contribution: categorical/mixed variables



- ▶ Estimating $p(\mathbf{x}_{\bar{S}} | \mathbf{x}_S = \mathbf{x}_S^*)$
 - For every subset S , fit a **(multivariate) regression tree** to $\mathbf{y} = \mathbf{x}_{\bar{S}}$ based on \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 $\mathbf{x}_S = \mathbf{x}_S^*$



How to use this is practice?

- ▶ All of this is implemented on our R-package *shapr* on GitHub (soon CRAN) github.com/NorskRegnesentral/shapr



- ▶ Paper (continuous variables): <https://arxiv.org/abs/1903.10464>
- ▶ Paper (categorical/mixed variables): <https://arxiv.org/abs/2007.01027>