

Ideas on Machine Learning Interpretability

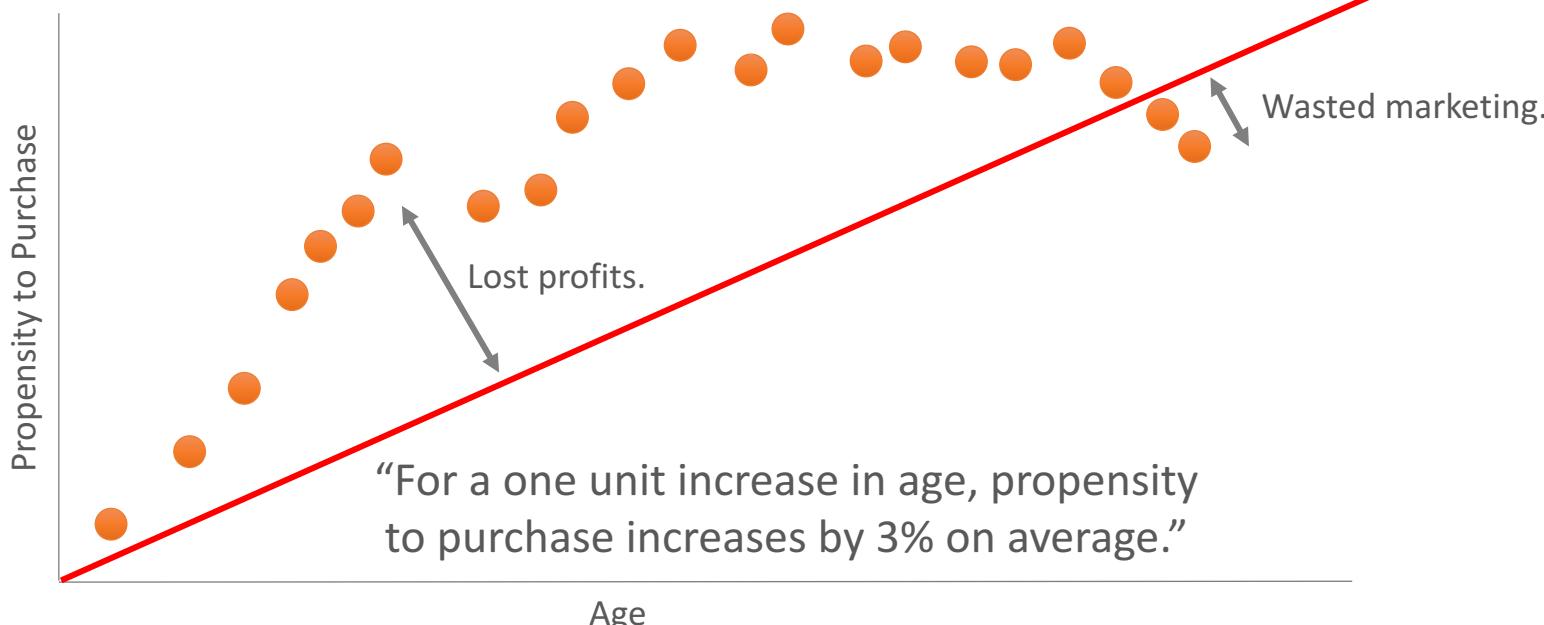
Patrick Hall, Wen Phan, SriSatish Ambati and the H2O.ai team

February 2017

Big Ideas

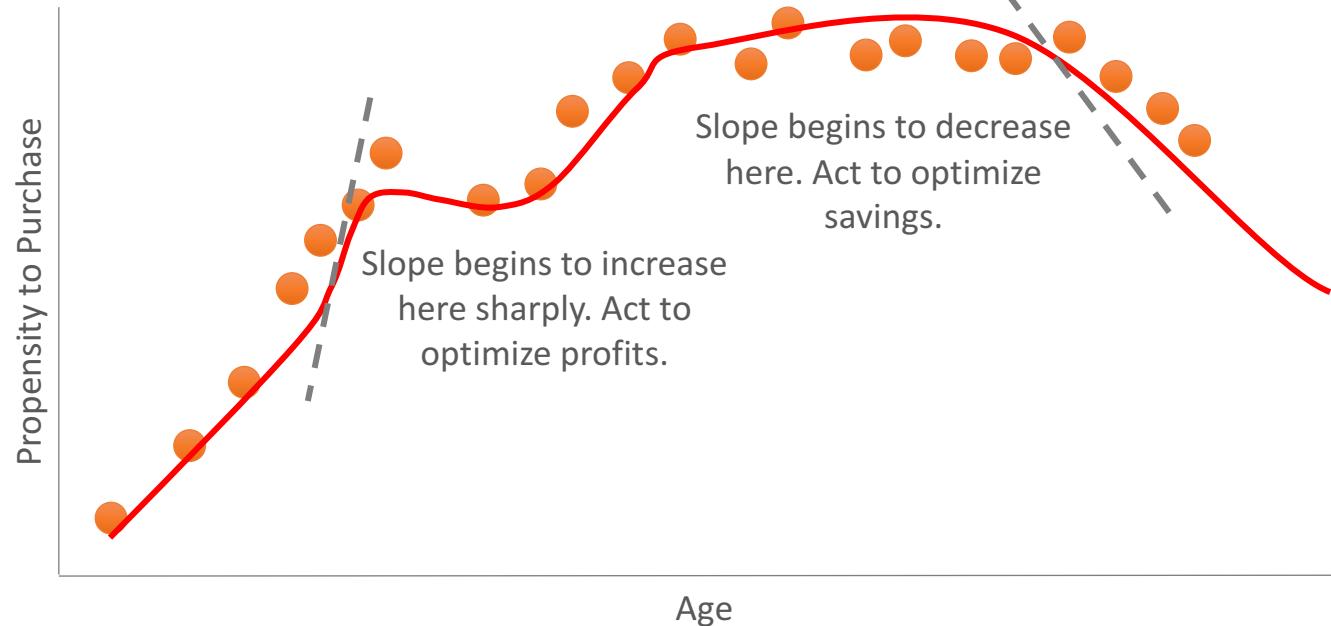
Linear Models

Exact explanations for
approximate models.



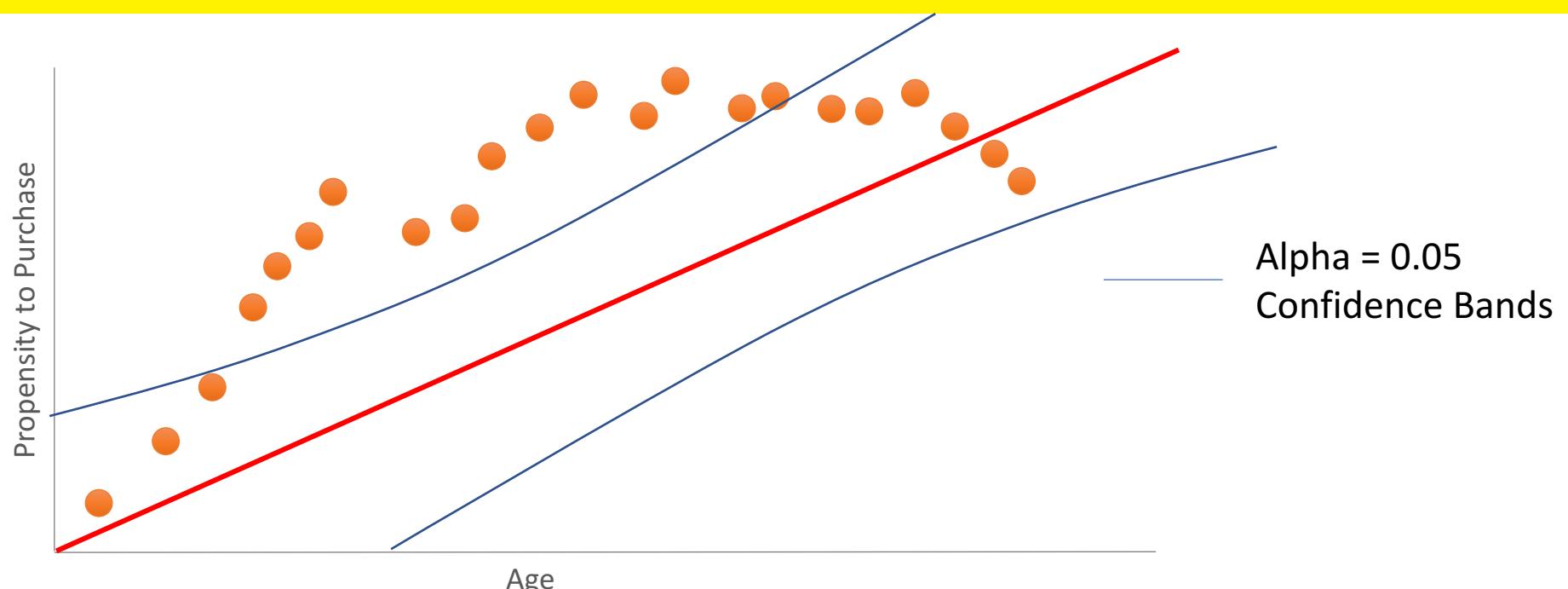
Machine Learning

Approximate explanations
for ***exact*** models.

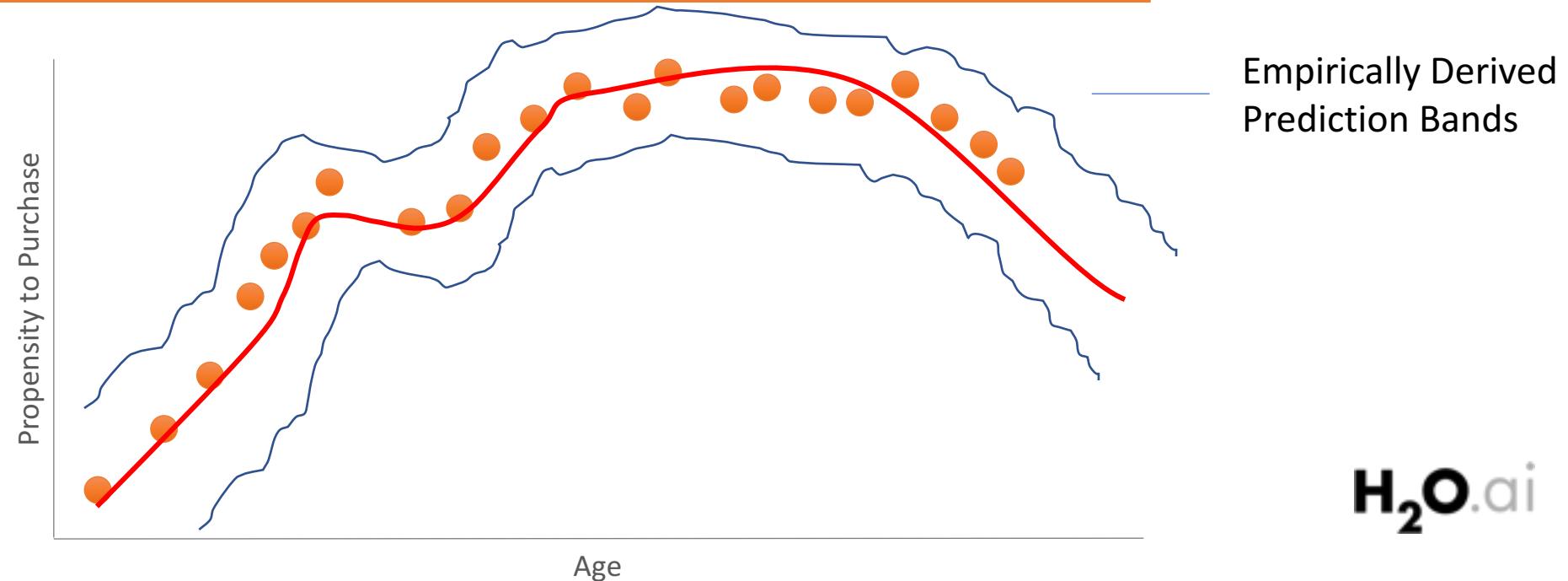


Linear Models

Risk is well defined ...
Theoretically ...
Based on strong
assumptions.

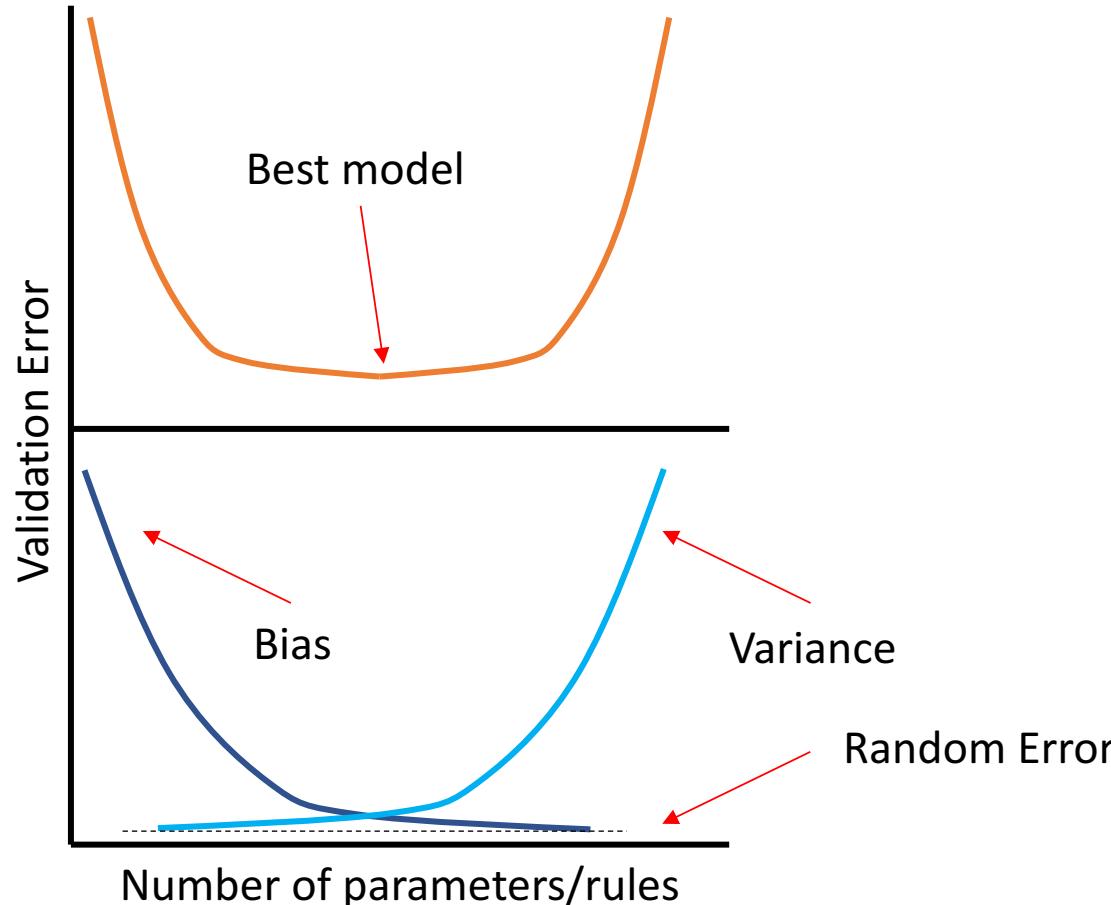


Machine Learning
Risk is empirically
quantifiable ...
But it's hard work.



Risk from Unwanted Bias and Prediction Variance

$$\text{Total Error} = \text{Bias} + \text{Variance} + \text{Random Error} = (\hat{f}(x) - f(x))^2$$

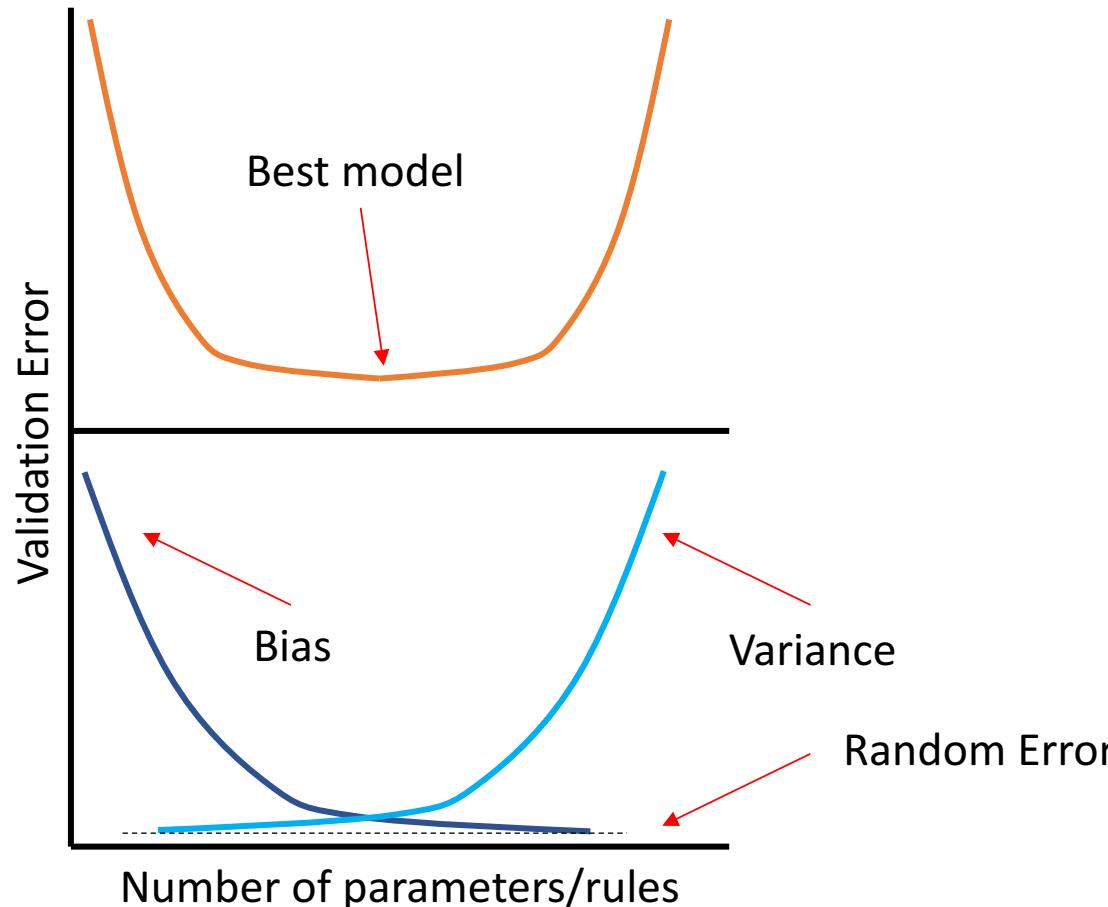


Bias = $E[\hat{f}(x)] - f(x)$ or the error that arises from a model's inability to replicate the fundamental phenomena represented by a data set.

Variance = $(\hat{f}(x) - E[\hat{f}(x)])^2$ or the error that arises from a model's ability to produce differing predictions from the values in a new data set.

Risk from Unwanted Bias and Prediction Variance

$$\text{Total Error} = \text{Bias} + \text{Variance} + \text{Random Error} = (\hat{f}(x) - f(x))^2$$



Risk from Unwanted Bias: Your model includes contributions from race, gender, disability status, marital status, or other unwanted latent features.

Risk from Prediction Variance: Your model is unpredictable outside of the training domain.

A framework for interpretability

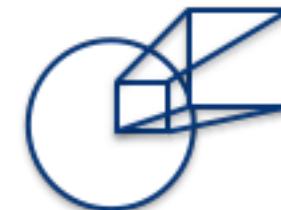
Complexity of learned functions:

- Linear, monotonic
- Nonlinear, monotonic
- Nonlinear, non-monotonic



Scope of interpretability:

Global vs. local



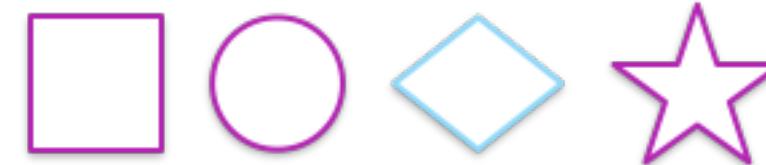
Enhancing trust and understanding:

the mechanisms and results of an interpretable model should be both transparent AND dependable.



Application domain:

Model-agnostic vs. model-specific



Contents

Part 1: Seeing your data

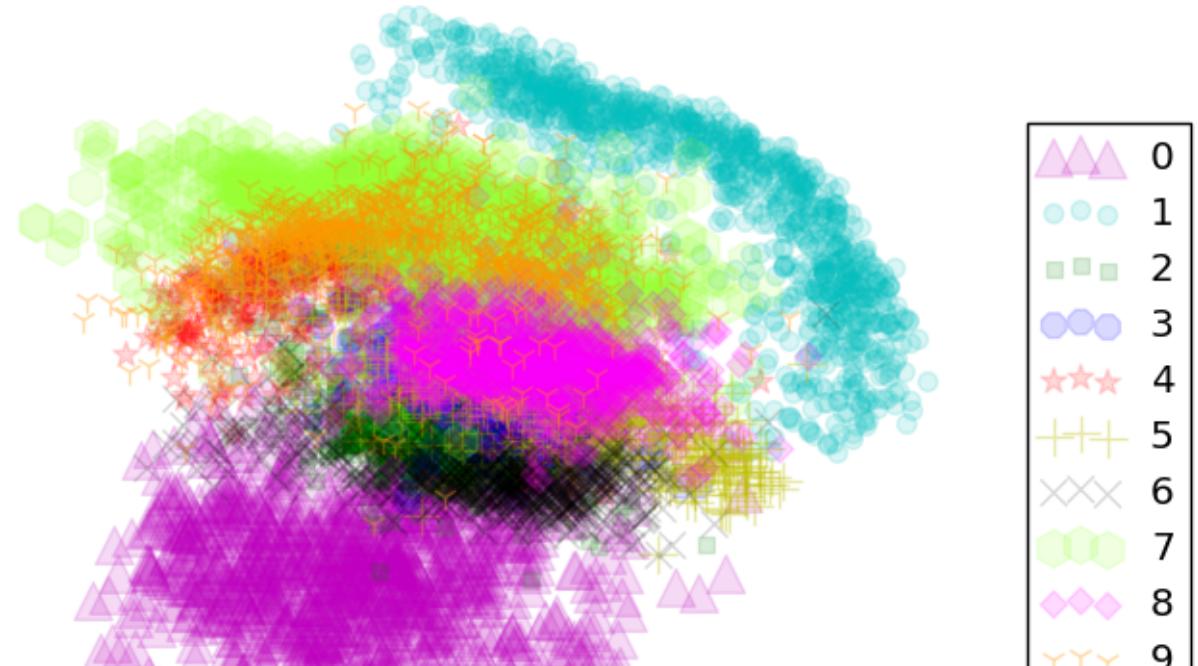
- Glyphs
- Correlation graphs
- 2-D projections
- Partial dependence plots
- Residual analysis

Part 2: Using machine learning in regulated industry

- OLS regression alternatives
- Build toward ML model benchmarks
- ML in traditional analytics processes
- Small, interpretable ensembles
- Monotonicity constraints
- Rule-Based Models

Part 3: Understanding complex ML models

- Surrogate models
- LIME
- Maximum activation analysis
- Sensitivity analysis
- Variable importance measures
- TreeInterpreter



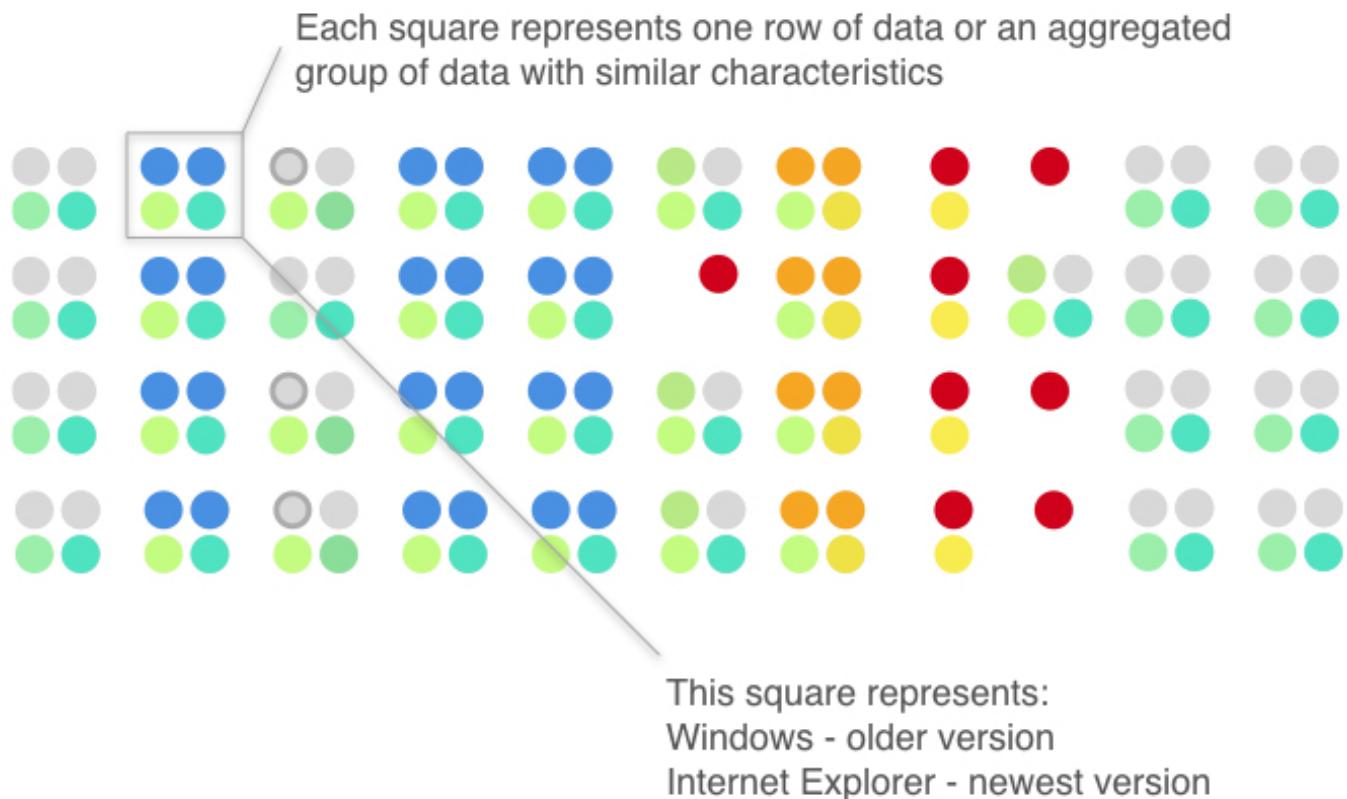
Part 1: Seeing your data

Glyphs

Variables and their values can be represented by small pictures with certain attributes, called “glyphs”.

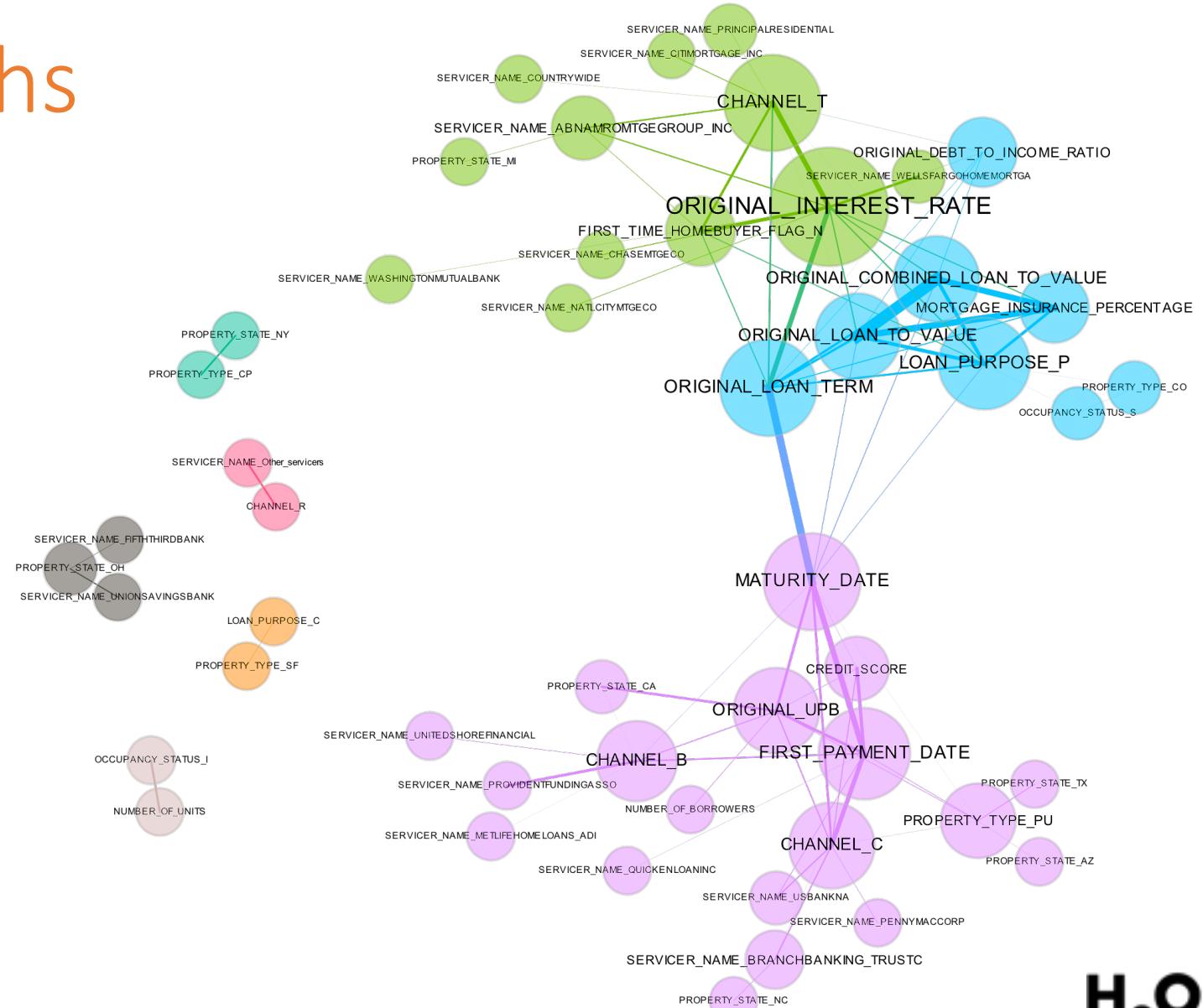
Here the four variables are represented by their position in a square and their values are represented by a color.

os_type					
	iphone	OSX	Win	Android	Linux
os_version					oldest
	newest				
agent_type					
	Opera	Safari	IE	Others	Firefox
agent_version					oldest
	newest				

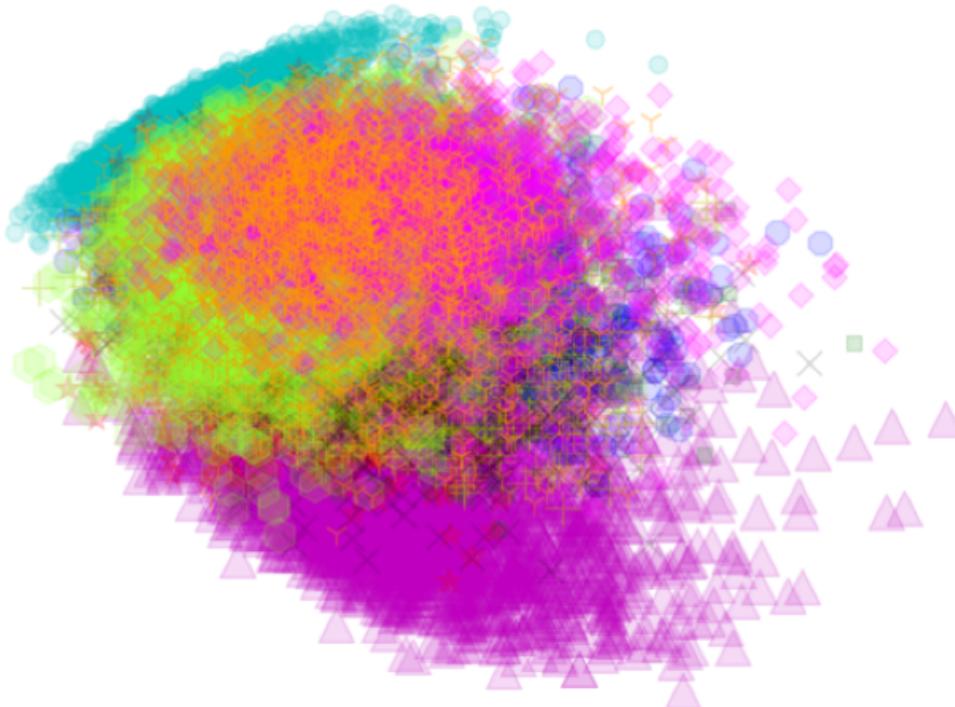


Correlation graphs

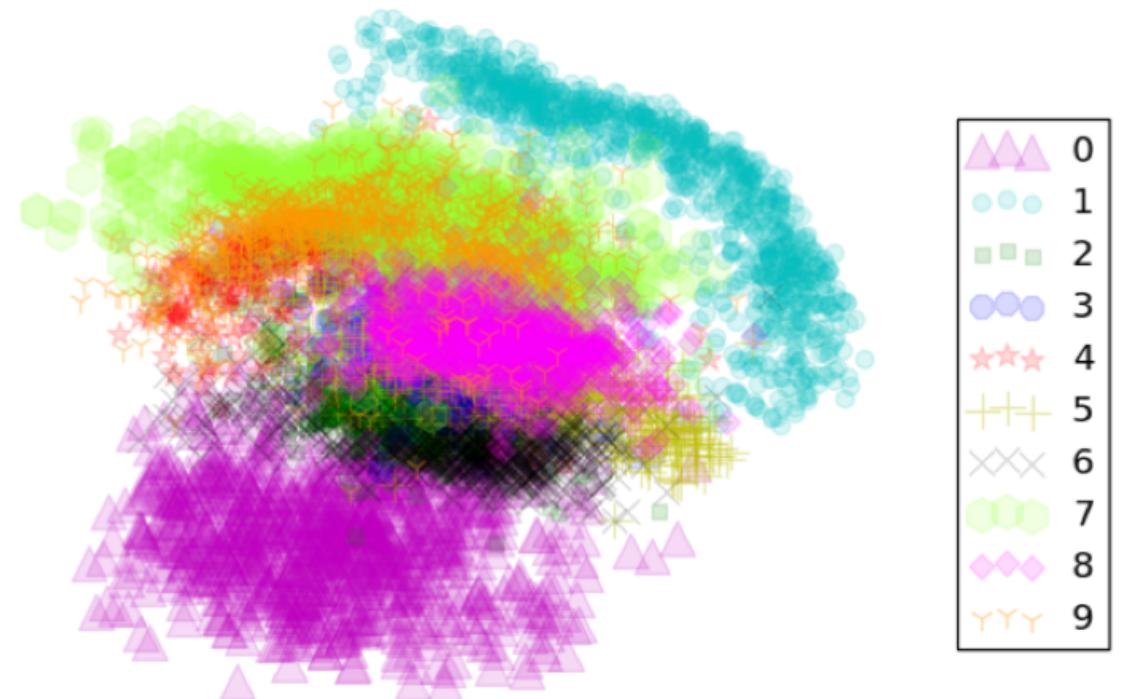
The nodes of this graph are the variables in a data set. The weights between the nodes are defined by the absolute value of their pairwise Pearson correlation.



2-D projections

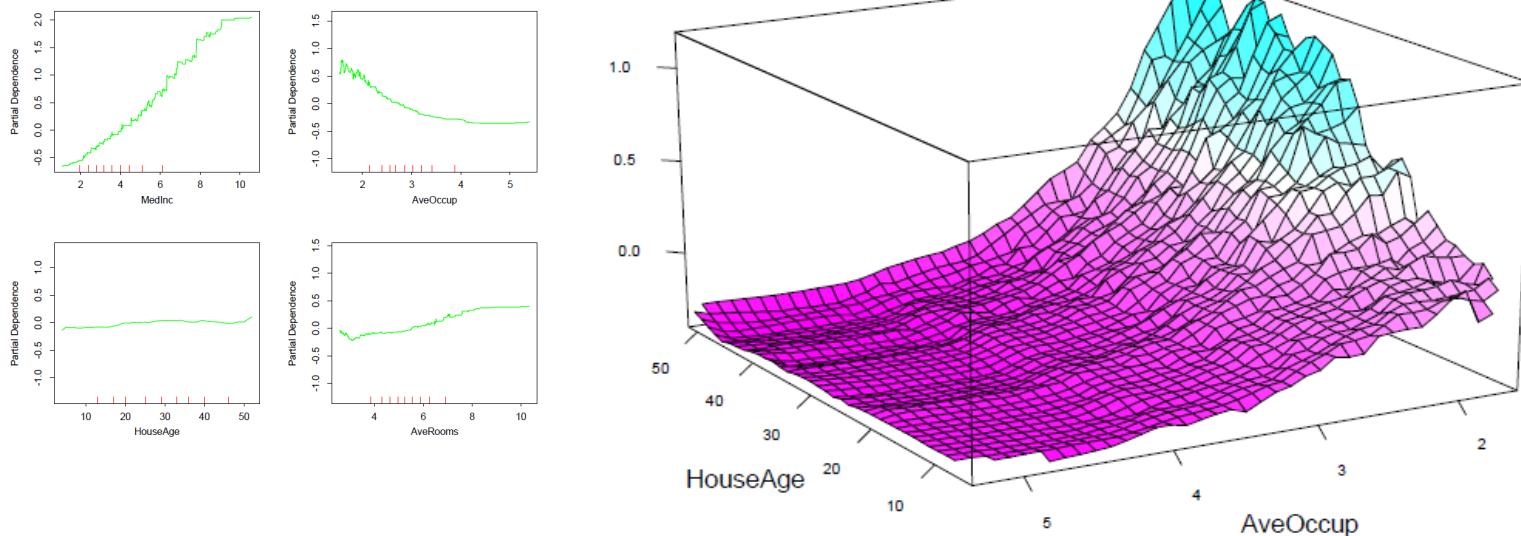


784 dimensions to 2 dimensions with PCA



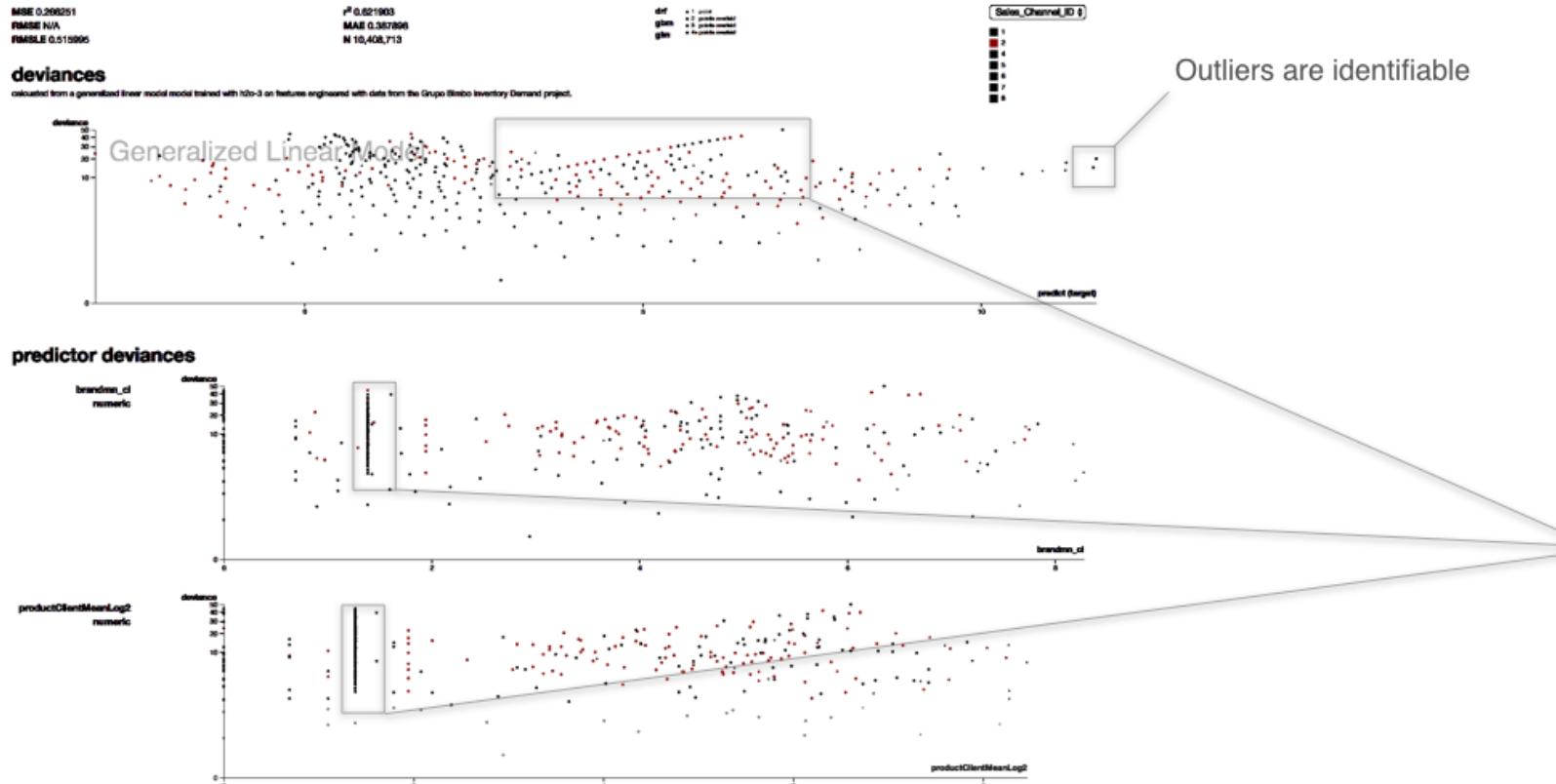
784 dimensions to 2 dimensions with
autoencoder network

Partial dependence plots



HomeValue ~ MedInc + AveOccup + HouseAge + AveRooms

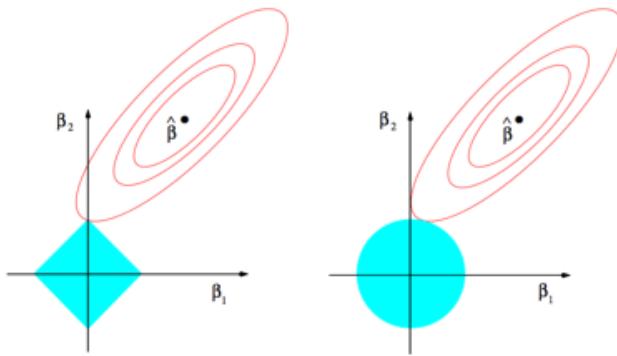
Residual analysis



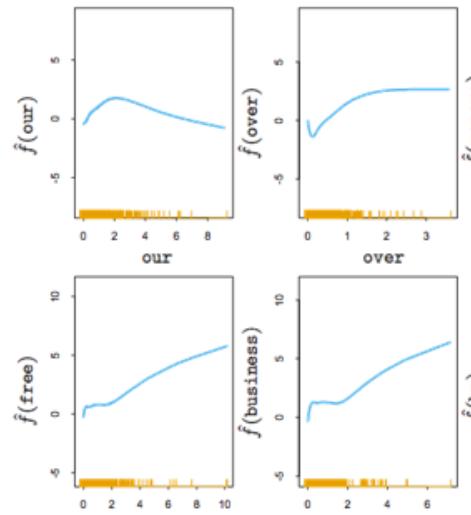
Residuals from a machine learning model should be randomly distributed
obvious patterns in residuals can indicate problems with data preparation or model specification

Part 2: Using ML in regulated industry

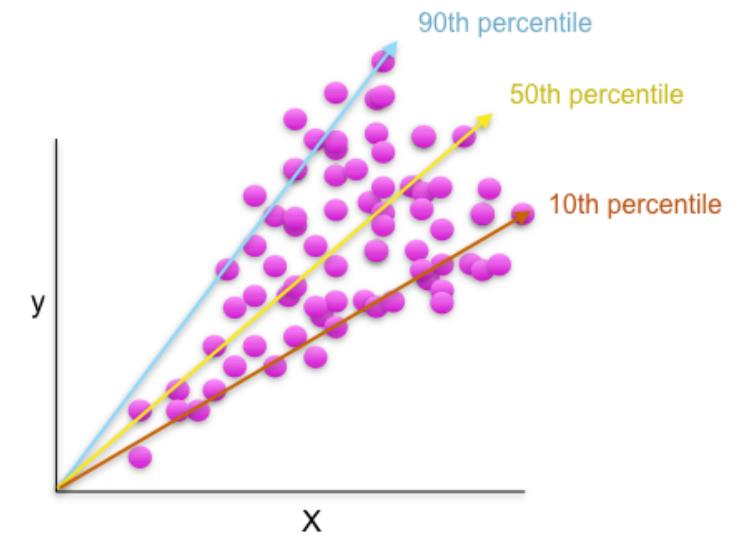
OLS regression alternatives



Penalized Regression



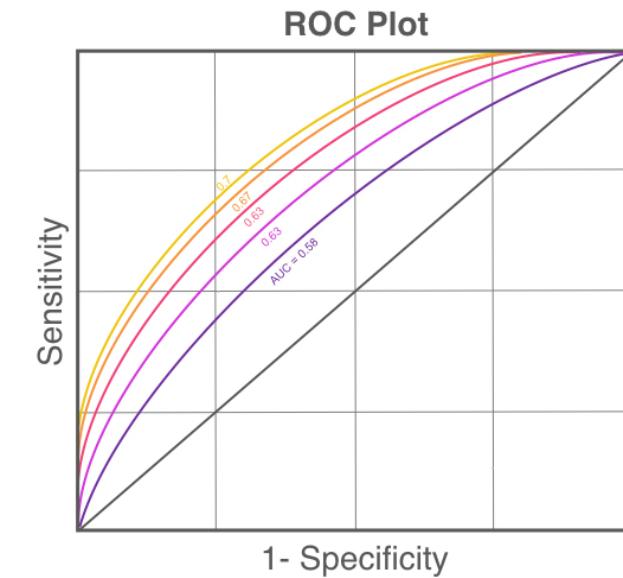
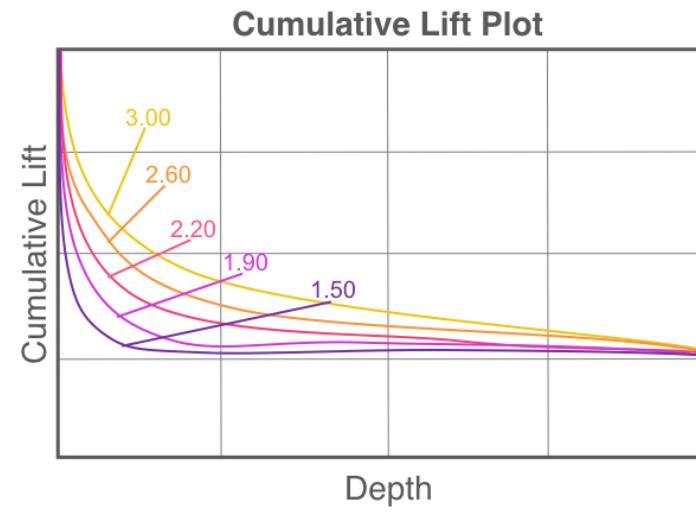
Generalized Additive Models



Quantile Regression

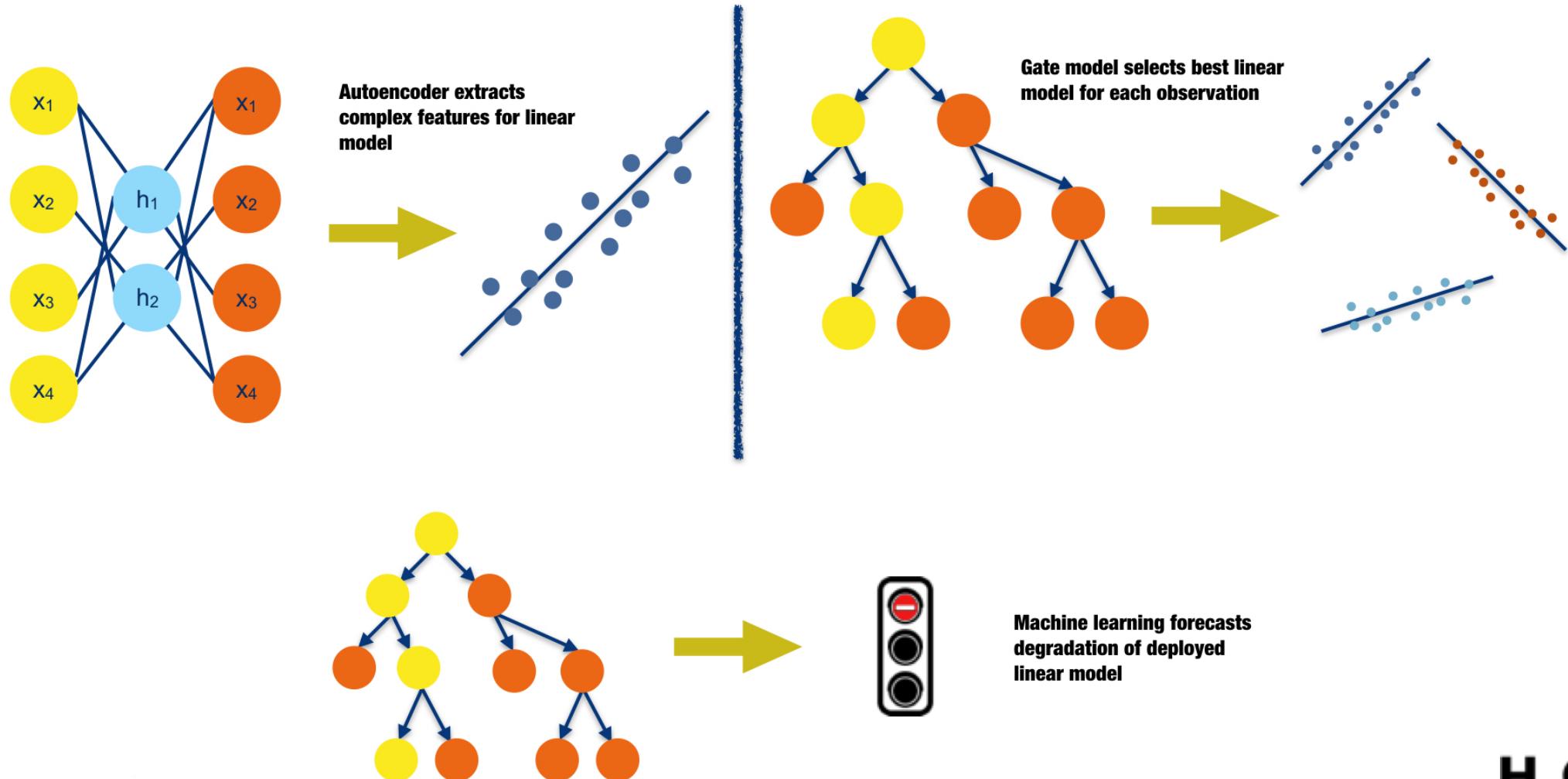
Build toward ML model benchmarks

Gradient Boosting	—
Neural Network	—
$y = x_1 + x_2 + x_3 + x_1 \cdot x_3 + x_2 \cdot x_3$	—
$y = x_1 + x_2 + x_3 + x_2 \cdot x_3$	—
$y = x_1 + x_2 + x_3$	—

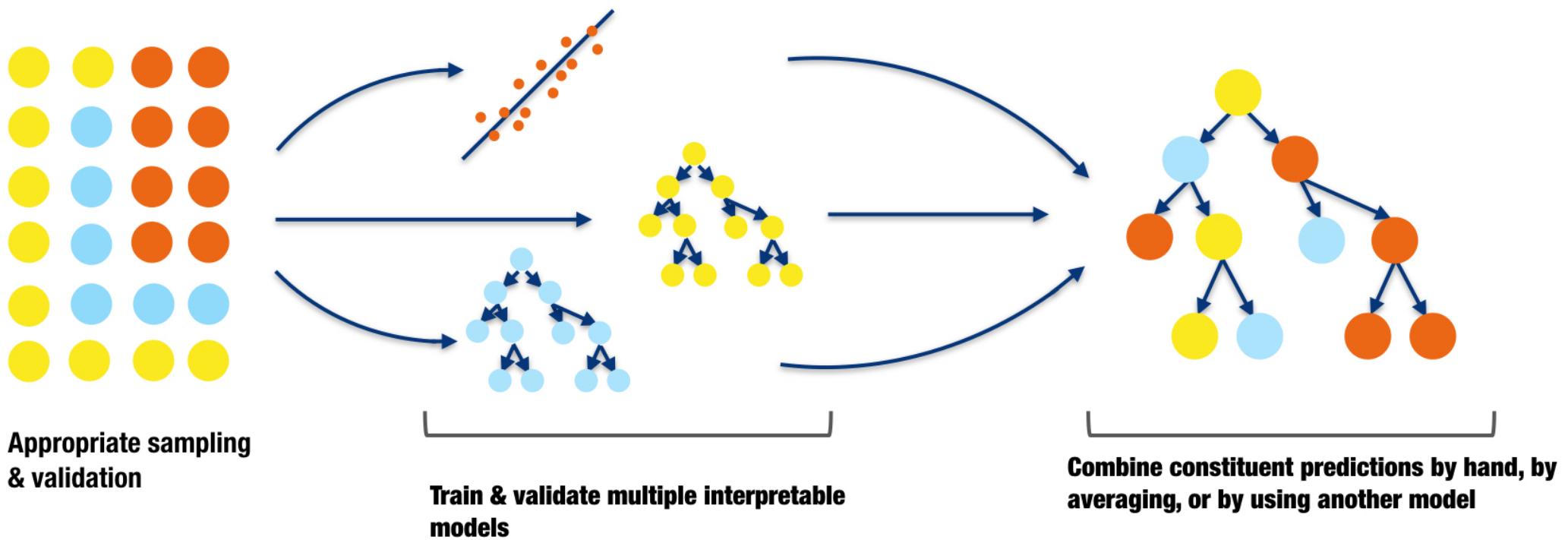


Incorporate interactions and piecewise linear components to increase the accuracy of linear models relative to machine learning models.

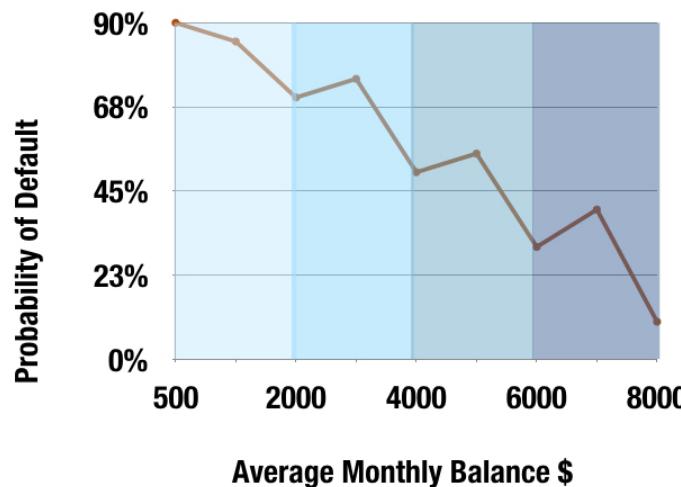
Machine learning in traditional analytics processes



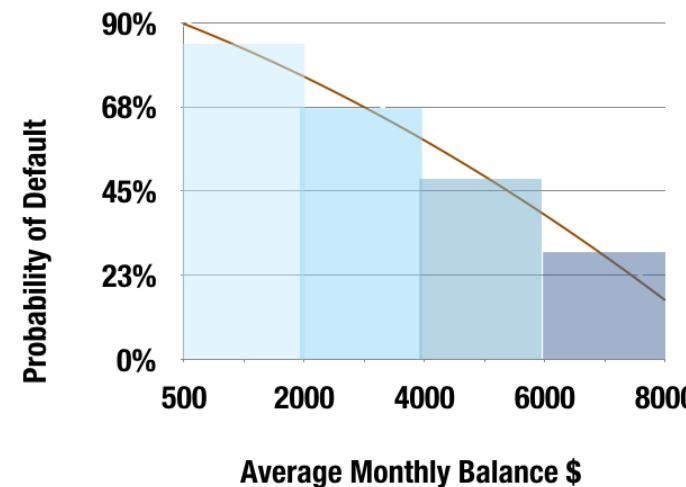
Small, interpretable ensembles



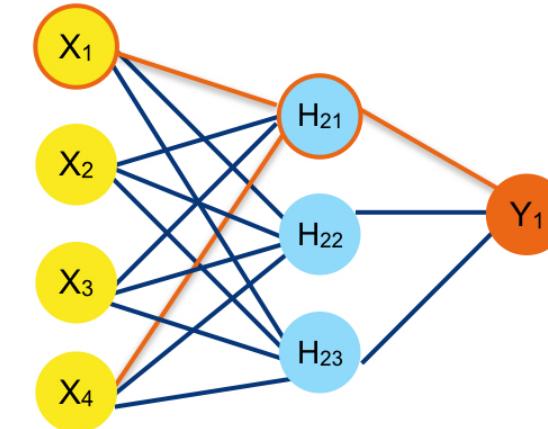
Monotonicity constraints



Average Monthly Balance is a nonnegative quantity, but is not monotonic with respect to Probability of Default.

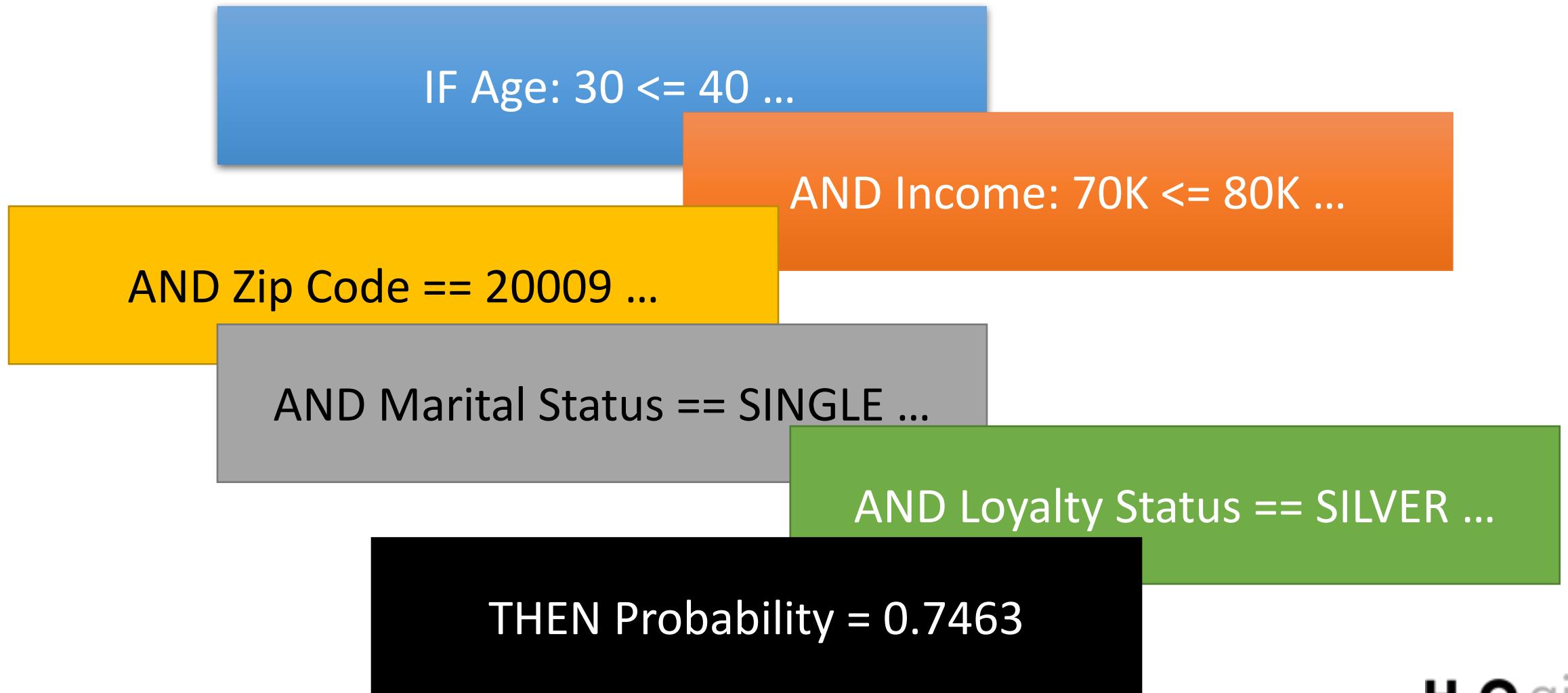


By discretization, the Average Monthly Balance can be transformed to be monotonic with respect to the target.



When all inputs are nonnegative and monotonic with respect to the target, and model weights are constrained to be nonnegative, it's easier understand the impact of individual features and to find interactions.

Rule-based models

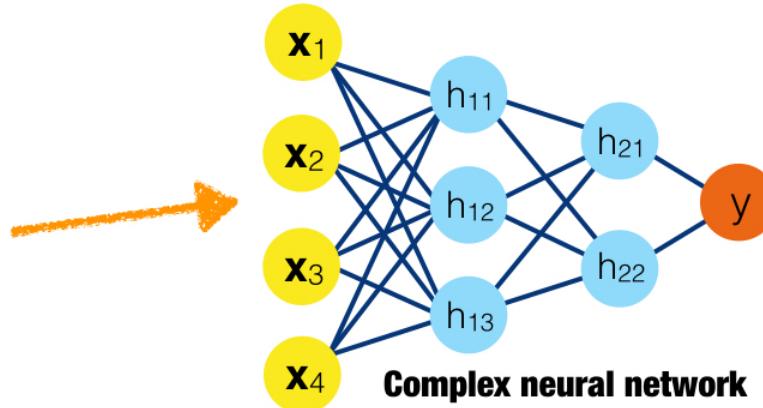


Part 3: Understanding complex machine learning models

Surrogate models

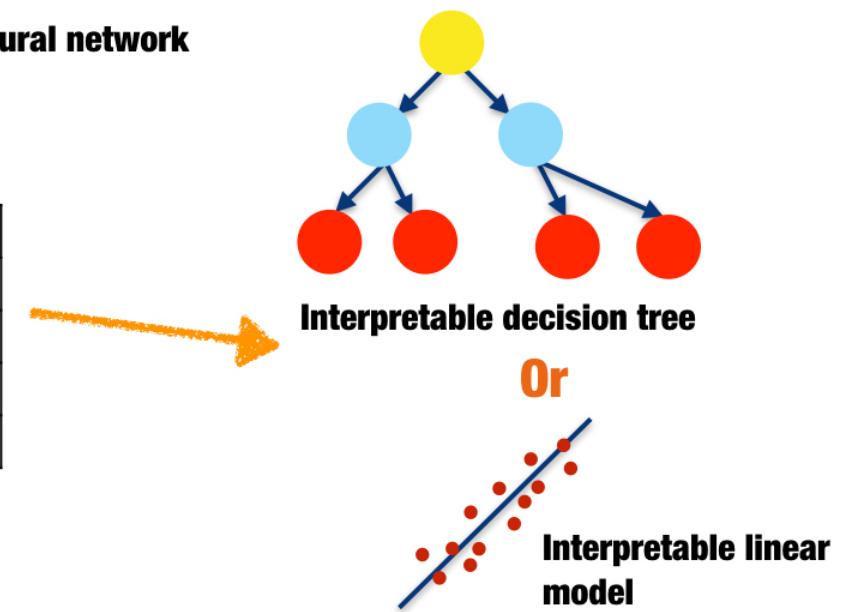
BAD	CUSTOMER_DTI	LOAN_PURPOSE	CHANNEL
0	0.18	MORT	7
1	0.42	HELOC	10
0	0.11	MORT	10
0	0.21	MORT	1

1. Train a complex machine learning model

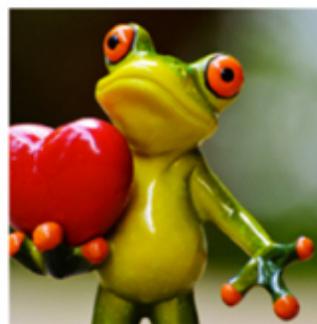


BAD	PREDICTED_BAD	CUSTOMER_DTI	LOAN_PURPOSE	CHANNEL
0	0.47	0.18	MORT	7
1	0.82	0.42	HELOC	10
0	0.18	0.11	MORT	10
0	0.12	0.21	MORT	1

2. Train an interpretable model on the original inputs and the predicted target values of the complex model



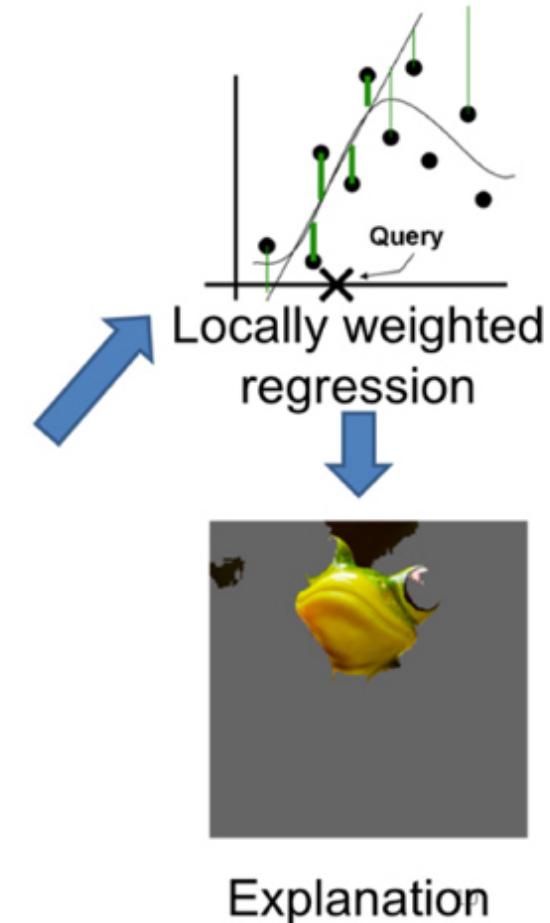
Local interpretable model-agnostic explanations



Original Image
 $P(\text{tree frog}) = 0.54$

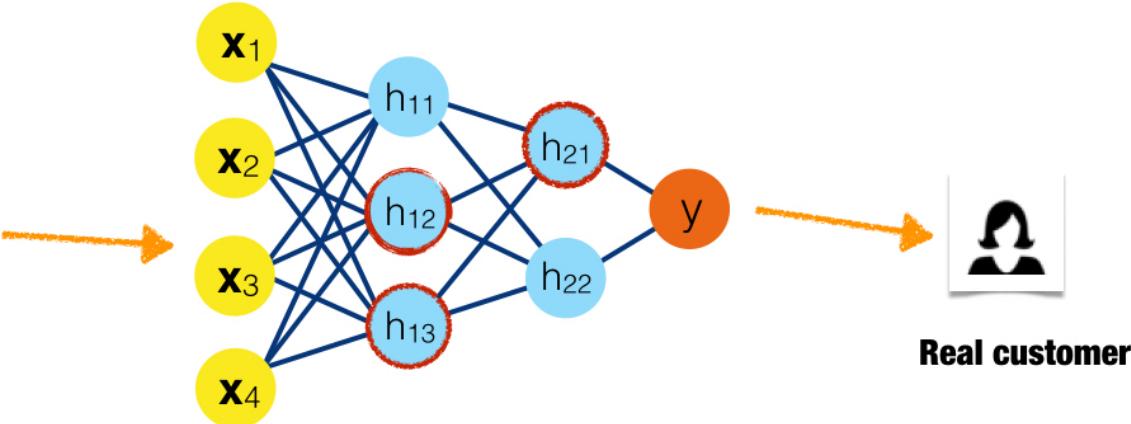


Perturbed Instances	$P(\text{tree frog})$
A photograph of a tree frog with several small red spots added to its body.	0.85
A photograph of a tree frog with several small yellow spots added to its body.	0.00001
A photograph of a tree frog with red flowers added to its body.	0.52



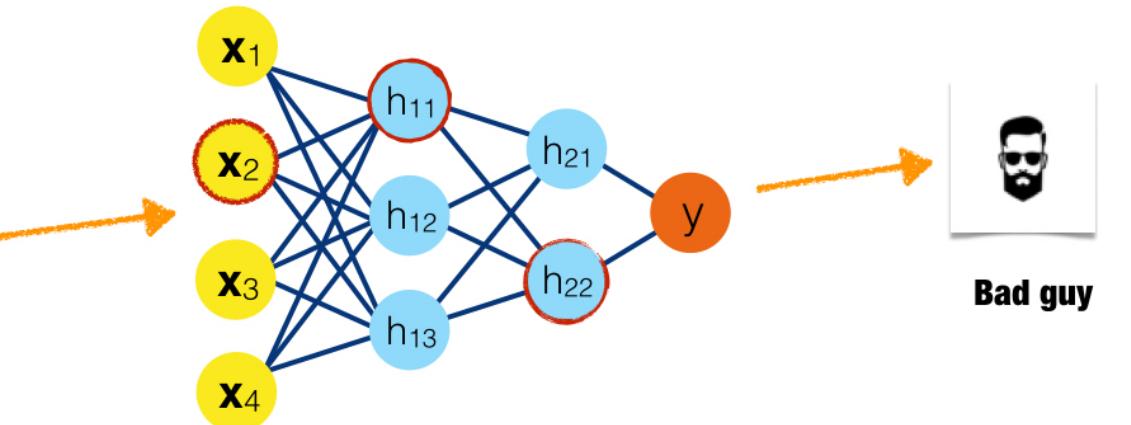
Maximum activation analysis

Name	Address	SSN	Phone
VICTOR HOLT	3456 7th St.	101-112-1314	(555)555-5555
GLORIANE BELL	123 Main St.	123-45-7689	(444)444-4444
VICTOR HOLMES	124 Main St.	910-11-1213	(333)333-3333
GLORY ANN BELL	123 Main St.	131-41-5167	(444)434-4334
GLORIANE BELL	7113 Third St.	123-45-7689	(444)444-4444
HECTOR HOLT	3456 7th St.	101-112-1331	(777)777-7777

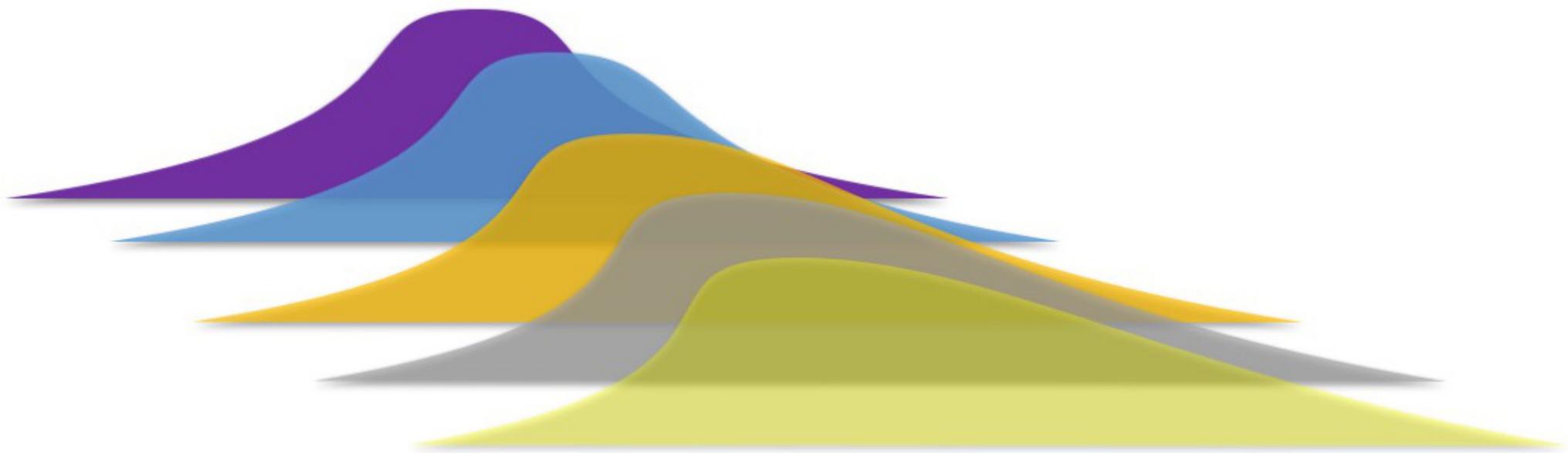


Different types of customers maximally activate different internal structures

Name	Address	SSN	Phone
VICTOR HOLT	3456 7th St.	101-112-1314	(555)555-5555
GLORIANE BELL	123 Main St.	123-45-7689	(444)444-4444
VICTOR HOLMES	124 Main St.	910-11-1213	(333)333-3333
GLORY ANN BELL	123 Main St.	131-41-5167	(444)434-4334
GLORIANE BELL	7113 Third St.	123-45-7689	(444)444-4444
HECTOR HOLT	3456 7th St.	101-112-1331	(777)777-7777



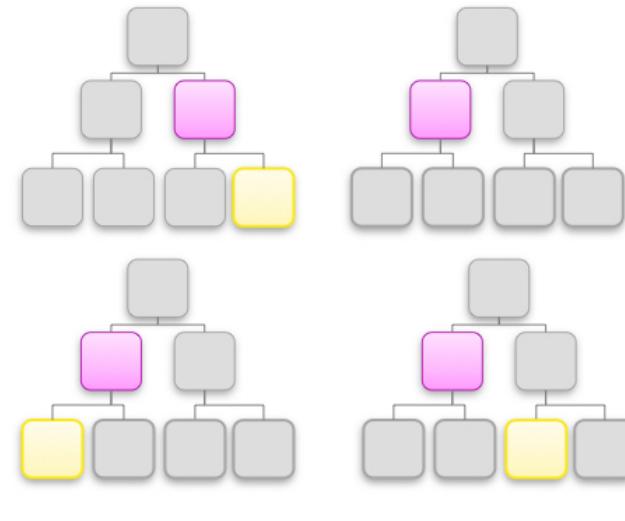
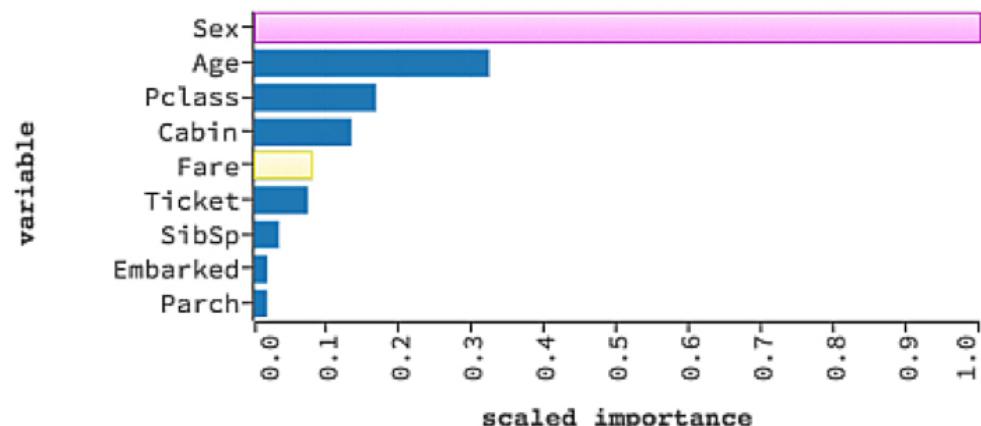
Sensitivity analysis



Data distributions shift over time. How will your model handle these shifts?

Variable importance measures

▼ VARIABLE IMPORTANCES



Global variable importance indicates the impact of a variable on the model for the entire training data set.

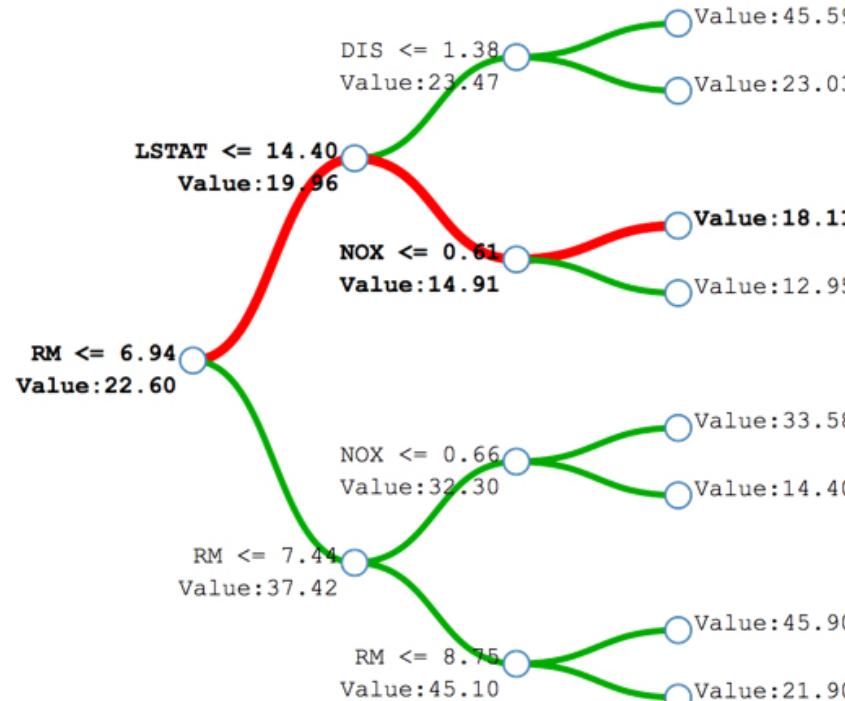
Sex	Age	...	Fare	\hat{y}	$\hat{y}_{(-\text{Sex})}$	$\hat{y}_{(-\text{Age})}$...	$\hat{y}_{(-\text{Fare})}$
M	11	...	8.45	0.2	0.01	0.1	...	0.21
F	34	...	51.86	0.8	0.6	0.65	...	0.78
M	26	...	21.08	0.5	0.2	0.3	...	0.53
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮

Local variable importance can indicate the impact of a variable for each decision a model makes – similar to reason codes.

Treeinterpreter

Tree interpreter decomposes decision tree and random forest predictions into bias (overall average) and component terms.

This slide portrays the decomposition of the decision path into bias and individual contributions for a single decision tree - treeinterpreter simply prints a ranked list of the bias and individual contributions for a given prediction.



Prediction: 18.11 ≈ 22.60 (trainset mean) - 2.64(loss from RM) - 5.04(loss from LSTAT) + 3.20(gain from NOX)