

Ideas on Machine Learning Interpretability

Navdeep Gill, H2O.ai

Big Ideas

UNKNOWN TARGET FUNCTION

$$f: \mathcal{X} \rightarrow \mathcal{Y}$$

(ideal credit approval function)

TRAINING EXAMPLES

$$(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$$

(historical records of credit customers)

LEARNING ALGORITHM



FINAL HYPOTHESIS

$$g \approx f$$

(final credit approval formula)

HYPOTHESIS SET

$$\mathcal{H}$$

(set of candidate formulas)

Learning from data ...

Adapted from:

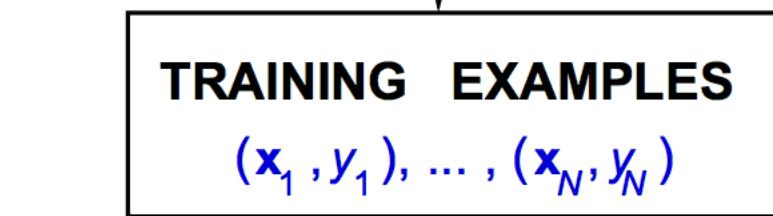
Learning from Data. <https://work.caltech.edu/textbook.html>

UNKNOWN TARGET FUNCTION

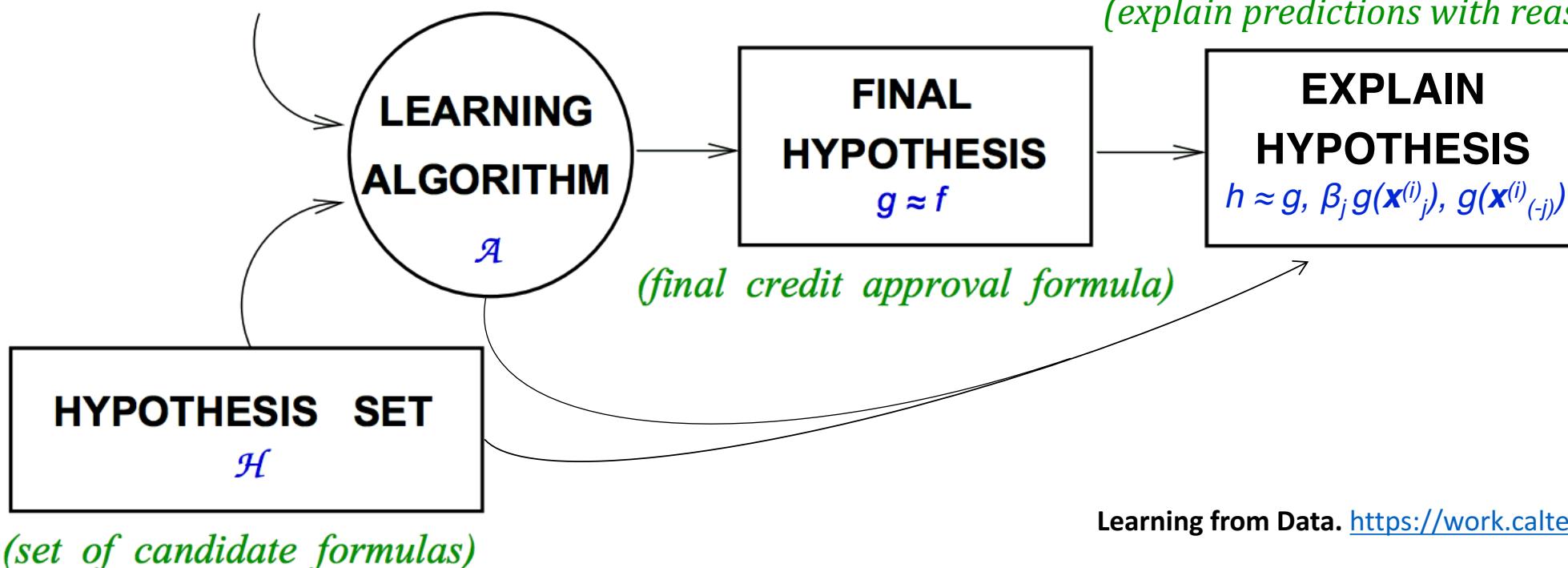
$$f: \mathcal{X} \rightarrow \mathcal{Y}$$

(ideal credit approval function)

Learning from data ...
transparently.



(historical records of credit customers)

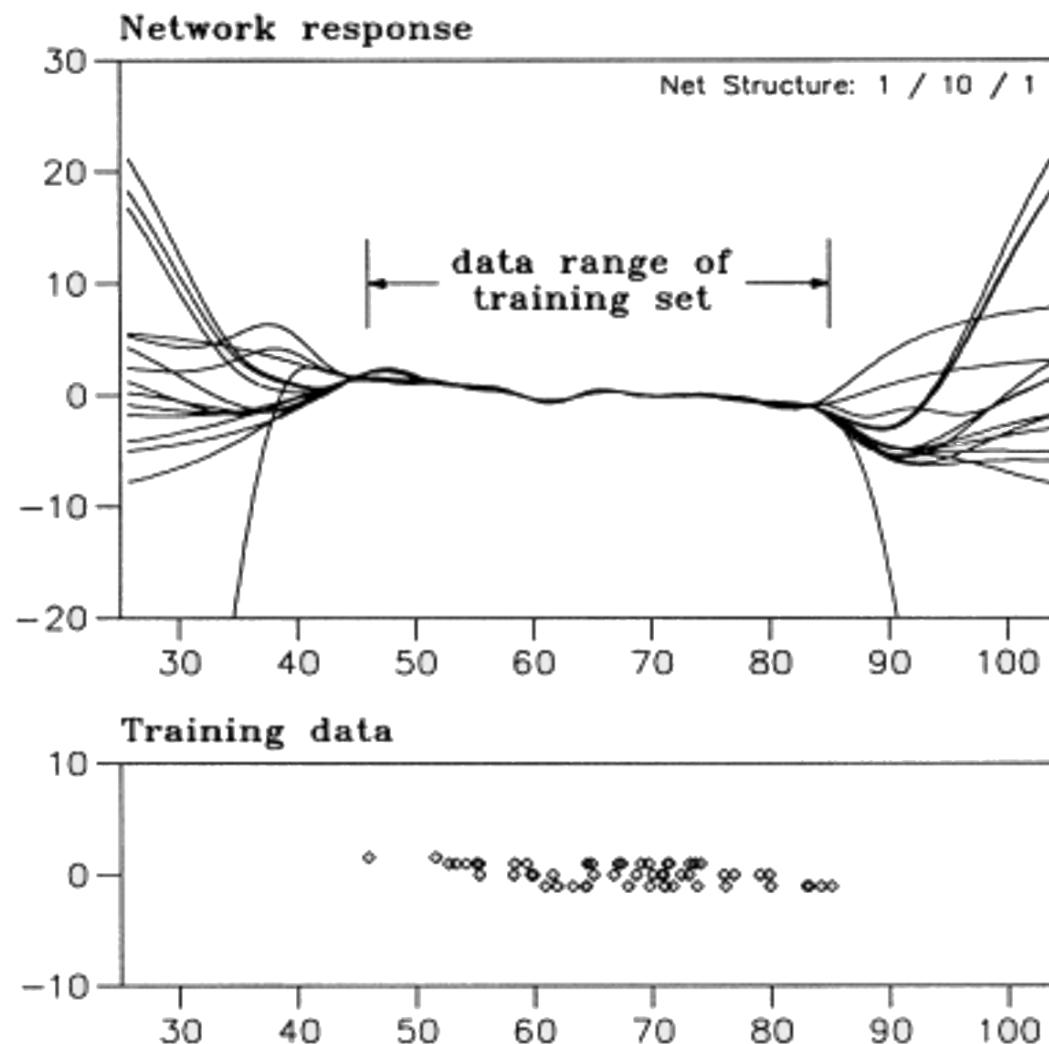


Adapted from:
Learning from Data. <https://work.caltech.edu/textbook.html>

Increasing fairness, accountability, and trust by
decreasing unwanted sociological biases



Increasing trust by quantifying prediction variance



A framework for interpretability

Complexity of learned functions:

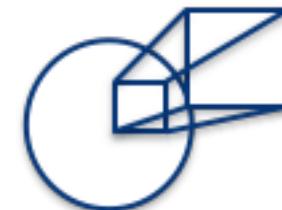
- Linear, monotonic
- Nonlinear, monotonic
- Nonlinear, non-monotonic

(~ Number of parameters/VC dimension)



Scope of interpretability:

Global vs. local



Enhancing trust and understanding:

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

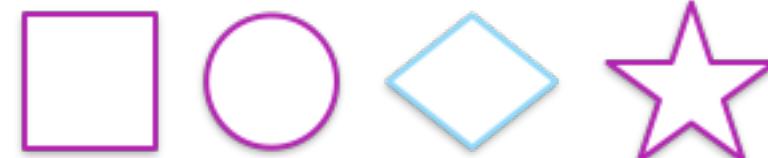


Understanding ~ transparency

Trust ~ fairness and accountability

Application domain:

Model-agnostic vs. model-specific

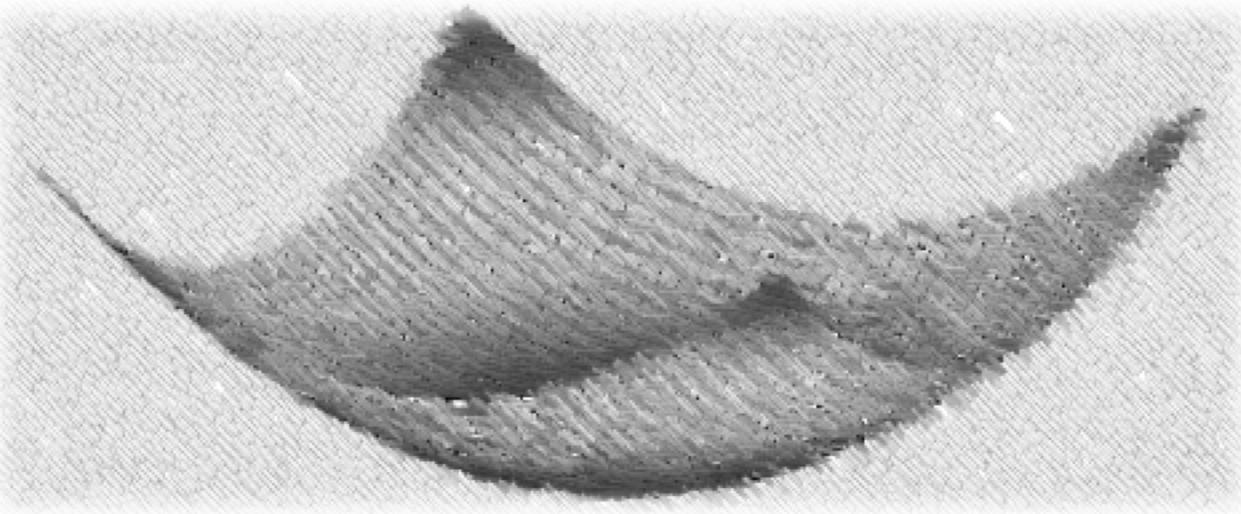


Big Challenges

Linear Models

Strong model locality

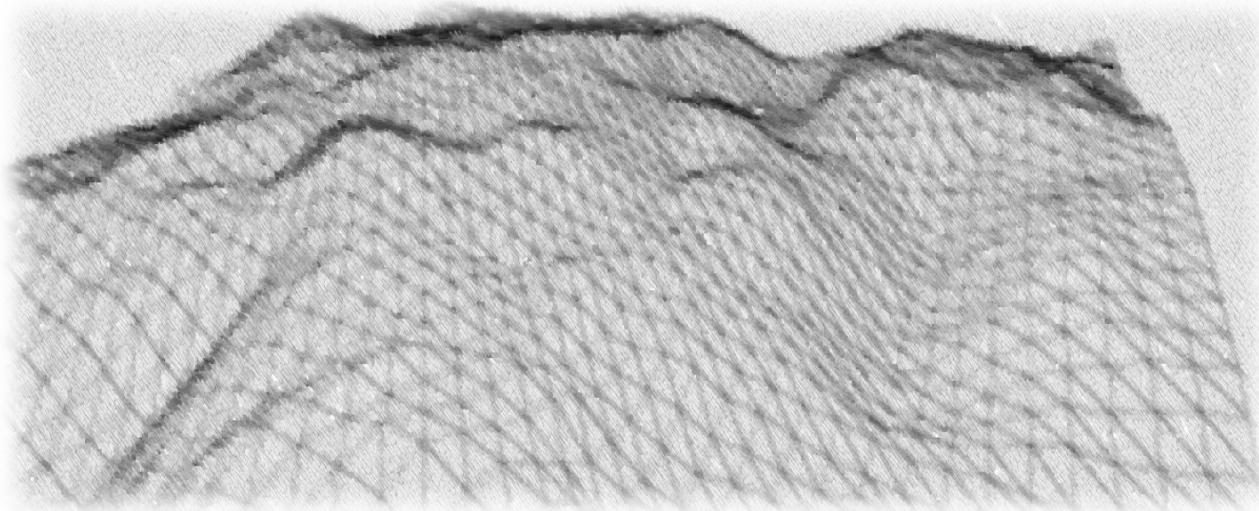
Usually stable models and explanations



Machine Learning

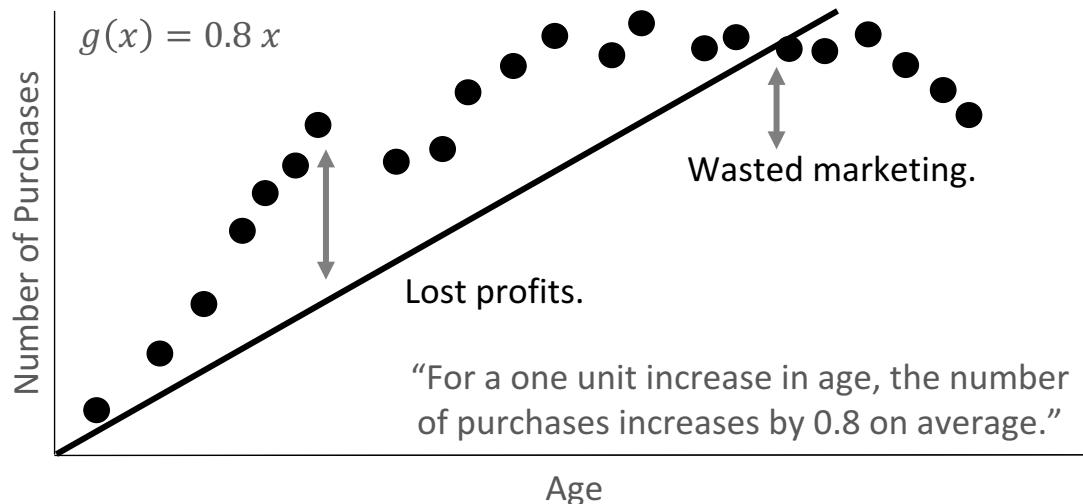
Weak model locality

Sometimes unstable models and explanations
(a.k.a. The Multiplicity of Good Models)



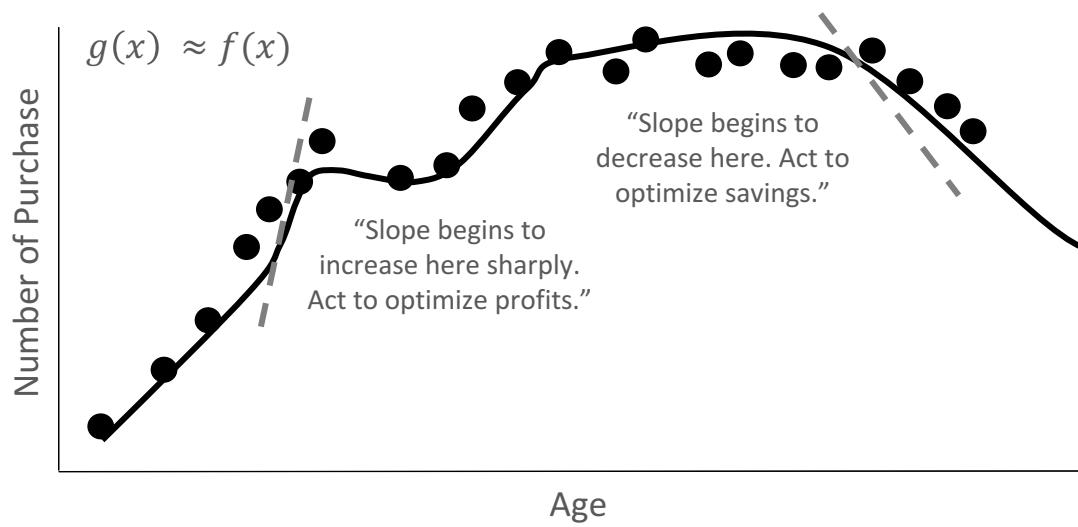
Linear Models

Exact explanations for
approximate models.



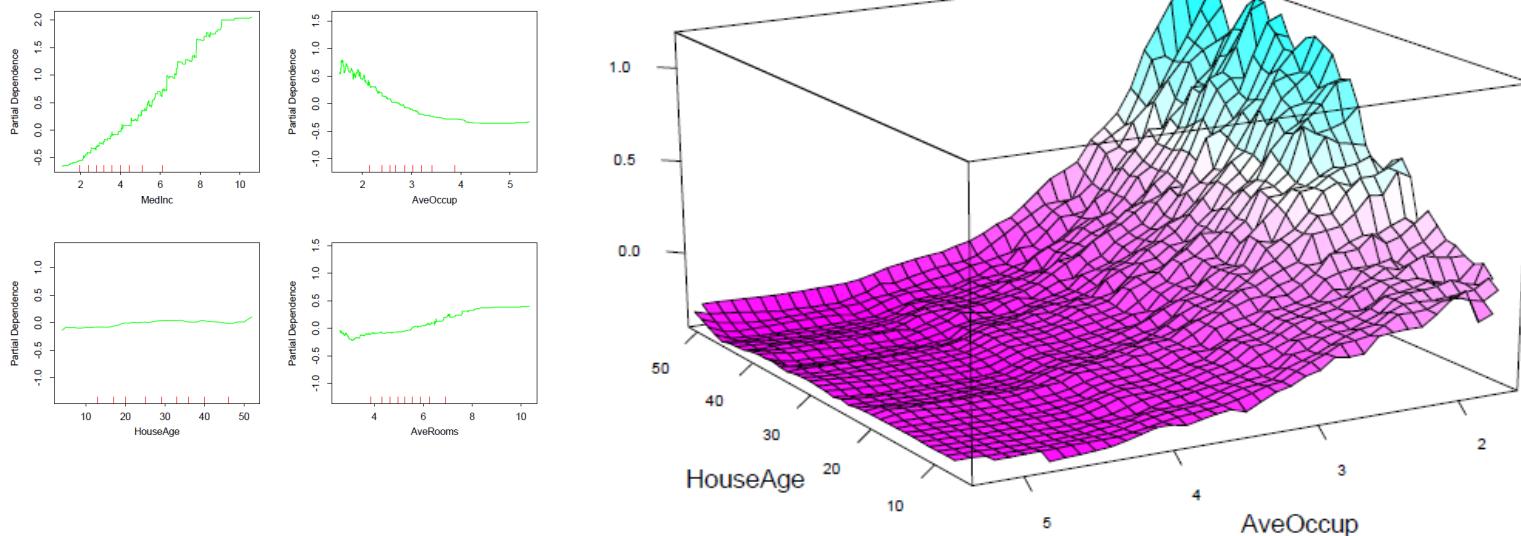
Machine Learning

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



A Few of Our Favorite Things

Partial dependence plots

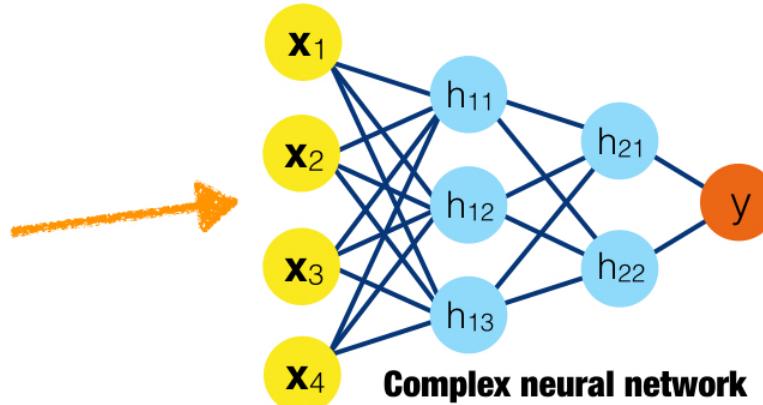


HomeValue ~ MedInc + AveOccup + HouseAge + AveRooms

Surrogate models

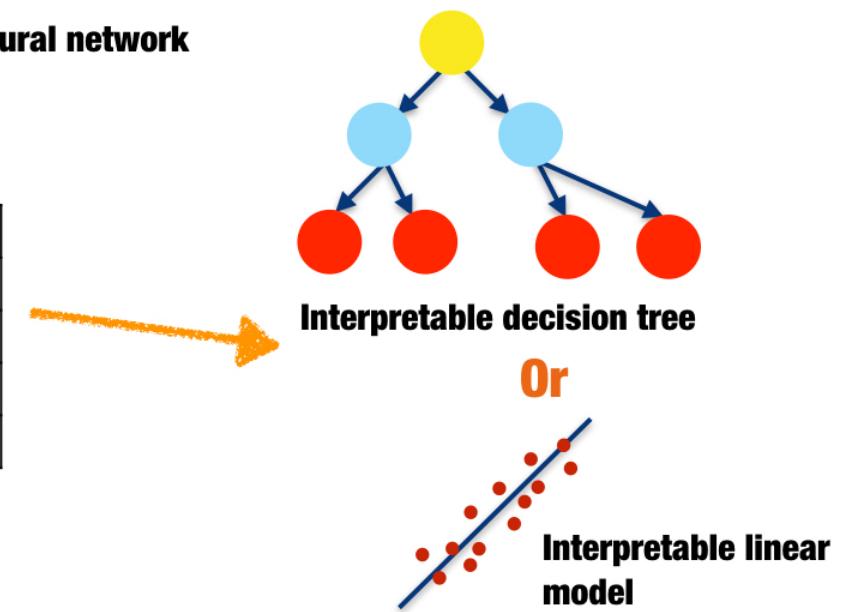
BAD	CUSTOMER_DTI	LOAN_PURPOSE	CHANNEL
0	0.18	MORT	7
1	0.42	HELOC	10
0	0.11	MORT	10
0	0.21	MORT	1

1. Train a complex machine learning model

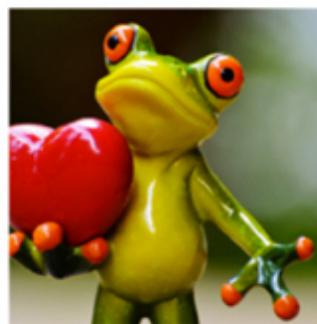


BAD	PREDICTED_BAD	CUSTOMER_DTI	LOAN_PURPOSE	CHANNEL
0	0.47	0.18	MORT	7
1	0.82	0.42	HELOC	10
0	0.18	0.11	MORT	10
0	0.12	0.21	MORT	1

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



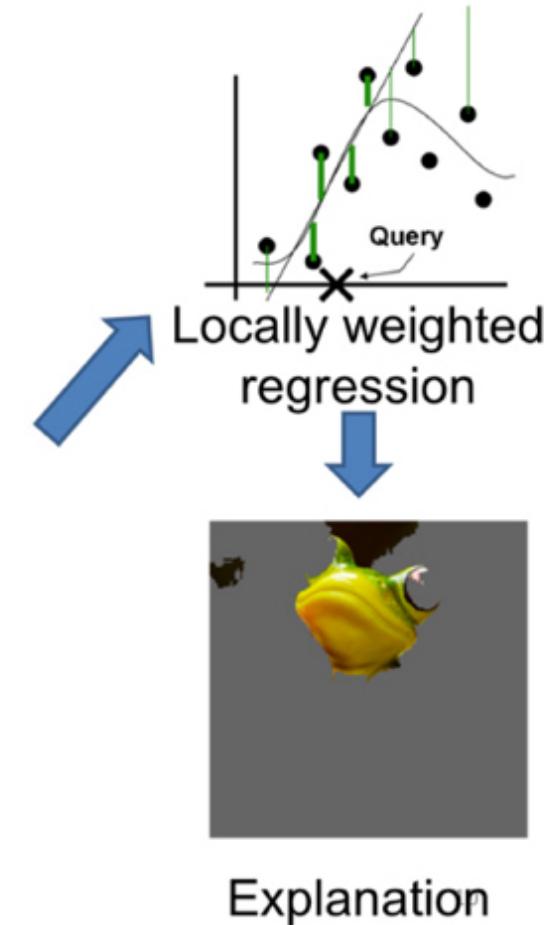
Local interpretable model-agnostic explanations



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

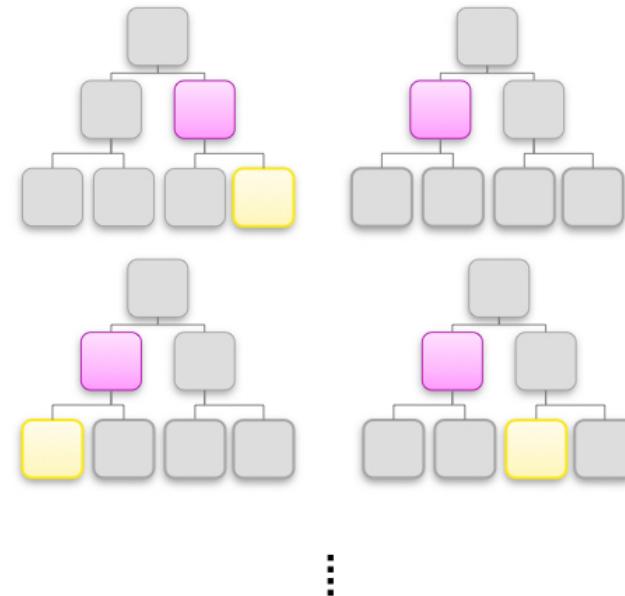
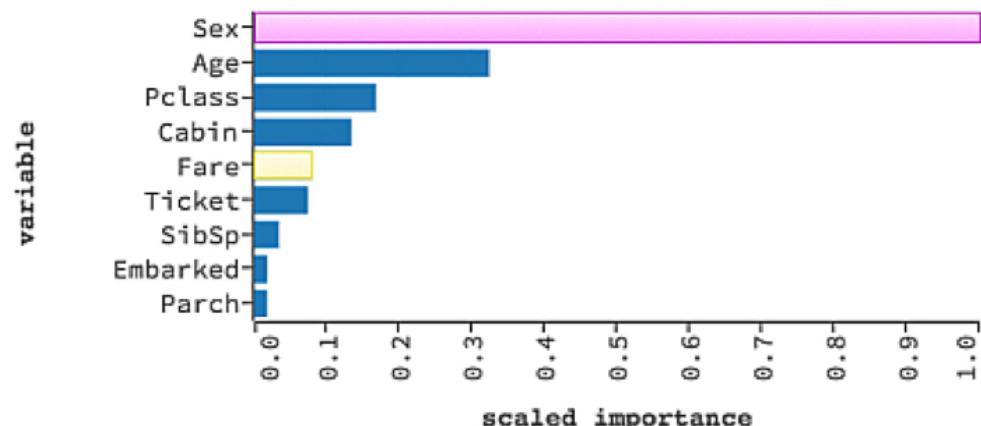


Perturbed Instances	$P(\text{tree frog})$
A photograph of the same frog with several small red spots added to its back.	0.85
A photograph of the same frog with its body color shifted towards yellow.	0.00001
A photograph of the same frog with red flowers added to its front legs.	0.52



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 Demo

0.9.5+LOCAL_584D771

TRAINING DATA

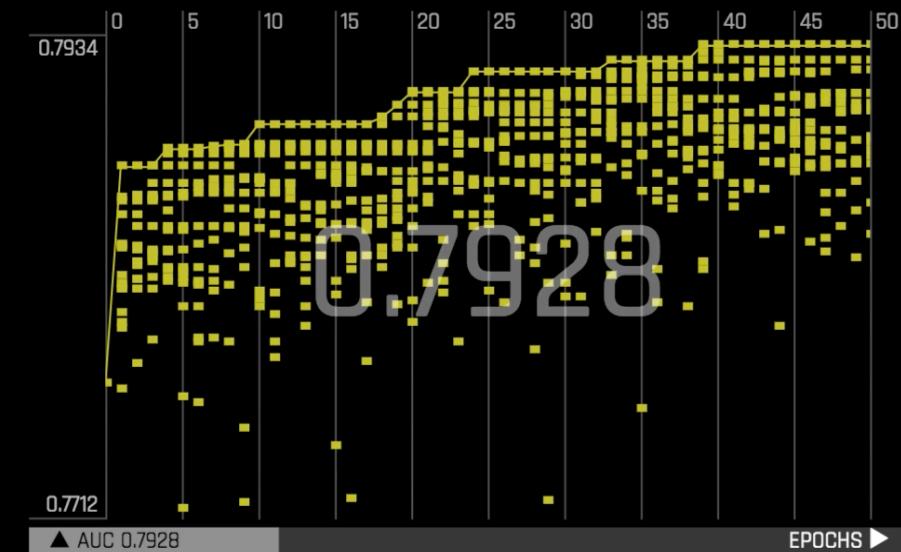
DATASET
creditcard_cat.csvROWS
24KCOLUMNS
25DROPPED COLS
1TEST DATASET
--

TARGET COLUMN

default.payment.next.month

TYPE
intCOUNT
23999UNIQUE
2FREQ
18630

ITERATION SCORES - INTERNAL VALIDATION



STATUS: COMPLETE

- [INTERPRET THIS MODEL](#)
- [SCORE ON ANOTHER DATASET](#)
- [DOWNLOAD \(HOLDOUT\) TRAINING PREDICTIONS](#)

[DOWNLOAD TRANSFORMED TRAINING DATA](#)[DOWNLOAD LOGS](#)[DOWNLOAD SCORING PACKAGE](#)

EXPERIMENT SETTINGS

ACCURACY
5TIME
5INTERPRETABILITY
5

CLASSIFICATION

SET RAND. SEED

ENABLE GPUS

SCORER
R2
MSE
RMSE
RMSLE
MAE
GINI
AUC
LOGLOSS

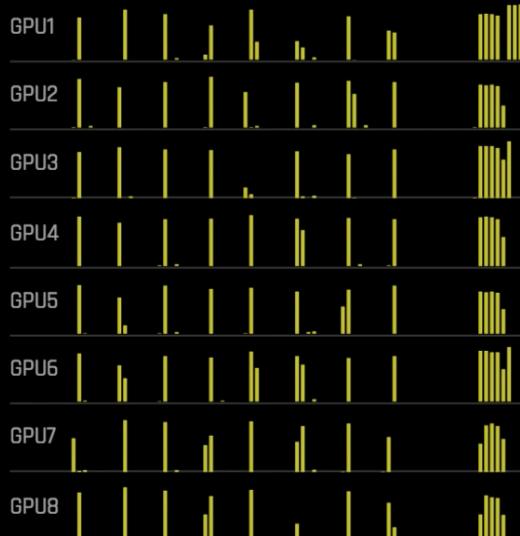
Trace

CPU / MEMORY

CPU

MEM

GPU USAGE

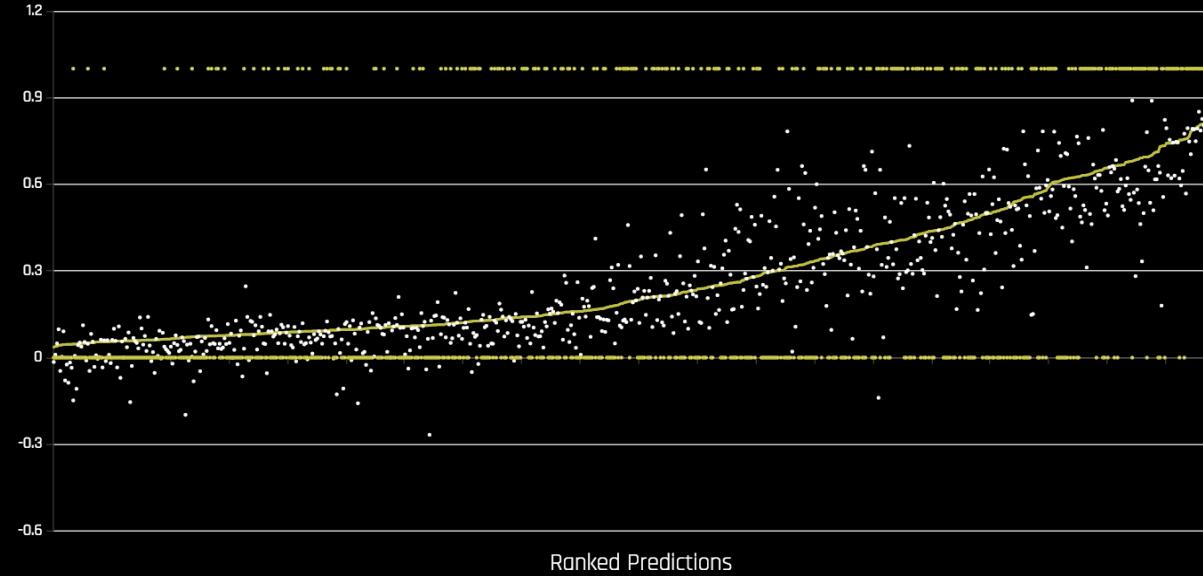


VARIABLE IMPORTANCE

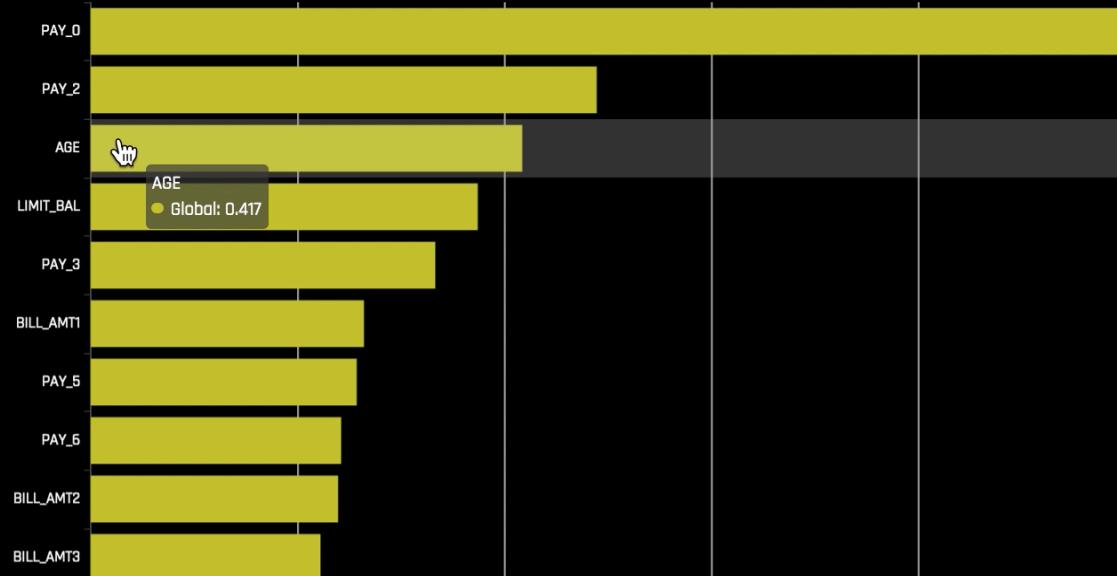
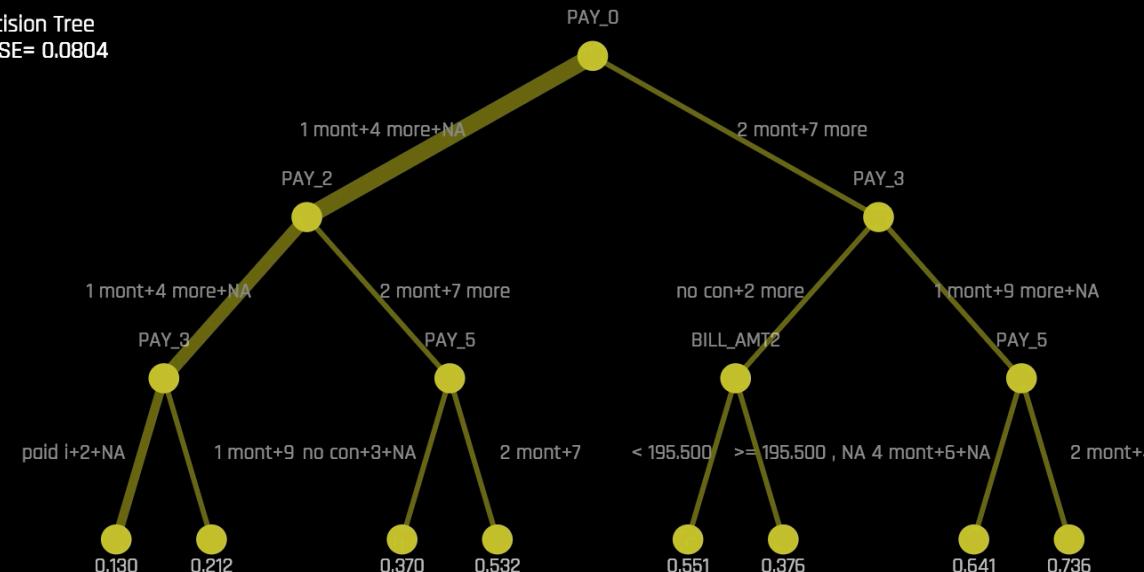
6_Freq_PAY_0	1.00
50_CV_TE_PAY_5_PAY_3_0	0.66
29_TruncSVD_PAY_AMT2_BILL_AMT1_BILL_AMT5_PAY_AMT1...	0.64
32_TruncSVD_BILL_AMT5_PAY_AMT2_LIMIT_BAL_0	0.55
36_TruncSVD_BILL_AMT1_BILL_AMT4_LIMIT_BAL_1	0.43
36_TruncSVD_BILL_AMT1_BILL_AMT4_LIMIT_BAL_0	0.41
46_TruncSVD_PAY_AMT6_BILL_AMT2_0	0.33
23_LIMIT_BAL	0.29
21_BILL_AMT1	0.29
37_TruncSVD_PAY_AMT1_PAY_AMT6_PAY_AMT3_0	0.29
17_PAY_AMT2	0.28
5_CV_TE_PAY_2_0	0.24
16_PAY_AMT1	0.24
44_CV_CatNumEnc_PAY_3_PAY_AMT5_mean	0.23

Global Interpretable Model Explanation Plot

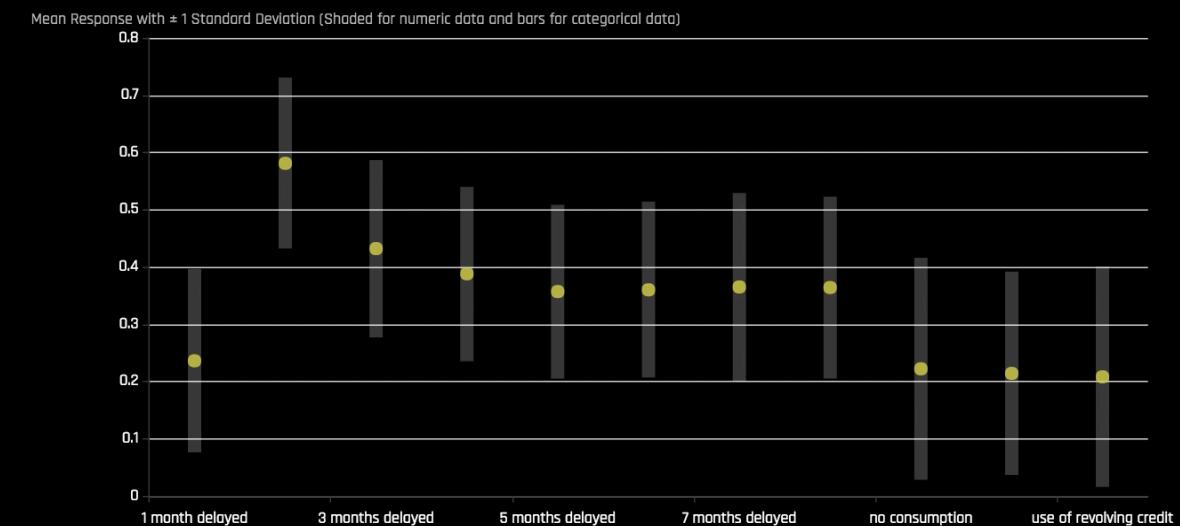
Model Prediction ● k-LIME Model Prediction ● Actual Target



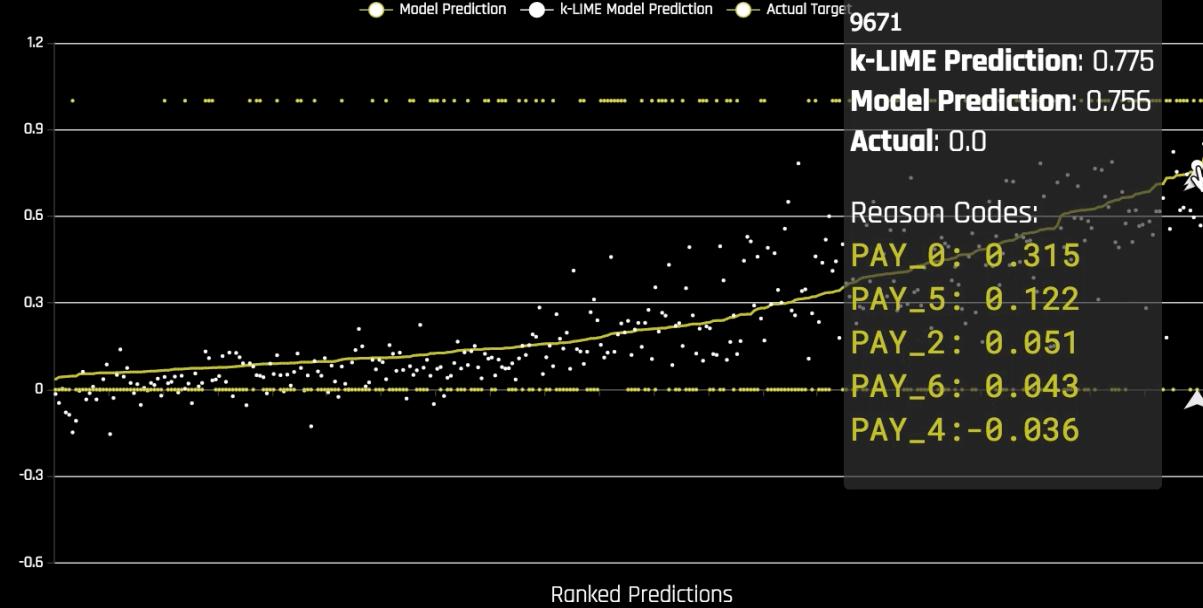
Variable Importance

Decision Tree
RMSE= 0.0804

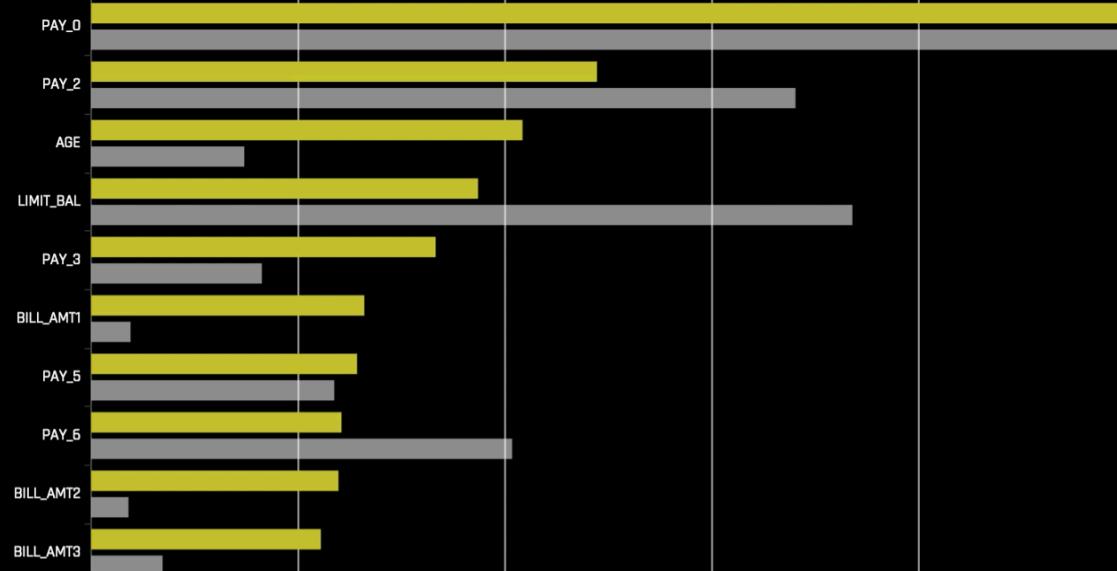
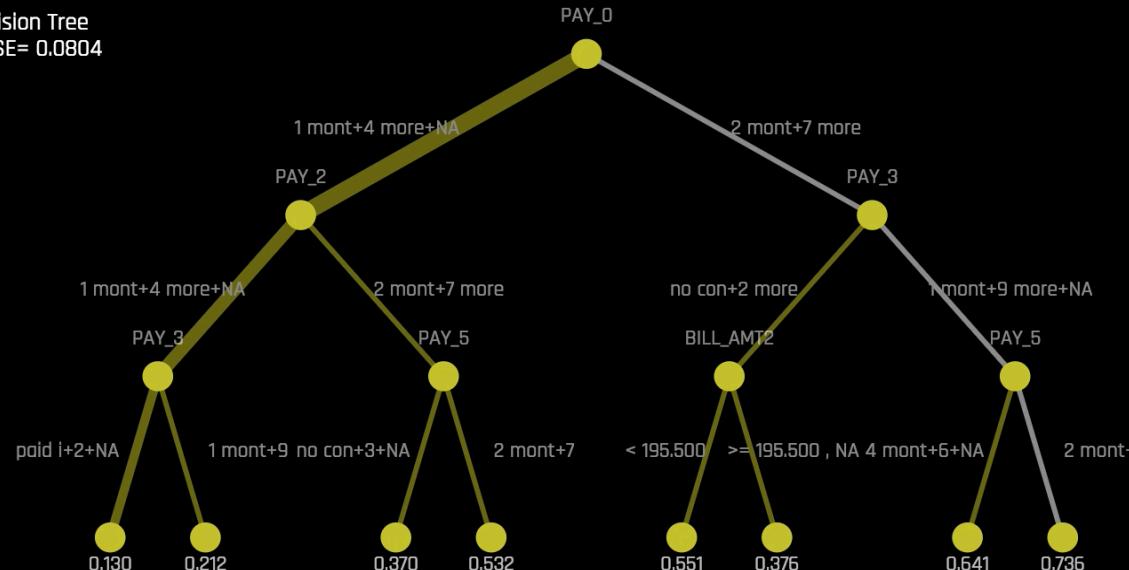
Partial Dependence



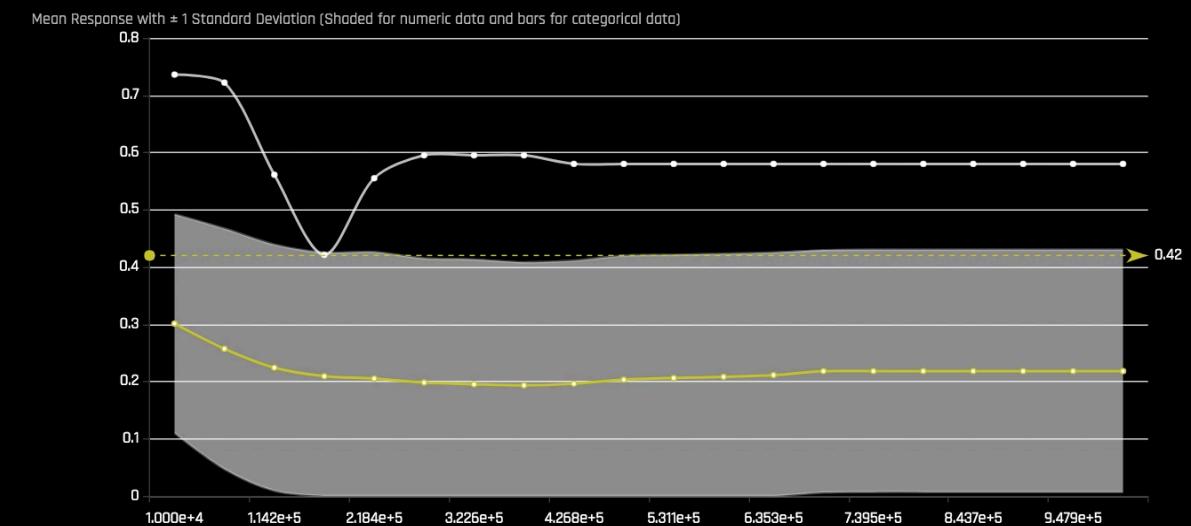
Cluster 0 k-LIME Plot



Variable Importance

Decision Tree
RMSE= 0.0804

Partial Dependence



Actual and Predicted Values

default.payment.next.month (Actual)	0
Model Prediction Value	0.7557
k-LIME Prediction Value	0.7746
k-LIME Prediction Accuracy (%)	97.5%

Local Reason Codes

k-LIME Local Attributions	Variable	with value	is associated with default.payment.next.month
Top Positive Local Attributions			
PAY_0	2	increase of	0.32
PAY_5	7	increase of	0.12
PAY_2	2	increase of	0.05
Skipped 7 additional attributions, click to view all ...			
Top Negative Local Attributions			
PAY_4	7	decrease of	0.04
LIMIT_BAL	170000	decrease of	0.03
MARRIAGE	NoN	decrease of	0.01
Skipped 10 additional attributions, click to view all ...			
k-LIME Cluster Baseline Attribution			0.2801
k-LIME Prediction (k-LIME Local + Cluster Baseline Attribtions)			0.7746

Cluster 0 Reason Codes

k-LIME explains 89.14% in default.payment.next.month for this cluster.

Variable	with value	is associated with default.payment.next.month	
Top Positive Cluster Attributions			
PAY_5	8	increase of	0.413
PAY_2	6	increase of	0.407
PAY_0	2	increase of	0.315

Resources

Machine Learning Interpretability with H2O Driverless AI

<https://www.h2o.ai/wp-content/uploads/2017/09/MLI.pdf>

Ideas on Interpreting Machine Learning

<https://www.oreilly.com/ideas/ideas-on-interpreting-machine-learning>

FAT/ML

<http://www.fatml.org/>

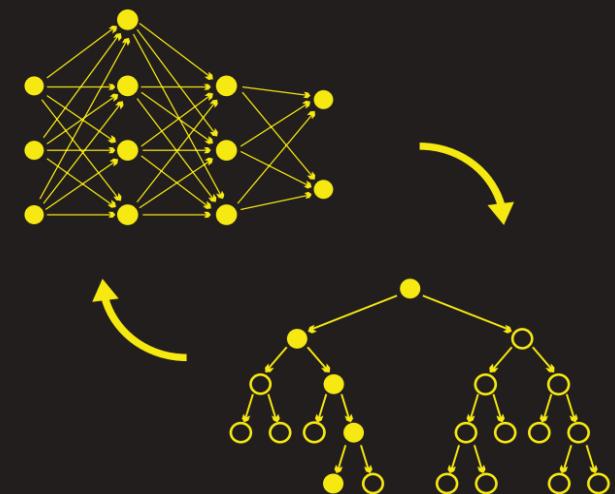
MLI Resources

<https://github.com/h2oai/mli-resources>

MACHINE LEARNING INTERPRETABILITY WITH H2O DRIVERLESS AI

Patrick Hall, Navdeep Gill, Megan Kurka & Wen Phan

Edited by Angela Bartz



Questions?