# mSHAP
## Explaining Two-Part Models

Spencer Matthews

Joint work with Brian Hartman
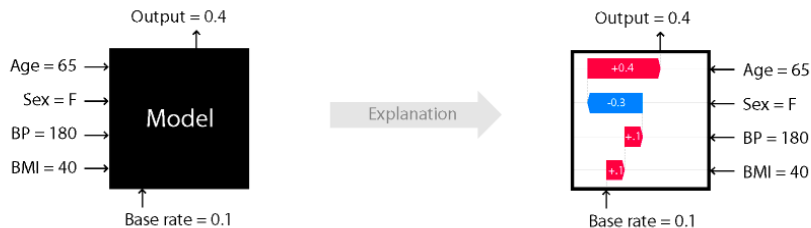
June 2021

Brigham Young University

- Two-part models are used by actuaries to set insurance rates, and therefore must be explainable
- Newer "black-box" methods (such as the gradient boosted forest) provide greater accuracy to these pricing models
- Although methods exist to explain individual models, there is not a good methodology to explain the predictions of a two-part model

https://github.com/slundberg/shap

## DEFINITIONS

- Three Models: $f, g,$ and $h$, where $h$ is the product of $f$ and $g$.
- Input Matrix: $A$ where $A_i$ is the $i$th row of $A$ and $A$ is $n \times p$ where $n$ is the number of observations and $p$ is the number of predictors.
- $f(A_i) = \hat{x}_i$, $g(A_i) = \hat{y}_i$, and $h(A_i) = \hat{z}_i$ and the contribution of the $j$th predictor to $\hat{x}_i$ as $s_{x_i j}$.
- $\mu_f, \mu_g, \mu_h$ signify the average model prediction over the data (known as the baseline term or expected model output)
- Based on the property of local accuracy:

$$\hat{x}_i = \mu_f + s_{x_i 1} + s_{x_i 2} + \ldots + s_{x_i p}$$

and

$$\hat{y}_i = \mu_g + s_{y_i 1} + s_{y_i 2} + \ldots + s_{y_i p}$$

# LOCAL ACCURACY

$$x_i = \mu_f + \sum_{j=1}^{p} s_{x_i j}$$

$\Longrightarrow$

The sum of the SHAP values and the expected model output must equal the model prediction
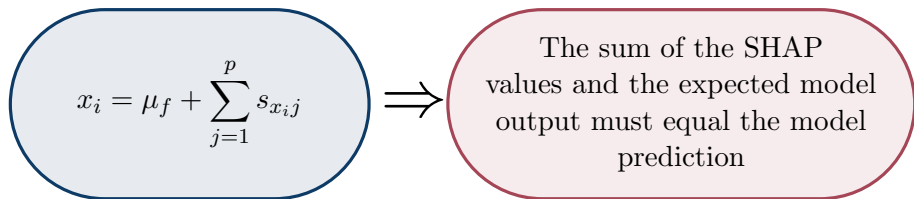
# Local Accuracy



$$x_i = \mu_f + \sum_{j=1}^{p} s_{x_i j}$$

$\implies$ The sum of the SHAP values and the expected model output must equal the model prediction

A Brief Example

# LOCAL ACCURACY

$$x_i = \mu_f + \sum_{j=1}^{p} s_{x_i j}$$

$\Longrightarrow$

The sum of the SHAP values and the expected model output must equal the model prediction

A Brief Example

$$\hat{x}_i = \mu_f + s_{x_i 1} + s_{x_i 2}$$

$$\hat{y}_i = \mu_g + s_{y_i 1} + s_{y_i 2}$$

# LOCAL ACCURACY

$$x_i = \mu_f + \sum_{j=1}^{p} s_{x_i j}$$

$\Longrightarrow$

The sum of the SHAP values and the expected model output must equal the model prediction

A Brief Example

$$\hat{x}_i = \mu_f + s_{x_i 1} + s_{x_i 2}$$

$$\hat{y}_i = \mu_g + s_{y_i 1} + s_{y_i 2}$$

$$\hat{x}_i \cdot \hat{y}_i$$

$$\neq$$

$$\mu_f \mu_g + s_{x_i 1} s_{y_i 1} + s_{x_i 2} s_{y_i 2}$$

# Expansion of Terms: $\hat{x}_i \times \hat{y}_i$

| | $s_{x_i 1}$ + | $s_{x_i 2}$ + | $s_{x_i 3}$ + $\ldots$ + | $s_{x_i p}$ + | $\mu_f$ |
|---|---|---|---|---|---|
| $s_{y_i 1}$ + | $s_{x_i 1} s_{y_i 1}$ | $s_{x_i 2} s_{y_i 1}$ | $s_{x_i 3} s_{y_i 1}$ $\quad \ldots$ | $s_{x_i p} s_{y_i 1}$ | $\mu_f s_{y_i 1}$ |
| $s_{y_i 2}$ + | $s_{x_i 1} s_{y_i 2}$ | $s_{x_i 2} s_{y_i 2}$ | $s_{x_i 3} s_{y_i 2}$ $\quad \ldots$ | $s_{x_i p} s_{y_i 2}$ | $\mu_f s_{y_i 2}$ |
| $s_{y_i 3}$ + | $s_{x_i 1} s_{y_i 3}$ | $s_{x_i 2} s_{y_i 3}$ | $s_{x_i 3} s_{y_i 3}$ $\quad \ldots$ | $s_{x_i p} s_{y_i 3}$ | $\mu_f s_{y_i 3}$ |
| $\vdots$ + | $\vdots$ | $\vdots$ | $\vdots$ $\quad \ddots$ | $\vdots$ | $\vdots$ |
| $s_{y_i n}$ + | $s_{x_i 1} s_{y_i p}$ | $s_{x_i 2} s_{y_i p}$ | $s_{x_i 3} s_{y_i p}$ $\quad \ldots$ | $s_{x_i p} s_{y_i p}$ | $\mu_f s_{y_i p}$ |
| $\mu_g$ | $s_{x_i 1} \mu_g$ | $s_{x_i 2} \mu_g$ | $s_{x_i 3} \mu_g$ $\quad \ldots$ | $s_{x_i p} \mu_g$ | $\mu_f \mu_g$ |

# EXPANSION OF TERMS: $\hat{x}_i \times \hat{y}_i$

|  | $s_{x_i1}$ + | $s_{x_i2}$ + | $s_{x_i3}$ + ... + | $s_{x_ip}$ + | $\mu_f$ |
|---|---|---|---|---|---|
| $s_{y_i1}$ | $s_{x_i1}s_{y_i1}$ | $s_{x_i2}s_{y_i1}$ | $s_{x_i3}s_{y_i1}$ ... | $s_{x_ip}s_{y_i1}$ | $\mu_f s_{y_i1}$ |
| + |  |  |  |  |  |
| $s_{y_i2}$ | $s_{x_i1}s_{y_i2}$ | $s_{x_i2}s_{y_i2}$ | $s_{x_i3}s_{y_i2}$ ... | $s_{x_ip}s_{y_i2}$ | $\mu_f s_{y_i2}$ |
| + |  |  |  |  |  |
| $s_{y_i3}$ | $s_{x_i1}s_{y_i3}$ | $s_{x_i2}s_{y_i3}$ | $s_{x_i3}s_{y_i3}$ ... | $s_{x_ip}s_{y_i3}$ | $\mu_f s_{y_i3}$ |
| + |  |  |  |  |  |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ $\ddots$ | $\vdots$ | $\vdots$ |
| + |  |  |  |  |  |
| $s_{y_in}$ | $s_{x_i1}s_{y_ip}$ | $s_{x_i2}s_{y_ip}$ | $s_{x_i3}s_{y_ip}$ ... | $s_{x_ip}s_{y_ip}$ | $\mu_f s_{y_ip}$ |
| + |  |  |  |  |  |
| $\mu_g$ | $s_{x_i1}\mu_g$ | $s_{x_i2}\mu_g$ | $s_{x_i3}\mu_g$ ... | $s_{x_ip}\mu_g$ | $\mu_f\mu_g$ |

# EXPANSION OF TERMS: $\hat{x}_i \times \hat{y}_i$

|  | $s_{x_i 1}$ + | $s_{x_i 2}$ + | $s_{x_i 3}$ + ... + | $s_{x_i p}$ + | $\mu_f$ |
|---|---|---|---|---|---|
| $s_{y_i 1}$ + | $s_{x_i 1} s_{y_i 1}$ | $s_{x_i 2} s_{y_i 1}$ | $s_{x_i 3} s_{y_i 1}$ ... | $s_{x_i p} s_{y_i 1}$ | $\mu_f s_{y_i 1}$ |
| $s_{y_i 2}$ + | $s_{x_i 1} s_{y_i 2}$ | $s_{x_i 2} s_{y_i 2}$ | $s_{x_i 3} s_{y_i 2}$ ... | $s_{x_i p} s_{y_i 2}$ | $\mu_f s_{y_i 2}$ |
| $s_{y_i 3}$ + | $s_{x_i 1} s_{y_i 3}$ | $s_{x_i 2} s_{y_i 3}$ | $s_{x_i 3} s_{y_i 3}$ ... | $s_{x_i p} s_{y_i 3}$ | $\mu_f s_{y_i 3}$ |
| $\vdots$ + | $\vdots$ | $\vdots$ | $\ddots$ | $\vdots$ | $\vdots$ |
| $s_{y_i n}$ + | $s_{x_i 1} s_{y_i p}$ | $s_{x_i 2} s_{y_i p}$ | $s_{x_i 3} s_{y_i p}$ ... | $s_{x_i p} s_{y_i p}$ | $\mu_f s_{y_i p}$ |
| $\mu_g$ | $s_{x_i} \mu_g$ | $s_{x_i 2} \mu_g$ | $s_{x_i 3} \mu_g$ ... | $s_{x_i p} \mu_g$ | $\mu_f \mu_g$ |

| | $s_{x_i1}$ + | $s_{x_i2}$ + | $s_{x_i3}$ + | ... + | $s_{x_ip}$ + | $\mu_f$ |
|---|---|---|---|---|---|---|
| $s_{y_i1}$ + | $s_{x_i1}s_{y_i1}$ | $s_{x_i2}s_{y_i1}$ | $s_{x_i3}s_{y_i1}$ | ... | $s_{x_ip}s_{y_i1}$ | $\mu_f s_{y_i1}$ |
| $s_{y_i2}$ + | $s_{x_i1}s_{y_i2}$ | $s_{x_i2}s_{y_i2}$ | $s_{x_i3}s_{y_i2}$ | ... | $s_{x_ip}s_{y_i2}$ | $\mu_f s_{y_i2}$ |
| $s_{y_i3}$ + | $s_{x_i1}s_{y_i3}$ | $s_{x_i2}s_{y_i3}$ | $s_{x_i3}s_{y_i3}$ | ... | $s_{x_ip}s_{y_i3}$ | $\mu_f s_{y_i3}$ |
| $\vdots$ + | $\vdots$ | $\vdots$ | $\vdots$ | $\ddots$ | $\vdots$ | $\vdots$ |
| $s_{y_in}$ + | $s_{x_i1}s_{y_ip}$ | $s_{x_i2}s_{y_ip}$ | $s_{x_i3}s_{y_ip}$ | ... | $s_{x_ip}s_{y_ip}$ | $\mu_f s_{y_ip}$ |
| $\mu_g$ | $s_{x_i1}\mu_g$ | $s_{x_i2}\mu_g$ | $s_{x_i3}\mu_g$ | ... | $s_{x_ip}\mu_g$ | $\mu_f\mu_g$ |

# Expansion of Terms: $\hat{x}_i \times \hat{y}_i$

| | $s_{x_i1}$ | $+$ | $s_{x_i2}$ | $+$ | $s_{x_i3}$ | $+ \ldots +$ | $s_{x_ip}$ | $+$ | $\mu_f$ |
|---|---|---|---|---|---|---|---|---|---|
| $s_{y_i1}$ $+$ | $s_{x_i1}s_{y_i1}$ | | $s_{x_i2}s_{y_i1}$ | | $s_{x_i3}s_{y_i1}$ | $\ldots$ | $s_{x_ip}s_{y_i1}$ | | $\mu_f s_{y_i1}$ |
| $s_{y_i2}$ $+$ | $s_{x_i1}s_{y_i2}$ | | $s_{x_i2}s_{y_i2}$ | | $s_{x_i3}s_{y_i2}$ | $\ldots$ | $s_{x_ip}s_{y_i2}$ | | $\mu_f s_{y_i2}$ |
| $s_{y_i3}$ $+$ | $s_{x_i1}s_{y_i3}$ | | $s_{x_i2}s_{y_i3}$ | | $s_{x_i3}s_{y_i3}$ | $\ldots$ | $s_{x_ip}s_{y_i3}$ | | $\mu_f s_{y_i3}$ |
| $\vdots$ $+$ | $\vdots$ | | $\vdots$ | | $\vdots$ | $\ddots$ | $\vdots$ | | $\vdots$ |
| $s_{y_in}$ $+$ | $s_{x_i1}s_{y_ip}$ | | $s_{x_i2}s_{y_ip}$ | | $s_{x_i3}s_{y_ip}$ | $\ldots$ | $s_{x_ip}s_{y_ip}$ | | $\mu_f s_{y_ip}$ |
| $\mu_g$ | $s_{x_i1}\mu_g$ | | $s_{x_i2}\mu_g$ | | $s_{x_i3}\mu_g$ | $\ldots$ | $s_{x_ip}\mu_g$ | | $\mu_f\mu_g$ |

|  | $s_{x_i1}$ + | $s_{x_i2}$ + | $s_{x_i3}$ + | ... + | $s_{x_ip}$ + | $\mu_f$ |
|---|---|---|---|---|---|---|
| $s_{y_i1}$ + | $s_{x_i1}s_{y_i1}$ | $s_{x_i2}s_{y_i1}$ | $s_{x_i3}s_{y_i1}$ | . . | $s_{x_ip}s_{y_i1}$ | $\mu_f s_{y_i1}$ |
| $s_{y_i2}$ + | $s_{x_i1}s_{y_i2}$ | $s_{x_i2}s_{y_i2}$ | $s_{x_i3}s_{y_i2}$ | . . | $s_{x_ip}s_{y_i2}$ | $\mu_f s_{y_i2}$ |
| $s_{y_i3}$ + | $s_{x_i1}s_{y_i3}$ | $s_{x_i2}s_{y_i3}$ | $s_{x_i3}s_{y_i3}$ | . . | $s_{x_ip}s_{y_i3}$ | $\mu_f s_{y_i3}$ |
| ⋮ + | . | . | . | . | . | . |
| $s_{y_in}$ + | $s_{x_i1}s_{y_ip}$ | $s_{x_i2}s_{y_ip}$ | $s_{x_i3}s_{y_ip}$ | . . | $s_{x_ip}s_{y_ip}$ | $\mu_f s_{y_ip}$ |
| $\mu_g$ | $s_{x_i1}\mu_g$ | $s_{x_i2}\mu_g$ | $s_{x_i3}\mu_g$ | . . | $s_{x_ip}\mu_g$ | $\mu_f\mu_g$ |

| | $s_{x_i 1}$ | + | $s_{x_i 2}$ | + | $s_{x_i 3}$ | + | ... | + | $s_{x_i p}$ | + | $\mu_f$ |
|---|---|---|---|---|---|---|---|---|---|---|---|
| $s_{y_i 1}$ + | $s_{x_i 1} s_{y_i 1}$ | | $s_{x_i 2} s_{y_i 1}$ | | $s_{x_i 3} s_{y_i 1}$ | | . | | $s_{x_i p} s_{y_i 1}$ | | $\mu_f s_{y_i 1}$ |
| $s_{y_i 2}$ + | $s_{x_i 1} s_{y_i 2}$ | | $s_{x_i 2} s_{y_i 2}$ | | $s_{x_i 3} s_{y_i 2}$ | | . | | $s_{x_i p} s_{y_i 2}$ | | $\mu_f s_{y_i 2}$ |
| $s_{y_i 3}$ + | $s_{x_i 1} s_{y_i 3}$ | | $s_{x_i 2} s_{y_i 3}$ | | $s_{x_i 3} s_{y_i 3}$ | | . | | $s_{x_i p} s_{y_i 3}$ | | $\mu_f s_{y_i 3}$ |
| $\vdots$ + | | | | | . | | . | | . | | |
| $s_{y_i n}$ + | $s_{x_i 1} s_{y_i p}$ | | $s_{x_i 2} s_{y_i p}$ | | $s_{x_i 3} s_{y_i p}$ | | . | | $s_{x_i p} s_{y_i p}$ | | $\mu_f s_{y_i p}$ |
| $\mu_g$ | $s_{x_i 1} \mu_g$ | | $s_{x_i 2} \mu_g$ | | $s_{x_i 3} \mu_g$ | | . | | $s_{x_i p} \mu_g$ | | $\mu_f \mu_g$ |

# EXPANSION OF TERMS: $\hat{x}_i \times \hat{y}_i$

| | $s_{x_i 1}$ | $+$ | $s_{x_i 2}$ | $+$ | $s_{x_i 3}$ | $+$ | $\ldots$ | $+$ | $s_{x_i p}$ | $+$ | $\mu_f$ |
|---|---|---|---|---|---|---|---|---|---|---|---|
| $s_{y_i 1}$ $+$ | $s_{x_i 1} s_{y_i 1}$ | | $s_{x_i 2} s_{y_i 1}$ | | $s_{x_i 3} s_{y_i 1}$ | | $\cdot$ | | $s_{x_i p} s_{y_i 1}$ | | $\mu_f s_{y_i 1}$ |
| $s_{y_i 2}$ $+$ | $s_{x_i 1} s_{y_i 2}$ | | $s_{x_i 2} s_{y_i 2}$ | | $s_{x_i 3} s_{y_i 2}$ | | $\cdot$ | | $s_{x_i p} s_{y_i 2}$ | | $\mu_f s_{y_i 2}$ |
| $s_{y_i 3}$ $+$ | $s_{x_i 1} s_{y_i 3}$ | | $s_{x_i 2} s_{y_i 3}$ | | $s_{x_i 3} s_{y_i 3}$ | | $\cdot$ | | $s_{x_i p} s_{y_i 3}$ | | $\mu_f s_{y_i 3}$ |
| $\vdots$ $+$ | | | $\cdot$ | | | | $\cdot$ | | | | $\cdot$ |
| $s_{y_i n}$ $+$ | $s_{x_i 1} s_{y_i p}$ | | $s_{x_i 2} s_{y_i p}$ | | $s_{x_i 3} s_{y_i p}$ | | $\cdot$ | | $s_{x_i p} s_{y_i p}$ | | $\mu_f s_{y_i p}$ |
| $\mu_g$ | $s_{x_i 1} \mu_g$ | | $s_{x_i 2} \mu_g$ | | $s_{x_i 3} \mu_g$ | | $\cdot$ | | $s_{x_i p} \mu_g$ | | $\mu_f \mu_g$ |

## Proposed Approach

Modified contributions from all variables

The mean prediction of the two-part model

$$\hat{z}_i = \left( \sum_{j=1}^{p} s'_{z_{ij}} \right) + \alpha + \mu_h$$

The output of the two-part model $(\hat{x}_i \cdot \hat{y}_i)$

The difference between $\mu_f \mu_g$ and $\mu_h$

where

$$s'_{z_{ij}} = \mu_f s_{y_{ij}} + s_{x_{ij}} \mu_g + \frac{1}{2} \sum_{a=1}^{p} (s_{x_{ij}} s_{y_{ia}} + s_{y_{ij}} s_{x_{ia}})$$

## Distributing $\alpha$

Uniformly Distributed:

$$s_{z_i j} = s'_{z_i j} + \frac{\alpha}{p}.$$

Raw Weights:

$$s_{z_i j} = s'_{z_i j} + \frac{s'_{z_i j}}{\hat{z}_i - \mu_f \mu_g}(\alpha).$$

Absolute Weights:

$$s_{z_i j} = s'_{z_i j} + \frac{|s'_{z_i j}|}{\sum_{k=1}^{p} |s'_{z_i k}|}(\alpha).$$

Squared Weights:

$$s_{z_i j} = s'_{z_i j} + \frac{(s'_{z_i j})^2}{\sum_{k=1}^{p}(s'_{z_i k})^2}(\alpha).$$

# SIMULATION STUDY

## Simulation Results

| Method | Score | Pct Same Sign | Pct Same Rank |
|---|---|---|---|
| Weighted by Absolute Value | **2.27** | **84.8%** | **62.5%** |
| Weighted by Squared Value | 2.21 | 81.8% | 60.8% |
| Uniformly Distributed | 2.20 | 83.7% | 59.4% |
| Weighted by Raw Value | 1.99 | 71.4% | 56.2% |

## Final mSHAP Equation

Thus, the final equation for the mSHAP value of the $j$th predictor on the $i$th observation can be written as

$$s_{z_ij} = \mu_f s_{y_ij} + s_{x_ij}\mu_g + \frac{1}{2}\left[\sum_{a=1}^{p}(s_{x_ij}s_{y_ia} + s_{y_ij}s_{x_ia})\right] + \frac{|s'_{z_ij}|}{\sum_{k=1}^{p}|s'_{z_ik}|}(\alpha).$$

And the overall prediction is

$$\hat{z_i} = \mu_h + \sum_{j=1}^{p} s_{z_ij}$$

This is implemented in the R package {mshap}, which is available on CRAN and on github at www.github.com/srmatth/mshap

Comparison of mSHAP and TreeSHAP
Score when Computed Against kernelSHAP with Similar Response Transformations

# Comparison to kernelSHAP

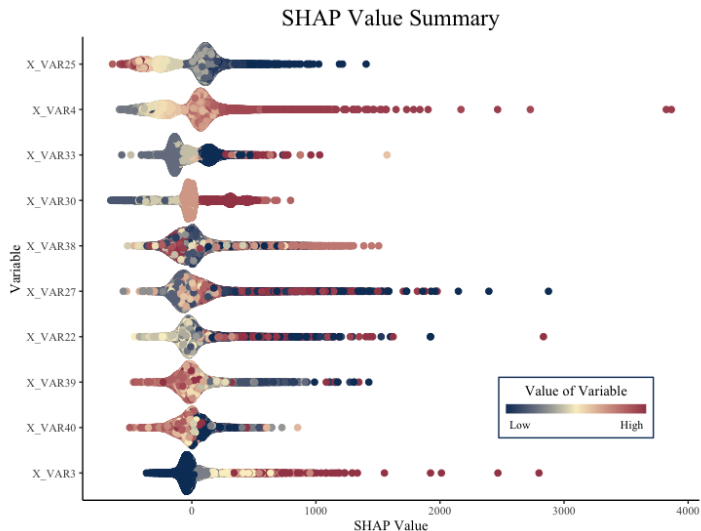A dramatic increase in speed and computational efficiency:



Practical Example:

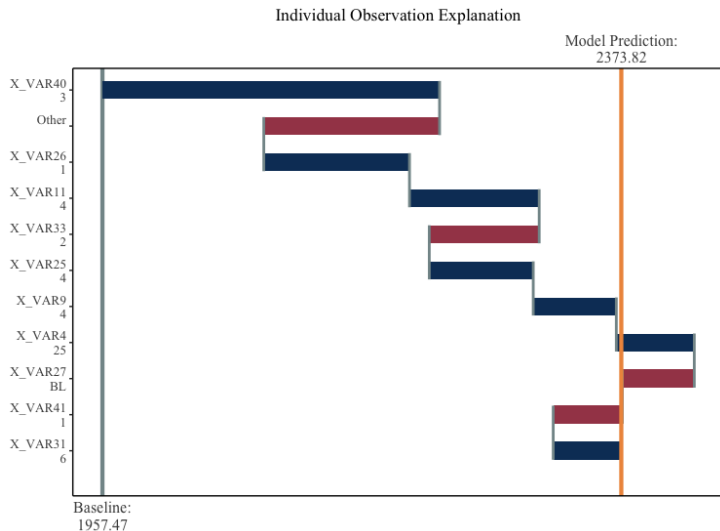- 5,000,000 Rows with 45 Covariates
- 131 Days vs. 3 Hours

## APPLICATION

- Obtained a property damage insurance data set which we then cleaned and split it into train, validation, and test sets (R)
- Trained a two-part model where both parts were random forests and the ultimate response of the model was the expected cost of a policy (Python)
- Computed the SHAP values for each individual model part using TreeSHAP (Python)
- Computed and visualized the contributions to the expected cost of a policy using mSHAP (R)

# APPLICATION



SHAP Value Summary

# APPLICATION



Individual Observation Explanation

## Conclusion

- kernelSHAP is unable to feasibly explain model predictions for two-part models
- mSHAP provides a framework for obtaining model explanations for two-part models, using the SHAP values of the individual model parts
- mSHAP will allow two-part models made up of tree-based models to be used in regulated industries such as insurance

## Acknowledgments

This work was supported by:

The paper is available on arxiv.org and srmatth.github.io
The code is available at www.github.com/srmatth/mshap

All plots were created with the {mshap} R package, which is available on CRAN and at www.github.com/srmatth/mshap

# mSHAP
### Explaining Two-Part Models

Spencer Matthews

Joint work with Brian Hartman

June 2021