

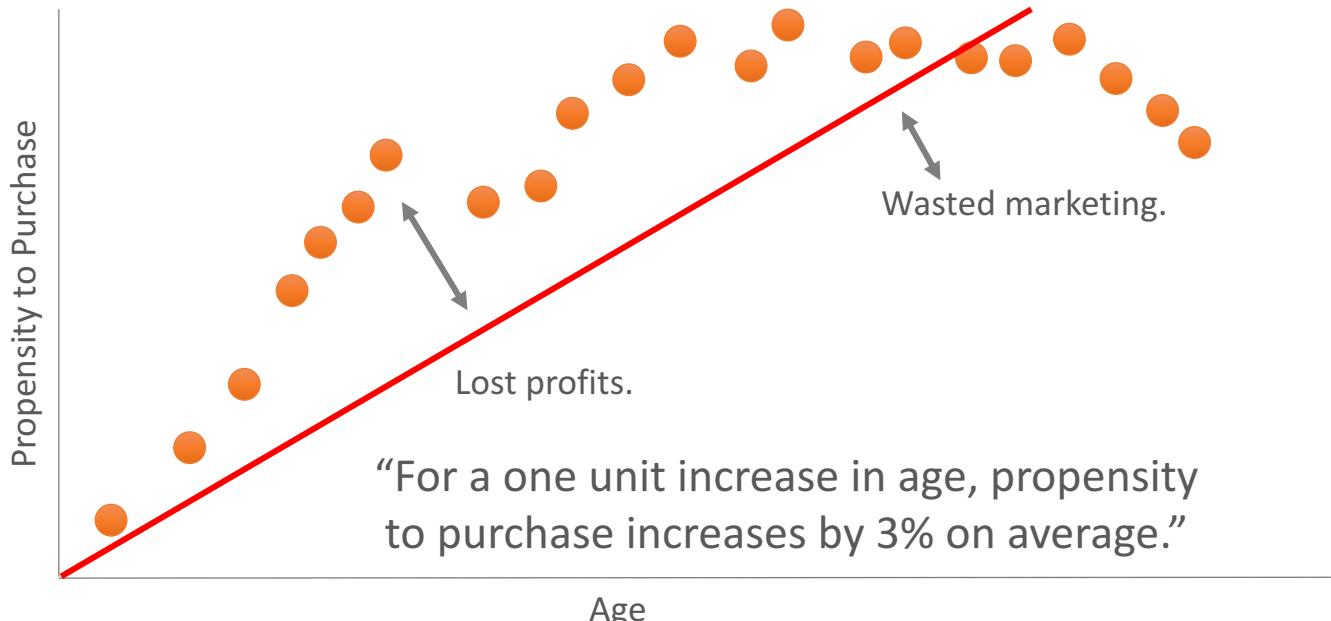
Ideas on Machine Learning Interpretability

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

Big Ideas

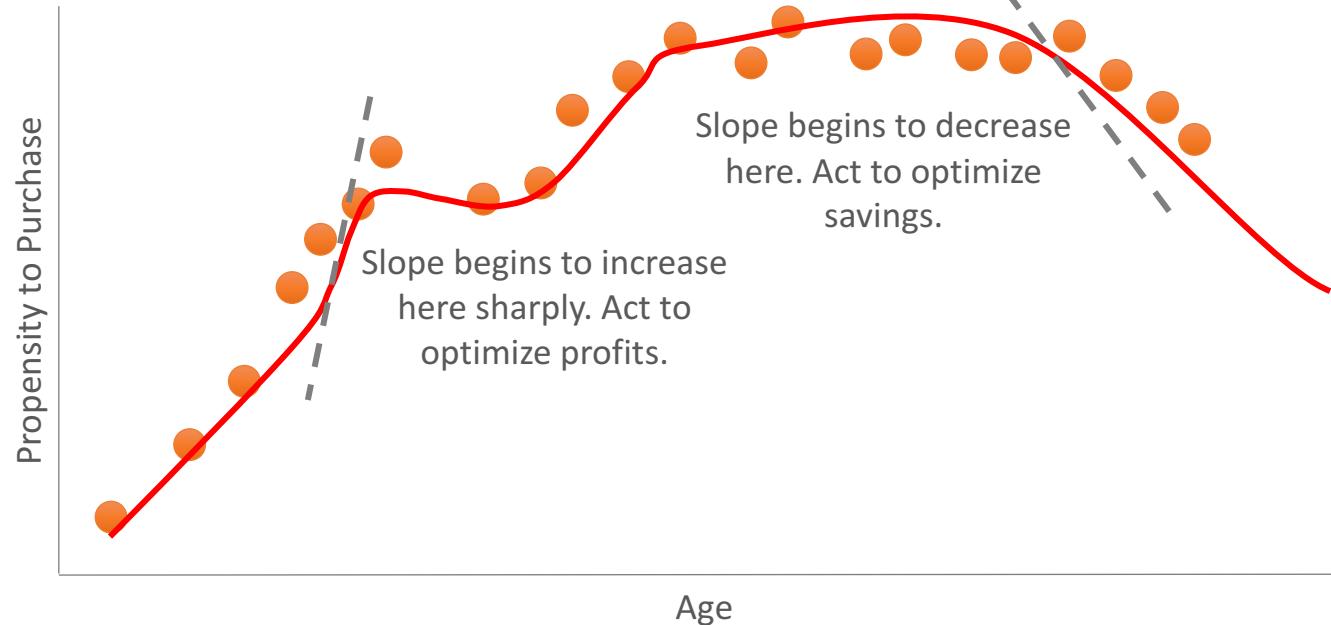
Linear Models

Exact explanations for
approximate models.



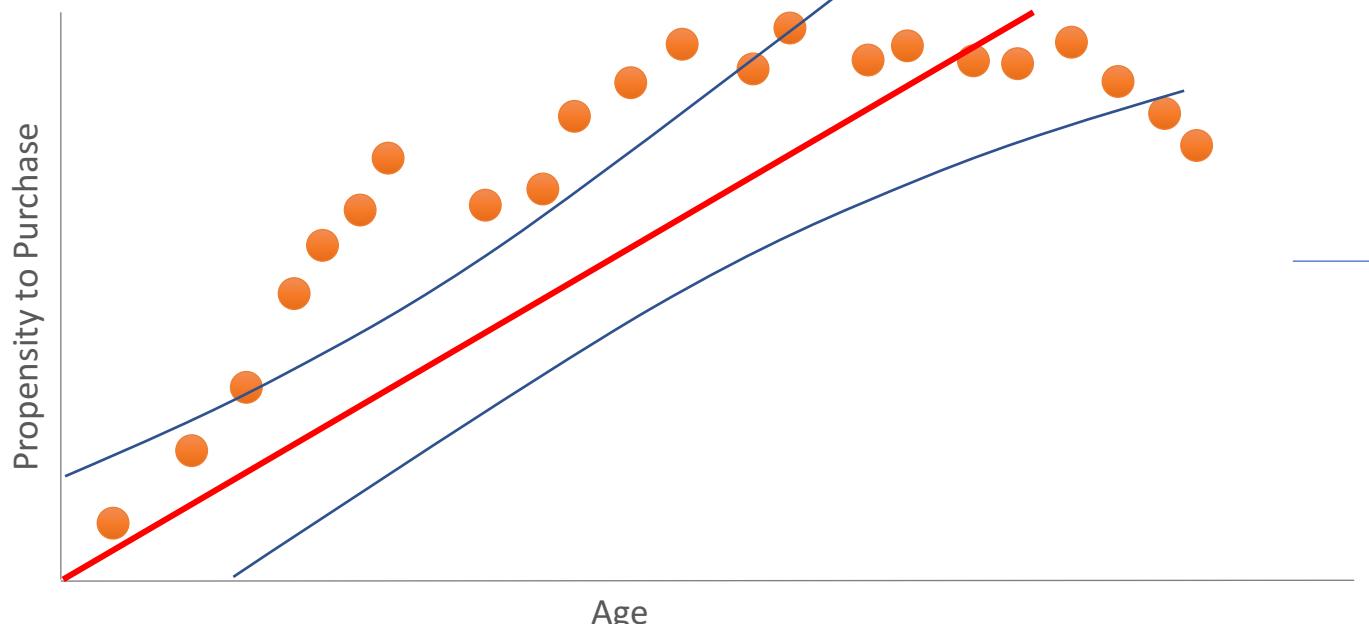
Machine Learning

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

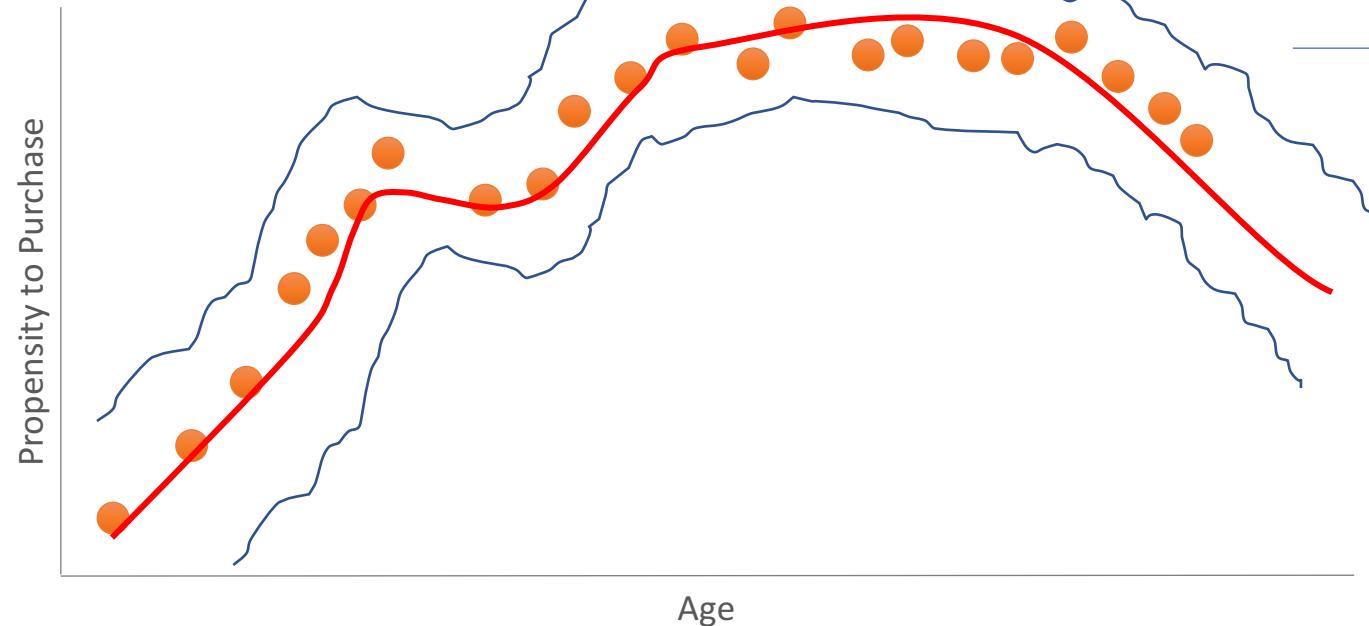


Linear Models

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



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

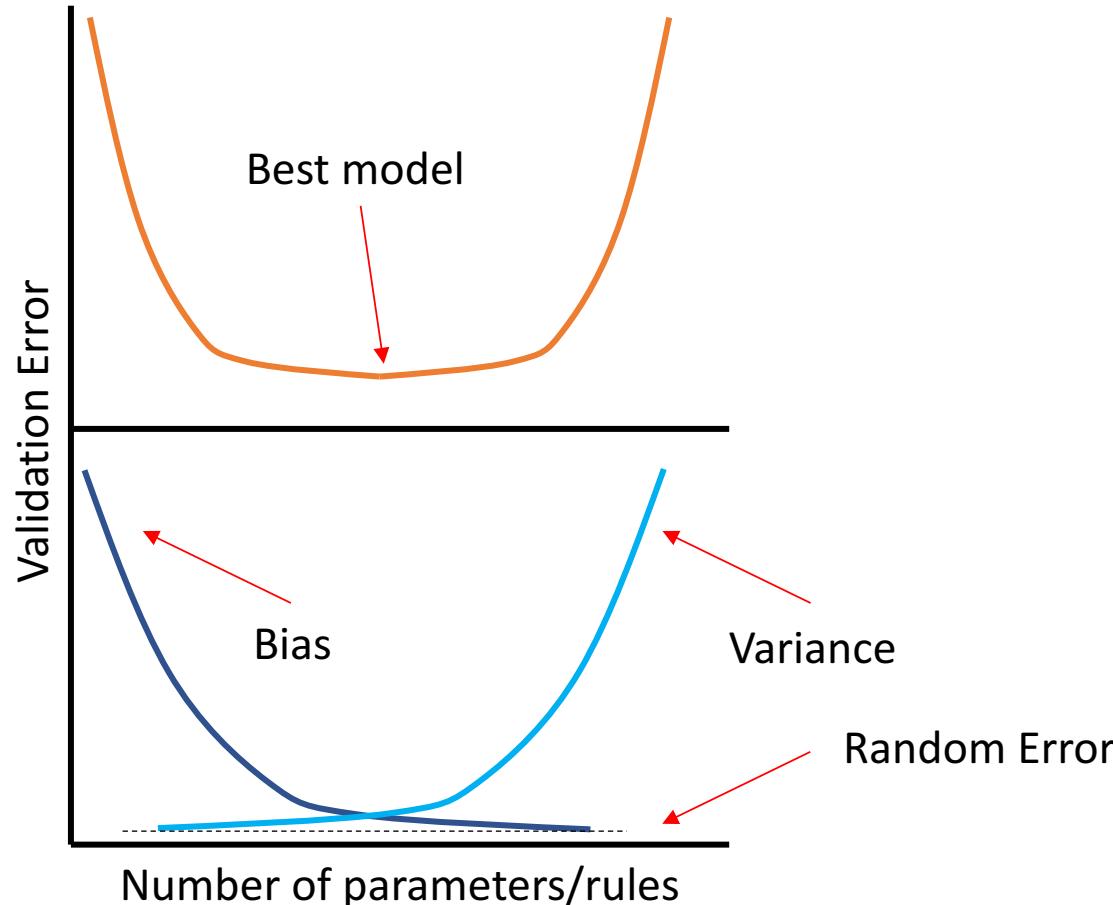


Nobody really believes that multivariate data is multivariate normal, but that data model occupies a large number of pages in every graduate textbook on multivariate statistical analysis.

-- Leo Breiman

Risk from Unwanted Bias and Prediction Variance

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

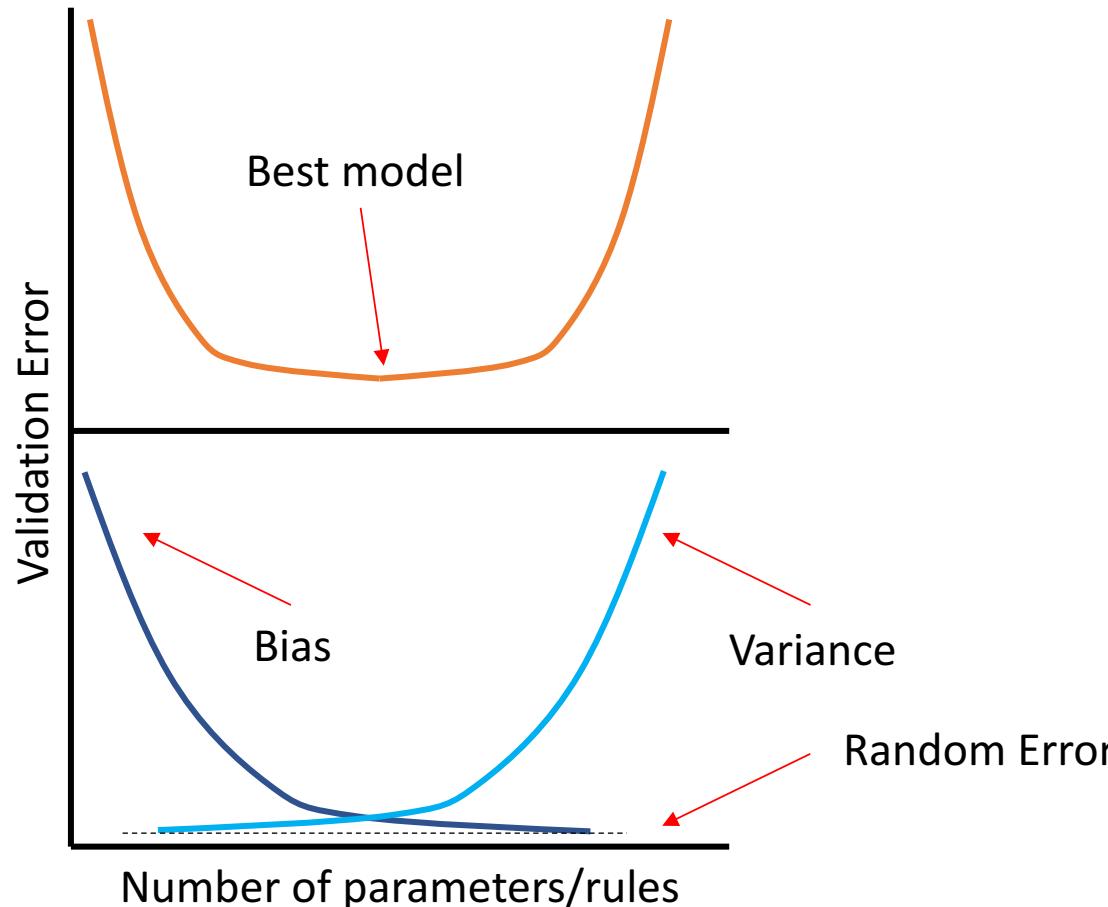


Bias = $E[\hat{f}(x)] - f(x)$ or the error that arises from a model's inability to replicate the fundamental phenomena represented by a data set.

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

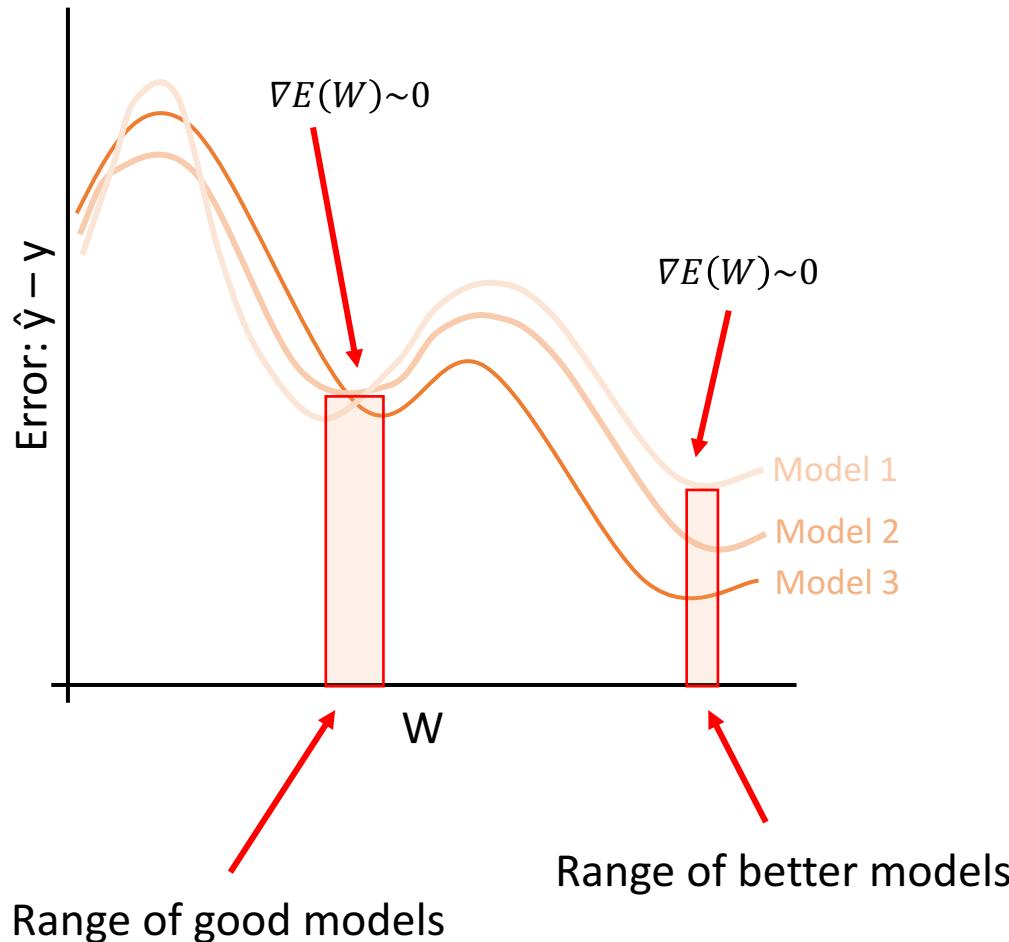
$$\text{Total Error} = \text{Bias} + \text{Variance} + \text{Random Error} = (\hat{f}(x) - f(x))^2$$



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.

The Multiplicity of Good Models



Training ML models often involves solving non-convex optimization problems with multiple local minima.

Different solutions for a good ML model produce the same distribution of model outputs, *but* the actual numeric predictions produced can be slightly different.

These small differences will result in slightly different explanations. Mathematically this is ok ... it's a different model after all, *but* it poses serious philosophical and regulatory problems.

A framework for interpretability

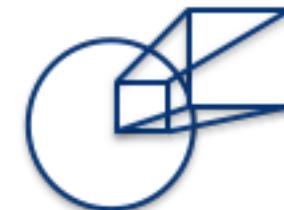
Complexity of learned functions:

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



Scope of interpretability:

Global vs. local



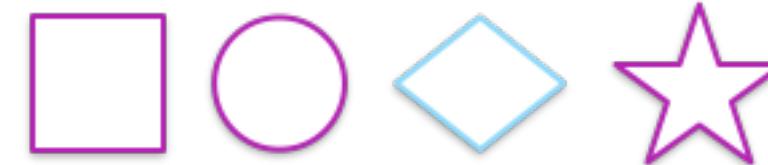
Enhancing trust and understanding:

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



Application domain:

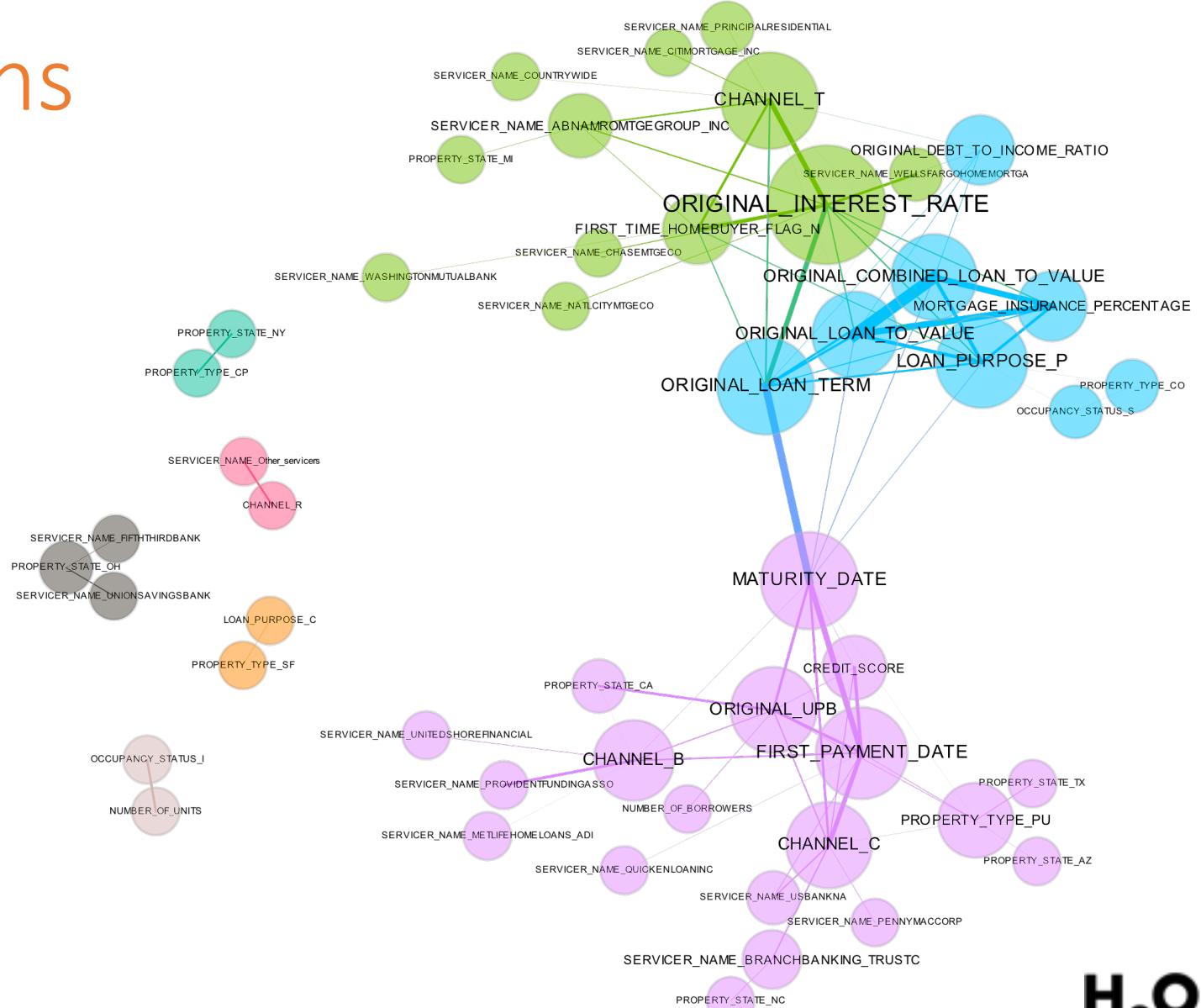
Model-agnostic vs. model-specific



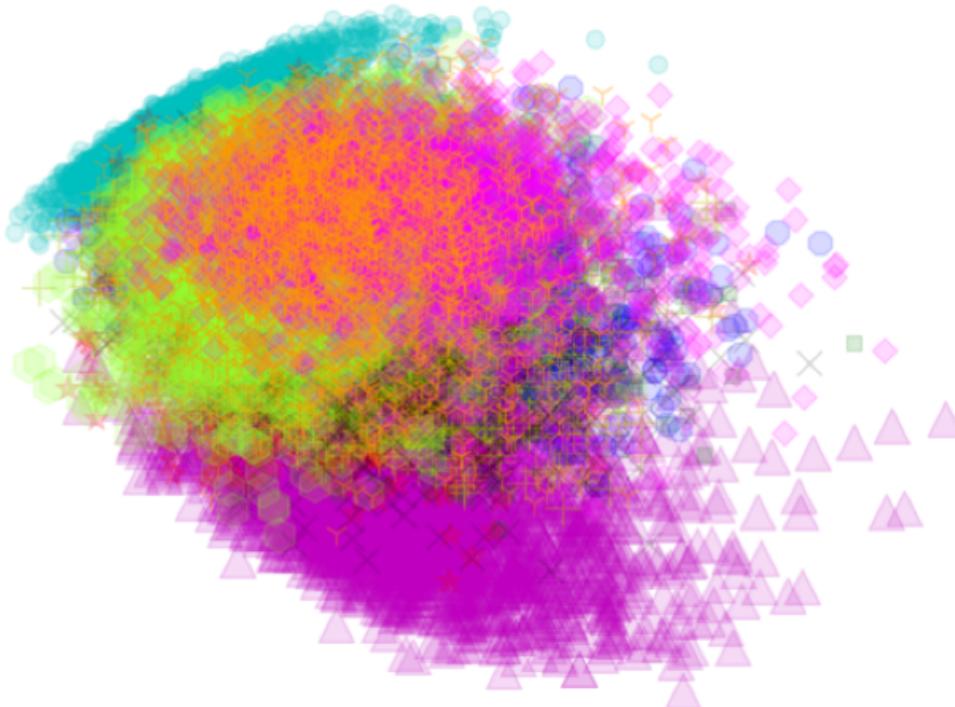
A Few of Our Favorite Things

Correlation graphs

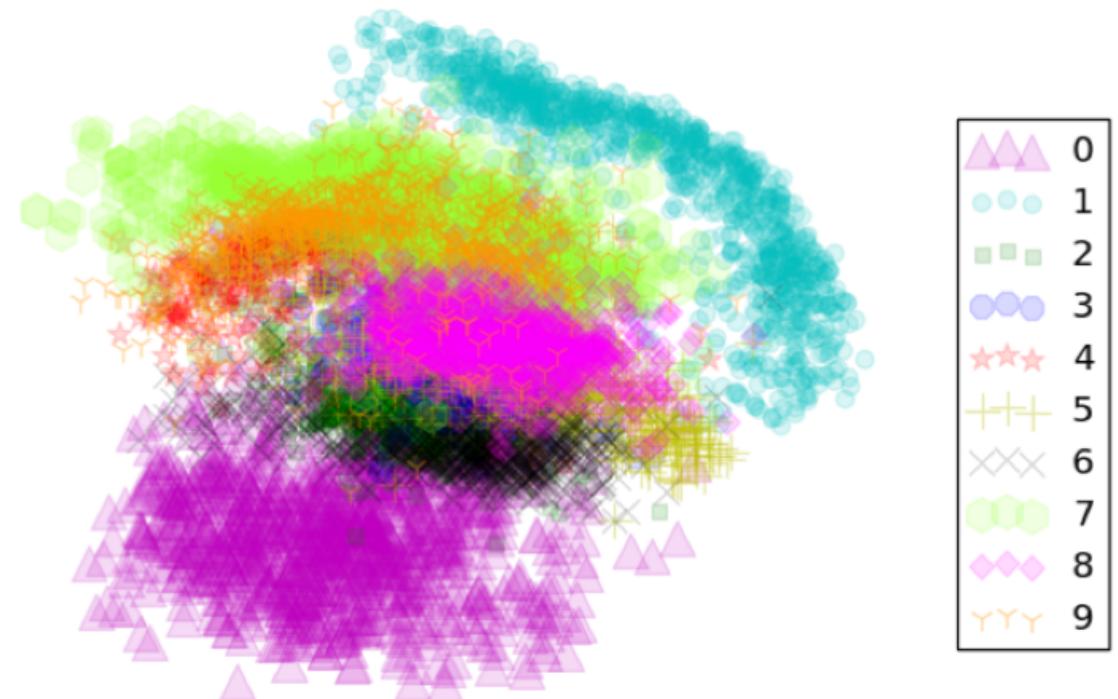
The nodes of this graph are the variables in a data set. The weights between the nodes are defined by the absolute value of their pairwise Pearson correlation.



2-D projections

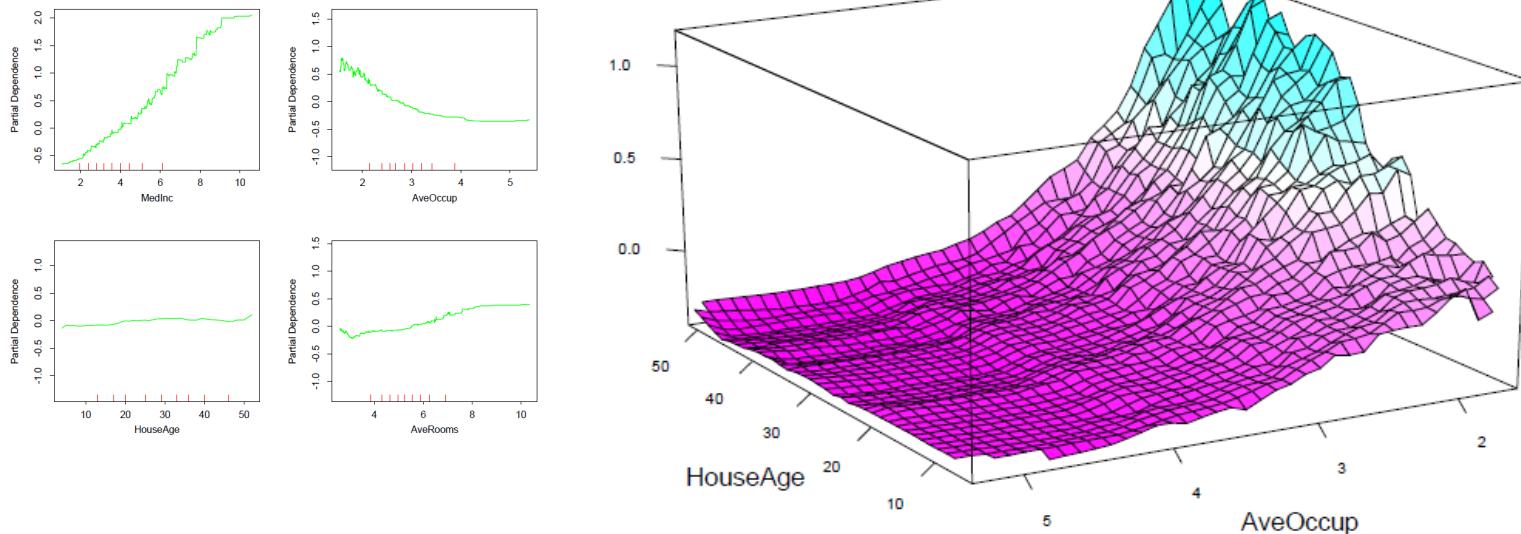


784 dimensions to 2 dimensions with PCA



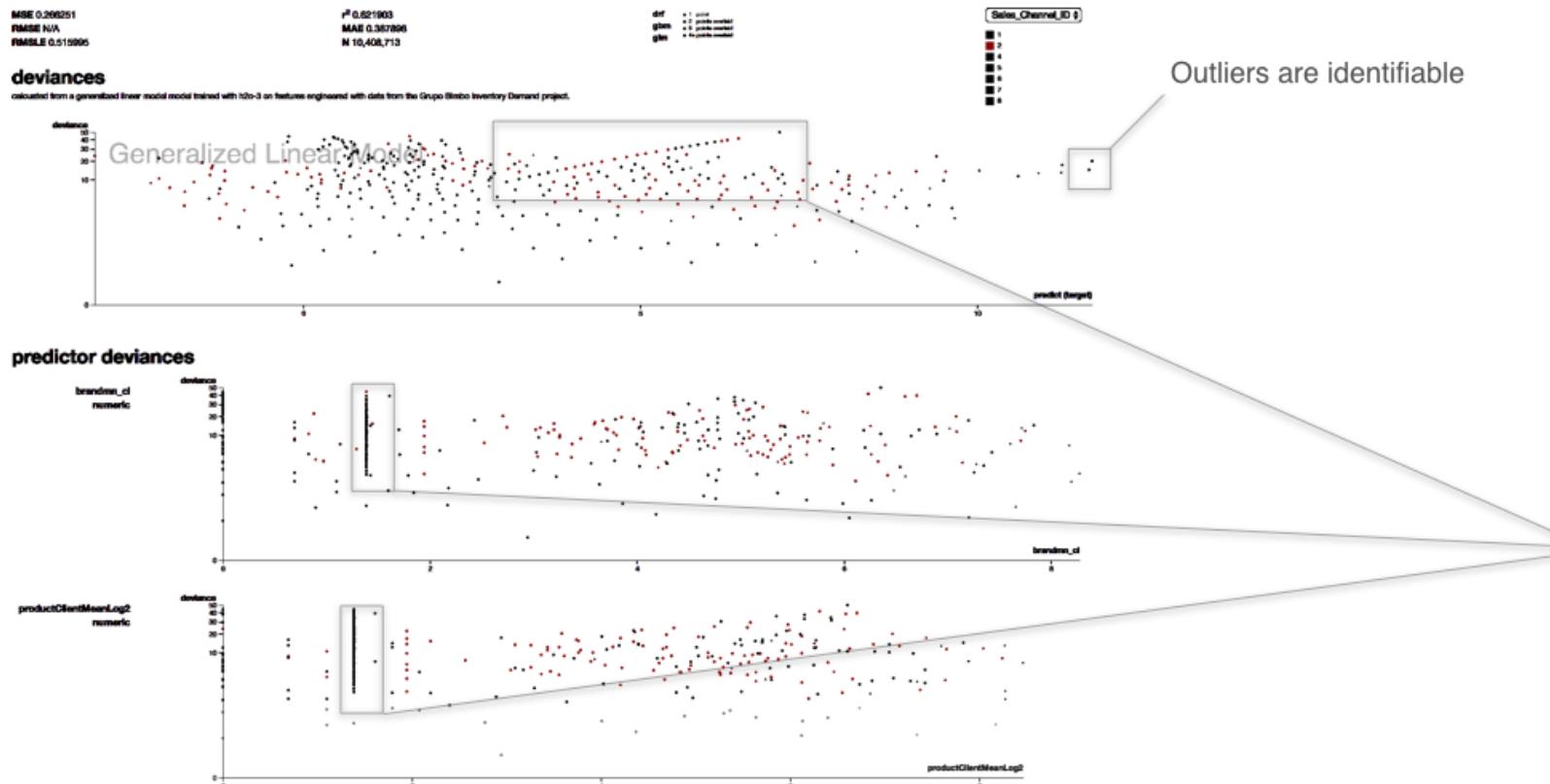
784 dimensions to 2 dimensions with
autoencoder network

Partial dependence plots



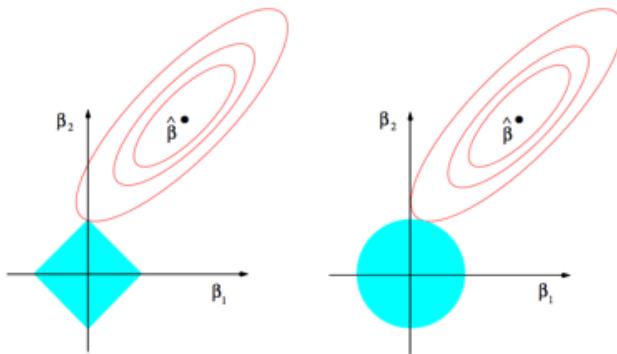
HomeValue ~ MedInc + AveOccup + HouseAge + AveRooms

Residual analysis

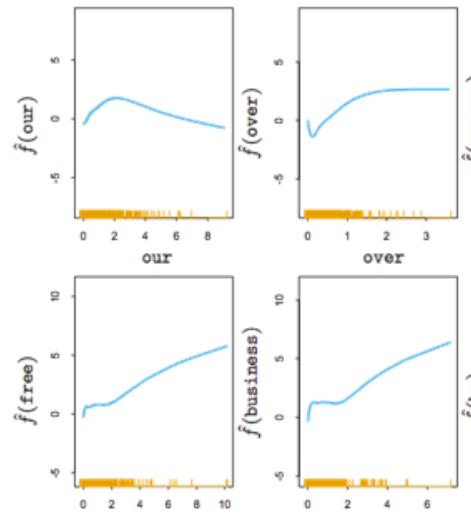


Residuals from a machine learning model should be randomly distributed
obvious patterns in residuals can indicate problems with data preparation or model specification

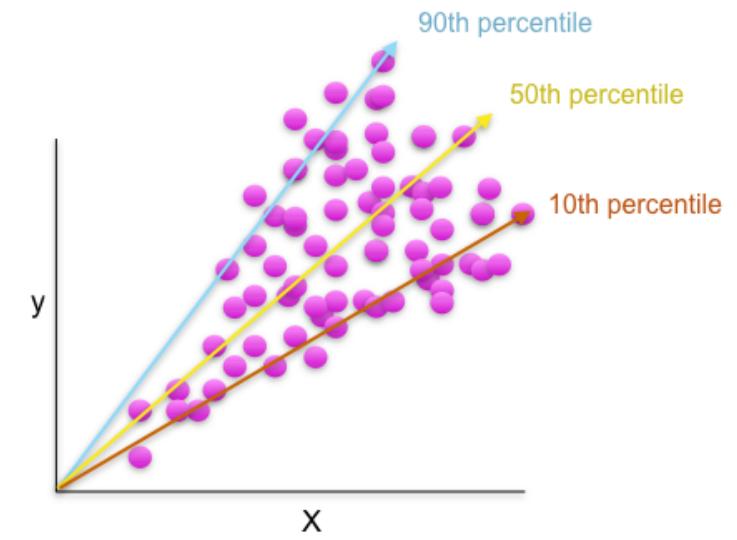
OLS regression alternatives



Penalized Regression



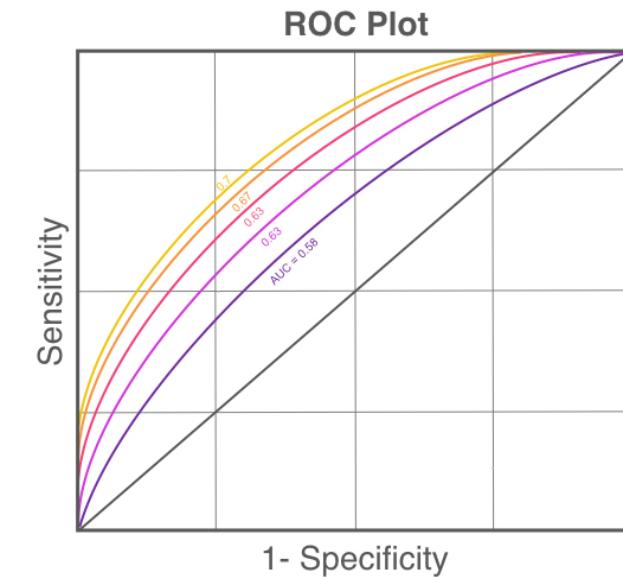
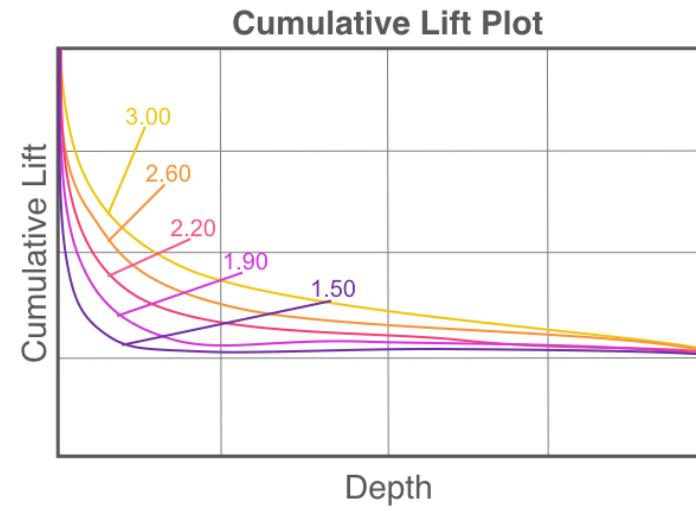
Generalized Additive Models



Quantile Regression

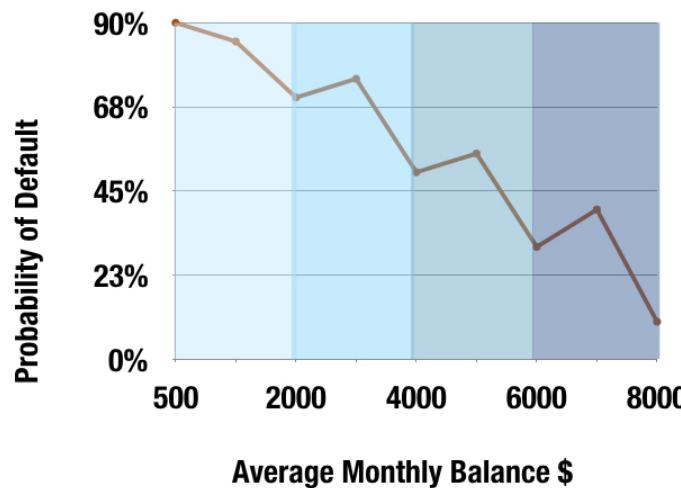
Build toward ML model benchmarks

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

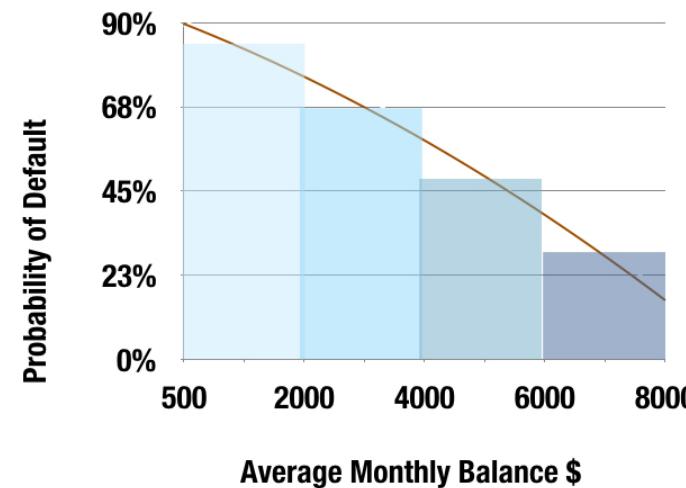


Incorporate interactions and piecewise linear components to increase the accuracy of linear models relative to machine learning models.

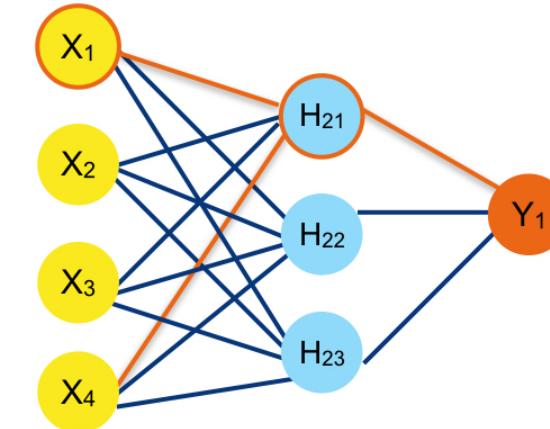
Monotonicity constraints



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

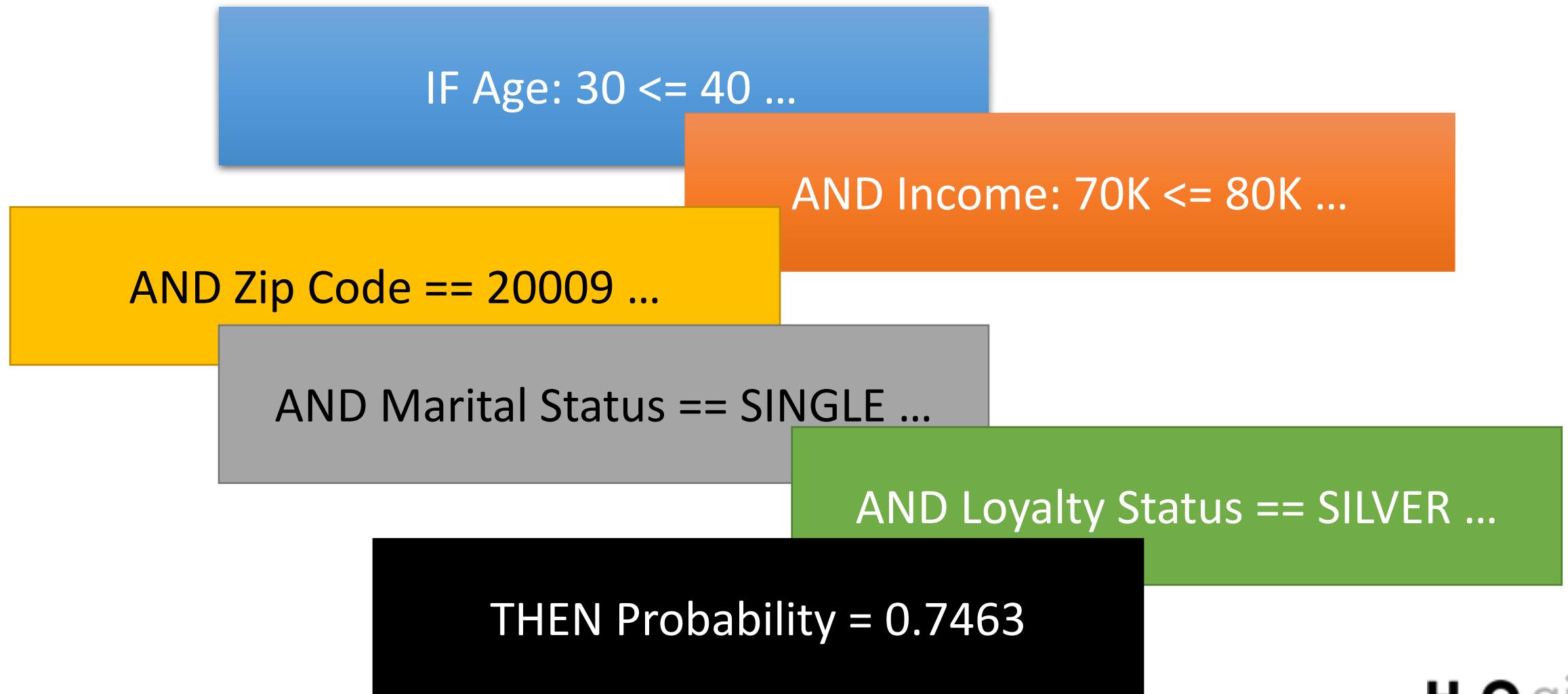


By discretization, the Average Monthly Balance can be transformed to be monotonic with respect to the target.



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.

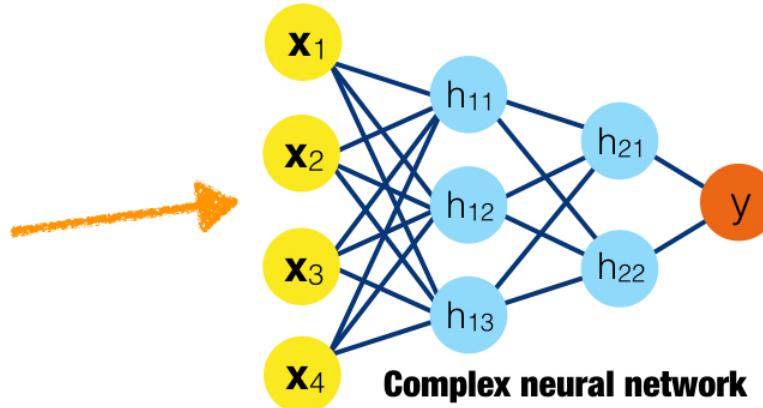
Rule-based models



Surrogate models

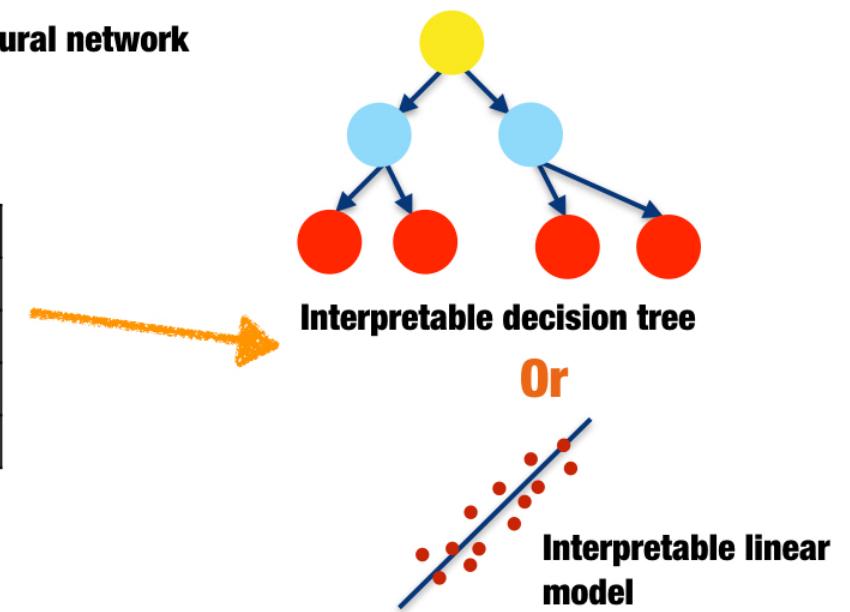
BAD	CUSTOMER_DTI	LOAN_PURPOSE	CHANNEL
0	0.18	MORT	7
1	0.42	HELOC	10
0	0.11	MORT	10
0	0.21	MORT	1

1. Train a complex machine learning model

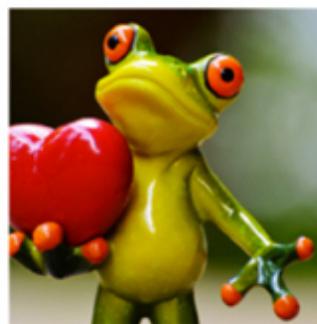


BAD	PREDICTED_BAD	CUSTOMER_DTI	LOAN_PURPOSE	CHANNEL
0	0.47	0.18	MORT	7
1	0.82	0.42	HELOC	10
0	0.18	0.11	MORT	10
0	0.12	0.21	MORT	1

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



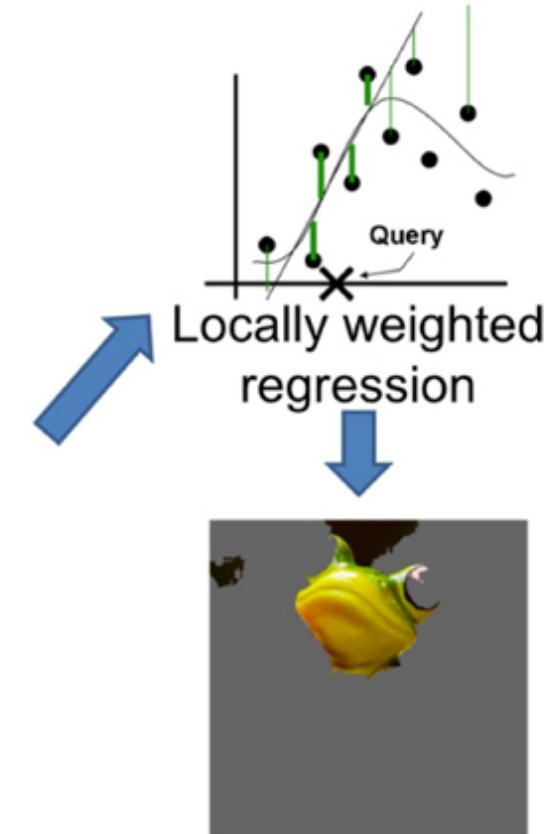
Local interpretable model-agnostic explanations



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

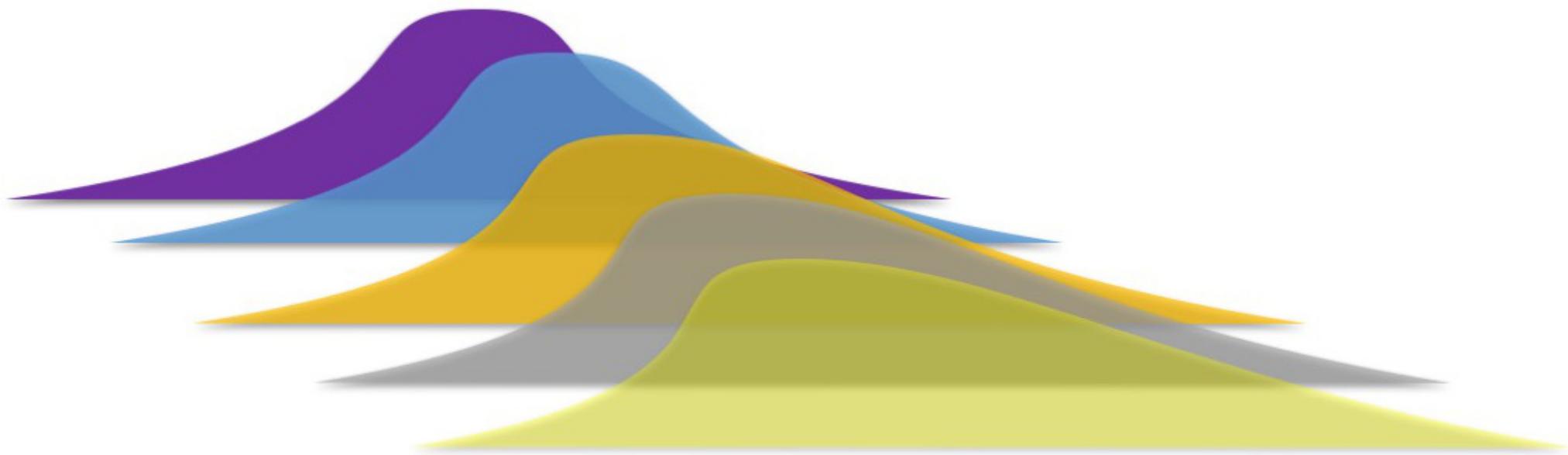


Perturbed Instances	$P(\text{tree frog})$
A photograph of the same tree frog with several small red spots added to its back.	0.85
A photograph of the same tree frog with its body turned almost entirely yellow.	0.00001
A photograph of the same tree frog with red flowers added to its front legs.	0.52



Explanation

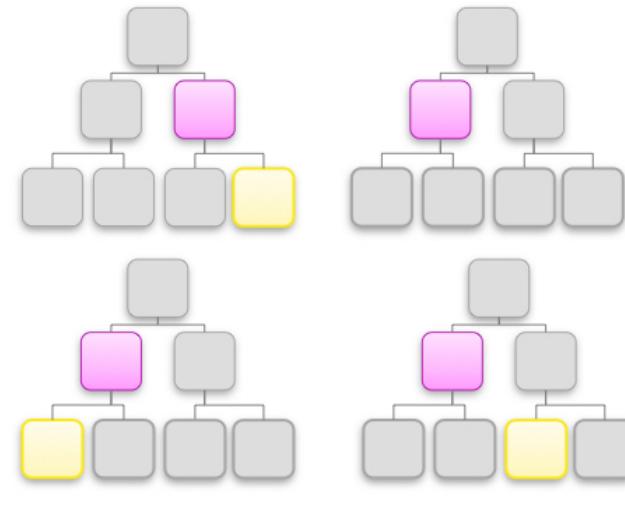
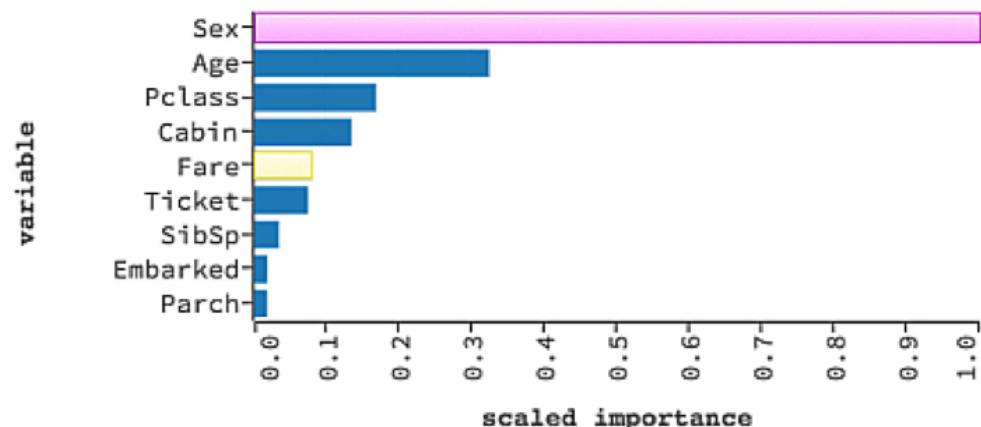
Sensitivity analysis



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.

H2O.ai Driverless AI



H₂O.ai



Featured Prediction Competition

Mercedes-Benz Greener Manufacturing

Can you cut the time a Mercedes-Benz spends on the test bench?

Daimler · 1,724 teams · a month to go (25 days to go until merger deadline)

\$25,000 Prize Money

Overview Data Kernels Discussion Leaderboard More My Submissions Submit Predictions

Your most recent submission				
Name	Submitted	Wait time	Execution time	Score
blend.csv	just now	1 seconds	0 seconds	0.56842
Complete				

Jump to your position on the leaderboard ▾

Public Leaderboard Private Leaderboard

This leaderboard is calculated with approximately 19% of the test data.
The final results will be based on the other 81%, so the final standings may be different.

Raw Data Refresh

■ In the money ■ Gold ■ Silver ■ Bronze

#	△1w	Team Name	Kernel	Team Members	Score ⓘ	Entries	Last
1	—	Utkarsh			0.57456	40	9h
2	▲ 531	plantsgo			0.57005	9	3h
3	▼ 1	Fred Navruzov			0.56904	23	4d
4	new	kongshuchen			0.56856	5	13h
5	▲ 115	Jughead			0.56844	27	5h
6	▲ 29	H2O.ai Driverless AI			0.56842	11	now

Your Best Entry ▾

Your submission scored 0.56842, which is an improvement of your previous score of 0.56605. Great job!

Tweet this!

7	▼ 4	happiness			0.56814	25	1d
8	▲ 389	Luis Moneda			0.56809	15	12h
9	—	Pan Tofelek			0.56800	35	12h
10	▼ 6	gotohell1			0.56790	15	16h

TRAINING DATA

DATASET
default_of_credit_card_clients.csv

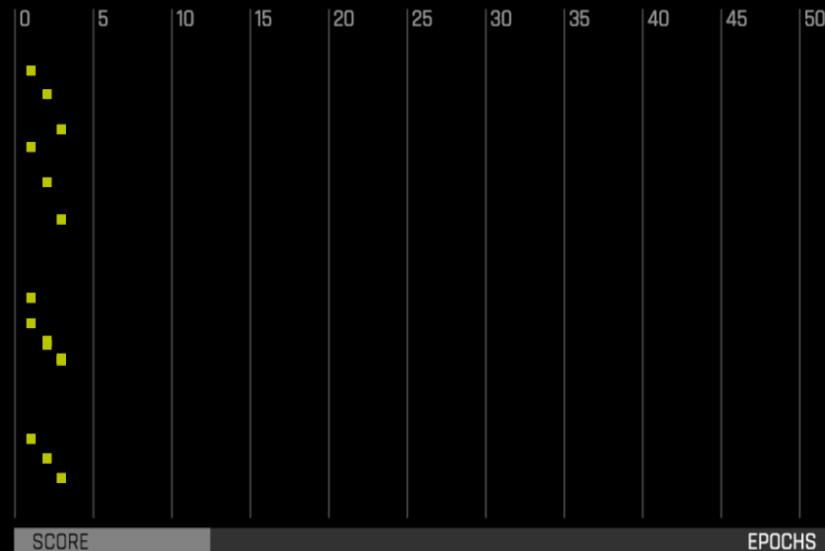
ROWS
35m COLUMNS
24 DROPPED
0 IGNORED
0

TARGET COLUMN

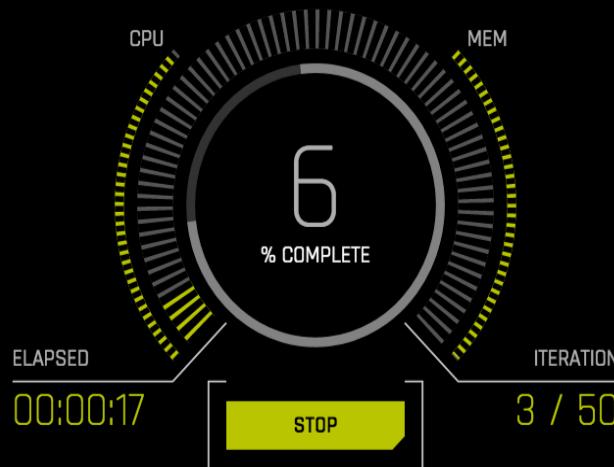
default_on_payment

TYPE
float MISSING
3% LEVELS
11 MEAN
56.7

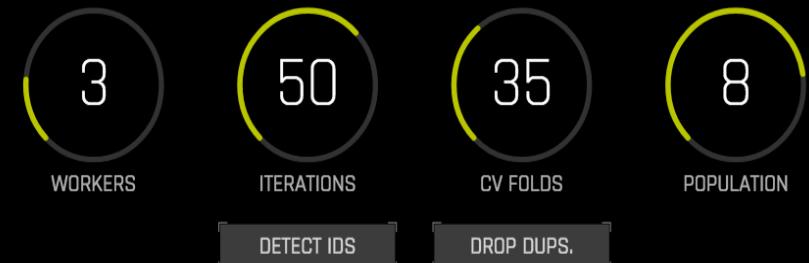
ITERATION SCORES



STATUS: RUNNING



EXPERIMENT SETTINGS



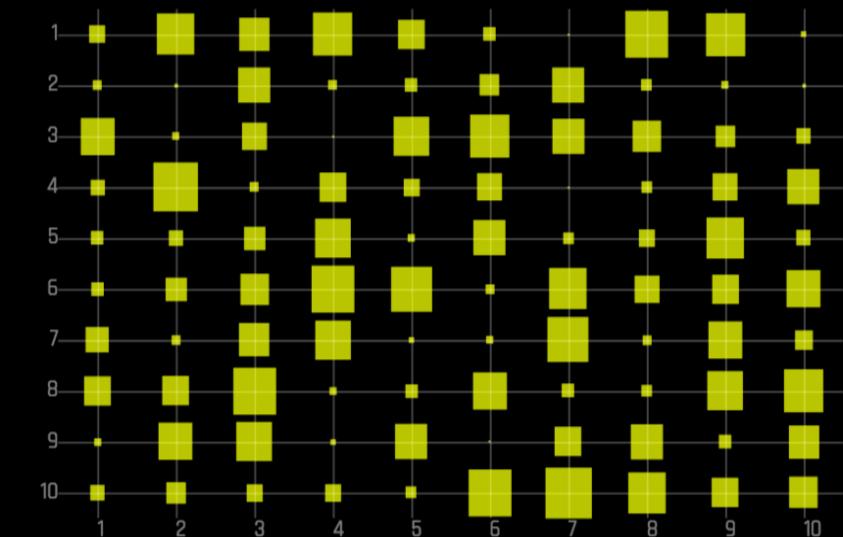
GPU STATS



VARIABLE IMPORTANCE

PAY_4	58
PAY_3	16
LIMIT_BAL	3
BILL_AMT5	46
PAY_AMT5	62
PAY_AMT3	88
BILL_AMT2	72
ID	77
AGE	47
PAY_0	35
SEX	76
BILL_AMT3	15
EDUCATION	94
PAY_6	20

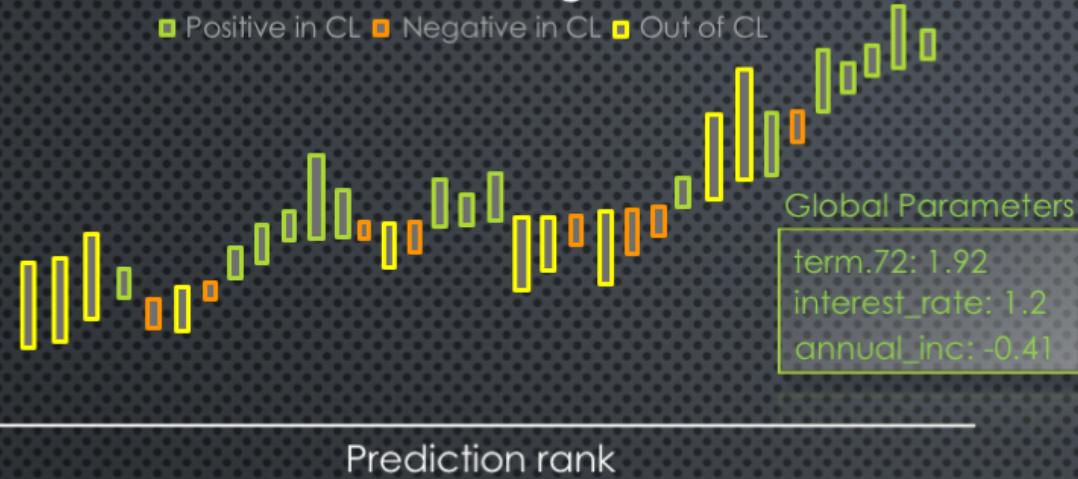
FEATURE TRANSFORMATIONS



Prediction for bad_loan

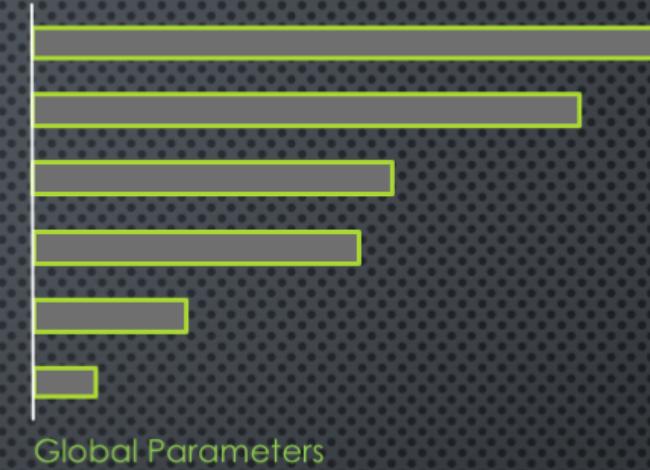
k-LIME Surrogate

■ Positive in CL ■ Negative in CL ■ Out of CL



Prediction rank

Variable Importance



interest_rate: 1.00

term: 0.89

annual_inc: 0.71

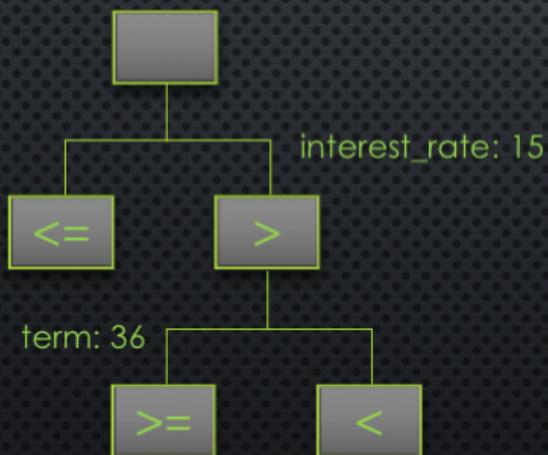
addr_state: 0.66

dti: 0.39

revol_util: 0.23

Global Parameters

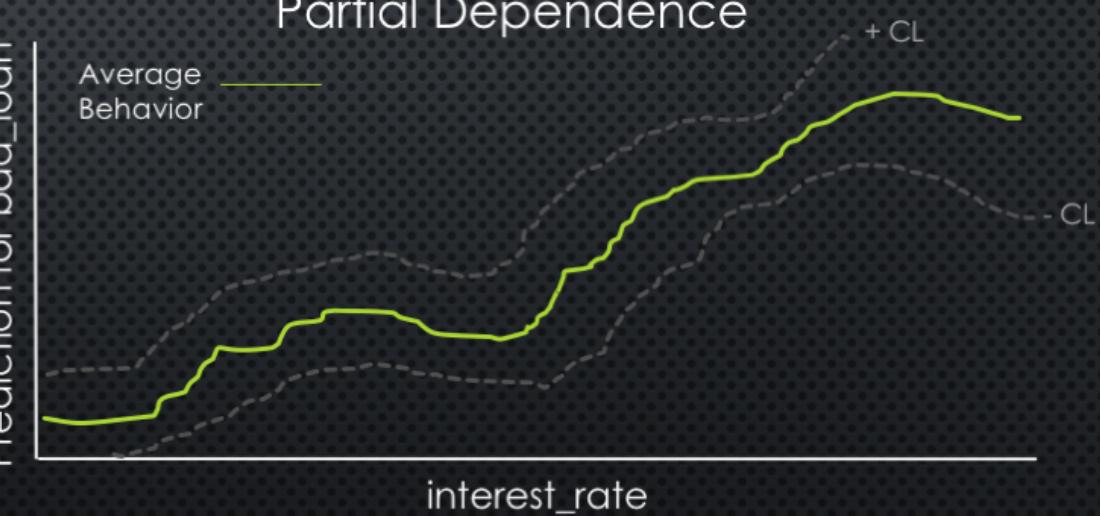
Decision Tree Surrogate



Prediction for bad_loan

Partial Dependence

Average Behavior



Globally, on average ...

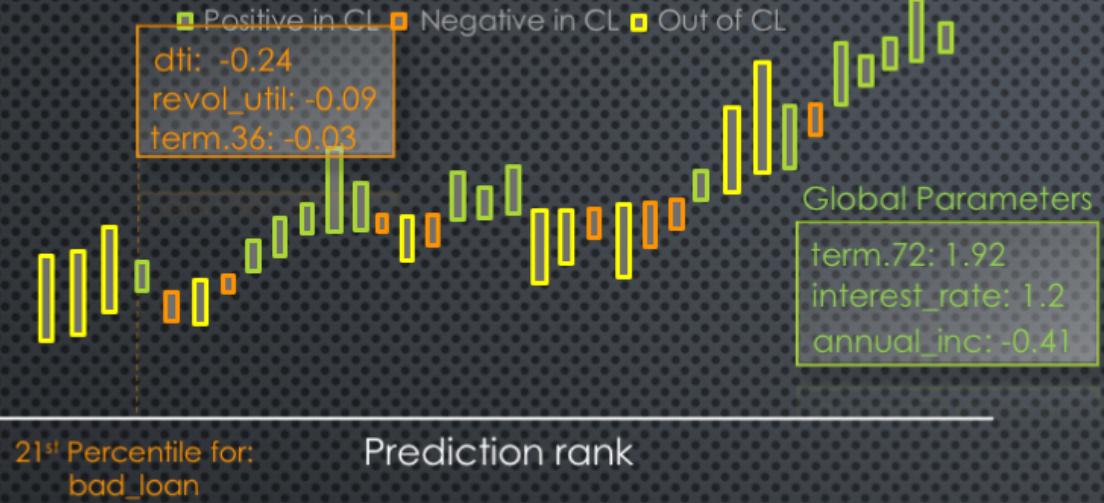
term = 72 causes a 1.92-unit **increase** in the probability of bad_loan.

A 1-unit change in interest_rate causes a 1.2-unit **increase** in the probability of bad_loan.

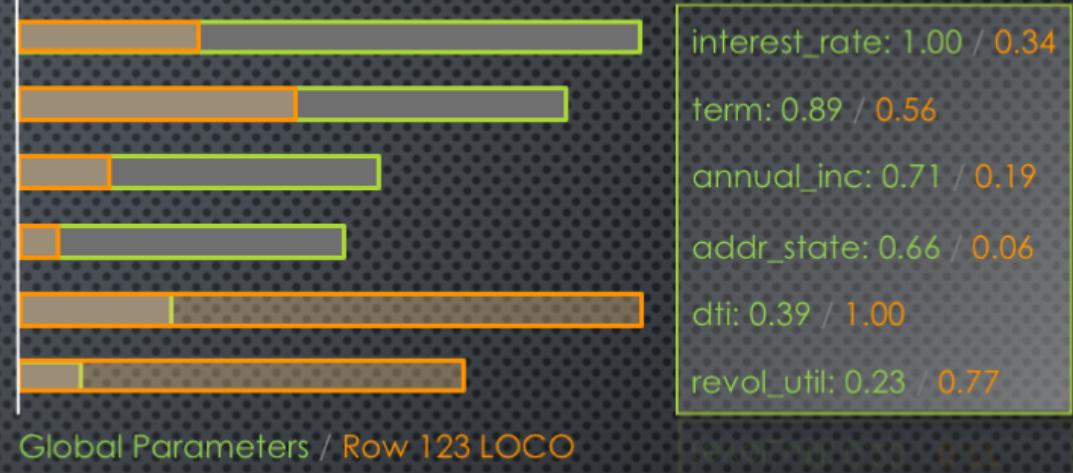
A 1-unit change in annual_inc causes a 0.41-unit **decrease** in the probability of bad_loan.

Prediction for bad_loan

k-LIME Surrogate

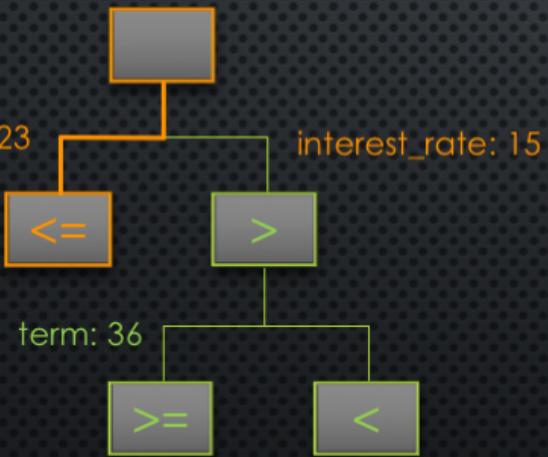


Variable Importance



Decision Tree Surrogate

Decision Path for Row 123



Globally, on average ...

term = 72 causes a 1.92-unit **increase** in the probability of bad_loan.

A 1-unit change in interest_rate causes a 1.2-unit **increase** in the probability of bad_loan.

A 1-unit change in annual_inc causes a 0.41-unit **decrease** in the probability of bad_loan.

Partial Dependence



Locally, for row 123 ...

dti = 2.0 causes a 0.24-unit **decrease** in the probability of bad_loan.

revol_util = 50.0 causes a 0.09-unit **decrease** in the probability of bad_loan.

term = 36 causes a 0.03-unit **decrease** in the probability of bad_loan.

Questions?