

# Responsible Machine Learning\*

## Lecture 5: Machine Learning Model Debugging

Patrick Hall

The George Washington University

June 18, 2020

---

\* This material is shared under a [CC By 4.0 license](#) which allows for editing and redistribution, even for commercial purposes. However, any derivative work should attribute the author.

Overview

- 
- 
- 

Why Debug

- 
- 
- 
- 

Methods of Debugging

- 
- 
- 
- 
- 

References

- 

# Contents

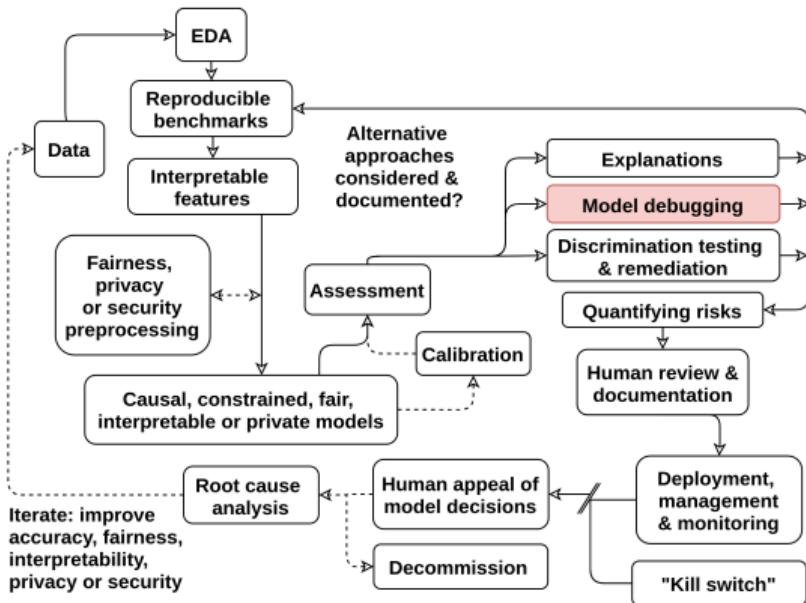
Overview

Why Debug

Methods of Debugging

References

# A Responsible Machine Learning Workflow<sup>†</sup>



<sup>†</sup>A Responsible Machine Learning Workflow

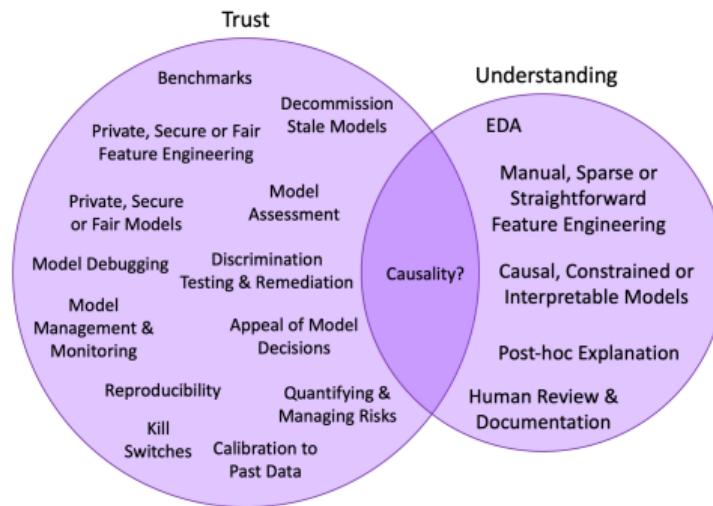
# What is Model Debugging?

- Model debugging is an emergent discipline focused on discovering and remediating errors in the internal mechanisms and outputs of machine learning models.<sup>‡</sup>
- Model debugging attempts to test machine learning models like code (because the models are code).
- Model debugging promotes trust directly and enhances interpretability as a side-effect.
- Model debugging is similar to regression diagnostics, but for nonlinear machine learning models.

---

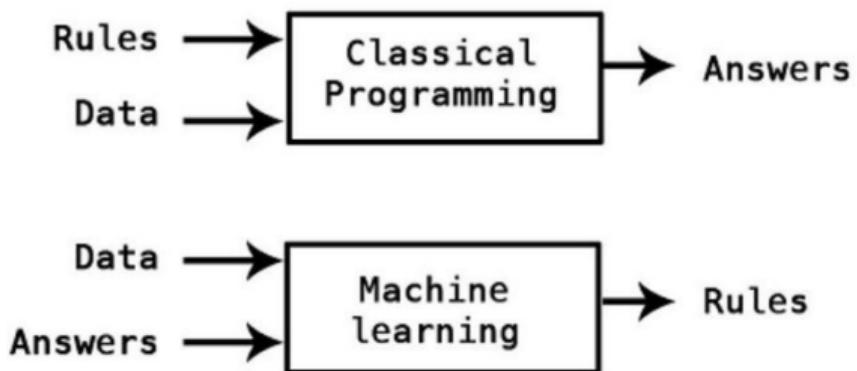
<sup>‡</sup>See <https://debug-ml-iclr2019.github.io/> for numerous model debugging approaches.

# Trust and Understanding



Trust and understanding in machine learning are different but complementary goals, and they are technically feasible *today*.

## Reminder: What is Machine Learning?



From: [Deep Learning with Python](#).

# Why Debug?

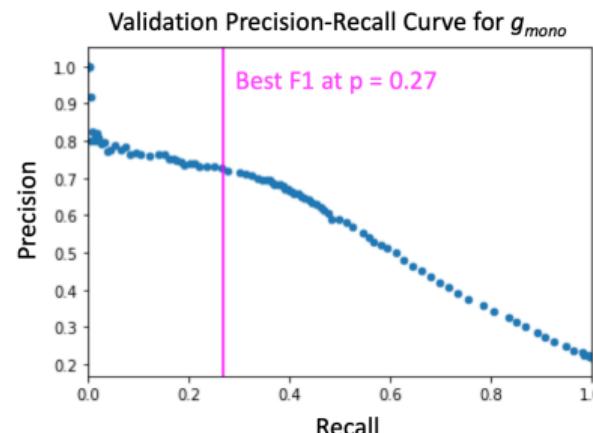
- Machine learning models can discriminate, leak private data, get hacked or make silly mistakes.
- Machine learning pipelines are not magically immune from the bugs and security vulnerabilities that affect all other software.
- AI incidents can cause harm to organizations, consumers, and the public.

*“If builders built houses the way programmers built programs, the first woodpecker to come along would destroy civilization.”*

— Gerald Weinberg

# Why Debug: Inadequate Assessment

- Constrained, monotonic GBM probability of default (PD) classifier,  $g_{mono}$ .
- Grid search over hundreds of models.
- Best model selected by validation-based early stopping.
- Seemingly well-regularized (row and column sampling, explicit specification of L1 and L2 penalties).
- No evidence of over- or under-fitting.
- Better validation logloss than benchmark GLM.
- Decision threshold selected by maximization of F1 statistic.
- BUT traditional assessment can be inadequate ...



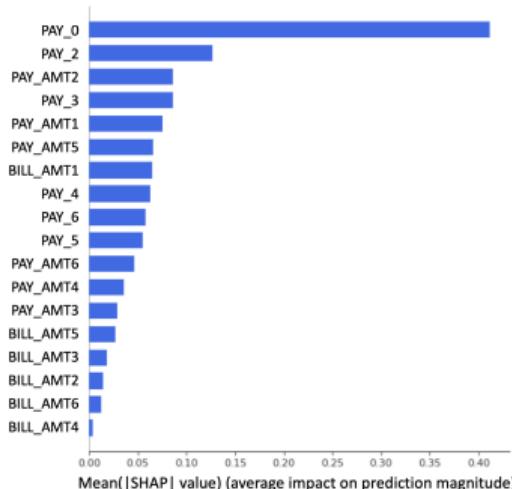
Validation Confusion Matrix at Threshold:

	Actual: 1	Actual: 0
Predicted: 1	1159	827
Predicted: 0	1064	6004

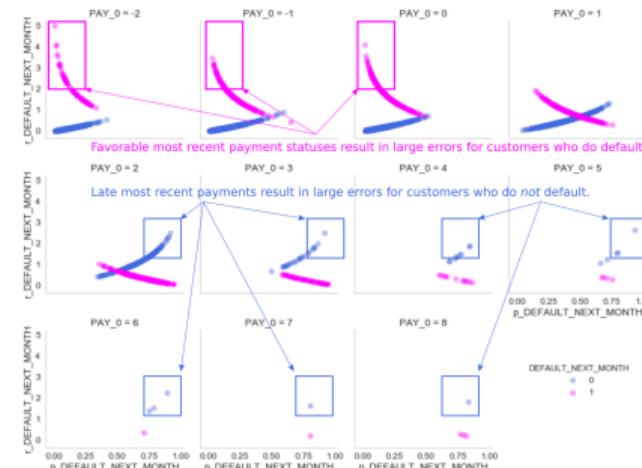


## Why Debug: Inaccuracy

Machine learning models can be **inaccurate**.



**gmono** PD classifier over-emphasizes the most important feature, a customer's most recent repayment status,  $PAY\_0$ .



**gmono** also struggles to predict default for favorable statuses,  $-2 \leq PAY\_0 < 2$ , and often cannot predict on-time payment when recent payments are late,  $PAY\_0 \geq 2$ .

## Why Debug: Discrimination

Machine learning models can perpetuate **sociological discrimination** [barocas-hardt-narayanan].

	Adverse Impact Disparity	Accuracy Disparity	True Positive Rate Disparity	Precision Disparity	Specificity Disparity
<b>single</b>	0.885	1.029	0.988	1.008	1.025
<b>divorced</b>	1.014	0.932	0.809	0.806	0.958
<b>other</b>	0.262	1.123	0.62	1.854	1.169

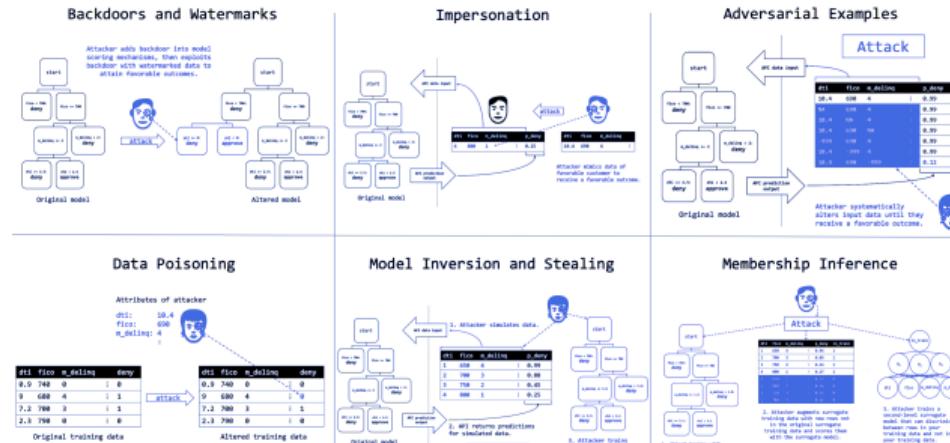
Group disparity metrics are out-of-range for  $g_{mono}$  across different marital statuses.



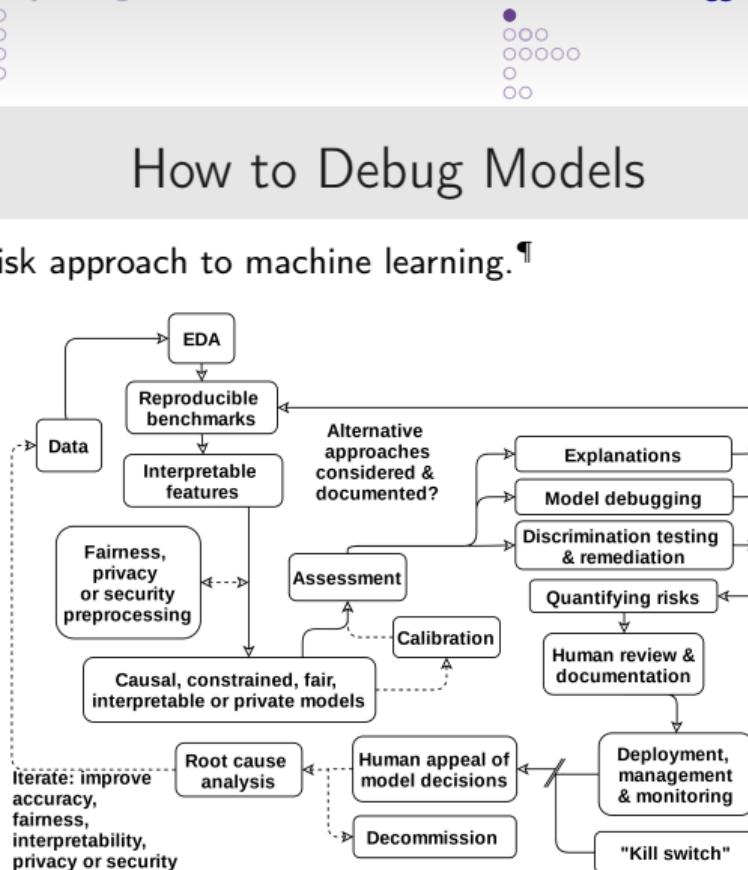
# Why Debug: Security Vulnerabilities

Machine learning models have **security vulnerabilities** [[security\\_of\\_ml](#)], [[membership\\_inference](#)], [[model\\_stealing](#)].<sup>§</sup>

## Machine Learning Attack Cheatsheet



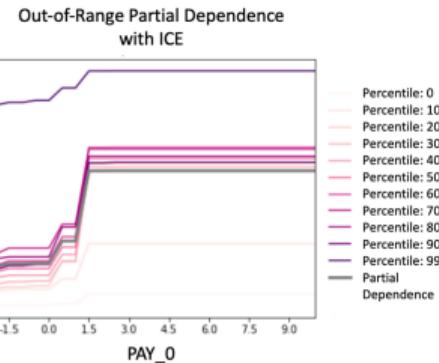
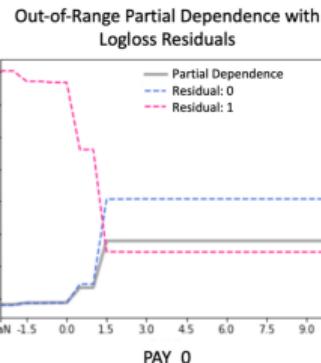
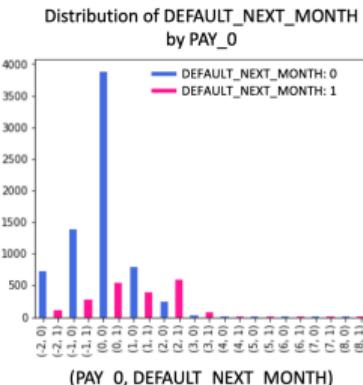
<sup>§</sup> See [https://github.com/jphall663/secure\\_ML\\_ideas](https://github.com/jphall663/secure_ML_ideas) for full size image and more information.



As part of a holistic, low-risk approach to machine learning. 

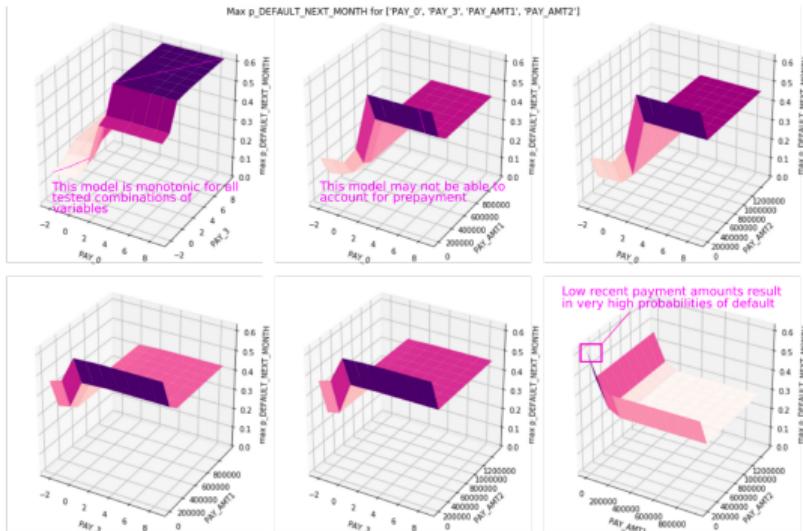
¶ See [https://github.com/iphall1663/hc\\_ml](https://github.com/iphall1663/hc_ml) for more information.

# Sensitivity Analysis: Partial Dependence and ICE



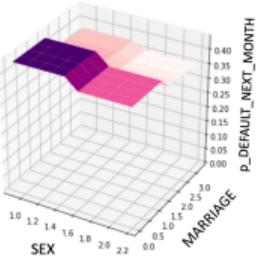
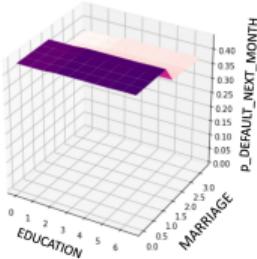
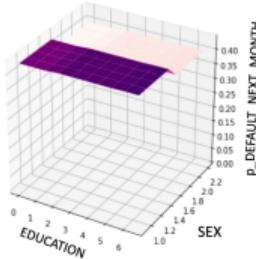
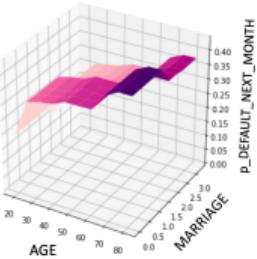
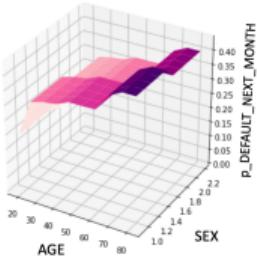
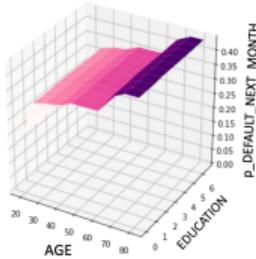
- Training data is very sparse for  $PAY\_0 > 2$ .
- Residuals of partial dependence confirm over-emphasis on  $PAY\_0$ .
- ICE curves indicate that partial dependence is likely trustworthy and empirically confirm monotonicity, but also expose adversarial attack vulnerabilities.
- Partial dependence and ICE indicate  $g_{mono}$  likely learned very little for  $PAY\_0 \geq 2$ .
- $PAY\_0 = \text{missing}$  gives lowest probability of default.

# Sensitivity Analysis: Search for Adversarial Examples



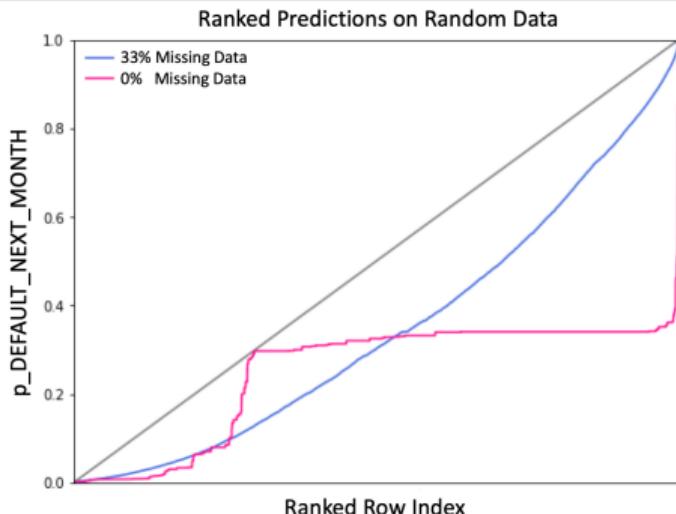
Adversary search confirms multiple avenues of attack and exposes a potential flaw in `gmono` scoring logic: default is predicted for customer's who make payments above their credit limit. (Try the `cleverhans` package for finding adversarial examples.)

# Sensitivity Analysis: Search for Adversarial Examples



`gmono` appears to prefer younger, unmarried customers, which should be investigated further with disparate impact analysis - see slide 8 - and could expose the lender to impersonation attacks. (Try the [AIF360](#), [aequitas](#), or [themis](#) packages for disparate impact audits.)

# Sensitivity Analysis: Random Attacks



- In general, random attacks are a viable method to identify software bugs in machine learning pipelines. **(Start here if you don't know where to start.)**
- Random data can apparently elicit all probabilities  $\in [0, 1]$  from  $g_{\text{mono}}$ .
- Around the decision threshold, lower probabilities can be attained simply by injecting missing values, yet another vulnerability to adversarial attack.

# Residual Analysis: Disparate Accuracy and Errors

Error Metrics for PAY\_0

	Prevalence	Accuracy	True Positive Rate	Precision	Specificity	Negative Predicted Value	False Positive Rate	False Discovery Rate	False Negative Rate	False Omissions Rate
<b>PAY_0</b>										
-2	0.124	0.864	0.099	0.333	0.972	0.884	0.028	0.667	0.901	0.116
-1	0.168	0.816	0.206	0.406	0.939	0.854	0.061	0.594	0.794	0.146
0	0.121	0.867	0.107	0.341	0.972	0.888	0.028	0.659	0.893	0.112
1	0.325	0.491	0.903	0.361	0.292	0.862	0.708	0.619	0.097	0.138
2	0.709	0.709	1	0.709	0	0.5	1	0.291	0	0.5
3	0.748	0.748	1	0.748	0	0.5	1	0.252	0	0.5
4	0.571	0.571	1	0.571	0	0.5	1	0.429	0	0.5
5	0.444	0.444	1	0.444	0	0.5	1	0.556	0	0.5
6	0.25	0.25	1	0.25	0	0.5	1	0.75	0	0.5
7	0.5	0.5	1	0.5	0	0.5	1	0.5	0	0.5
8	0.75	0.75	1	0.75	0	0.5	1	0.25	0	0.5

Error Metrics for SEX

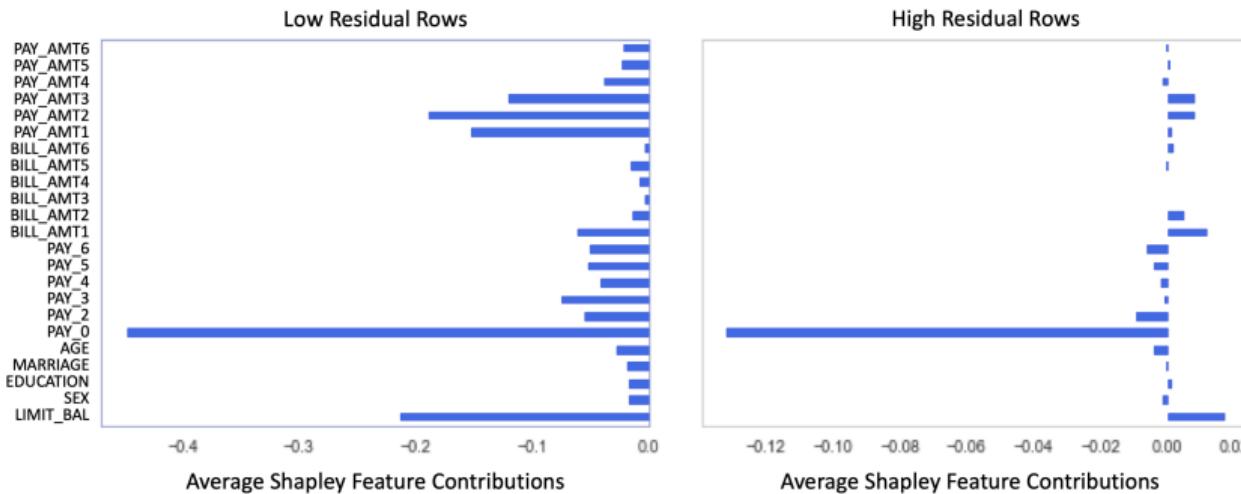
	Prevalence	Accuracy	True Positive Rate	Precision	Specificity	Negative Predicted Value	False Positive Rate	False Discovery Rate	False Negative Rate	False Omissions Rate
<b>SEX</b>										
Male	0.235	0.782	0.626	0.531	0.83	0.879	0.17	0.469	0.374	0.121
Female	0.209	0.797	0.552	0.514	0.862	0.879	0.138	0.486	0.448	0.121

For PAY\_0:

- Notable change in accuracy and error characteristics for PAY\_0  $\geq 2$ .
- 100% false omission rate for PAY\_0  $\geq 2$ . (Every prediction of non-default is incorrect!)

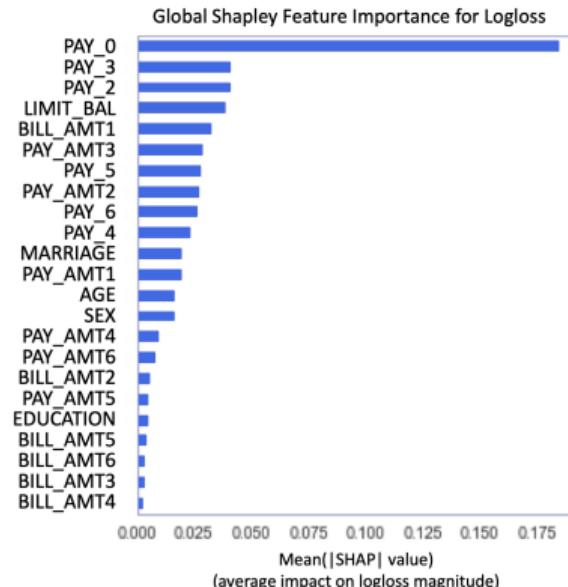
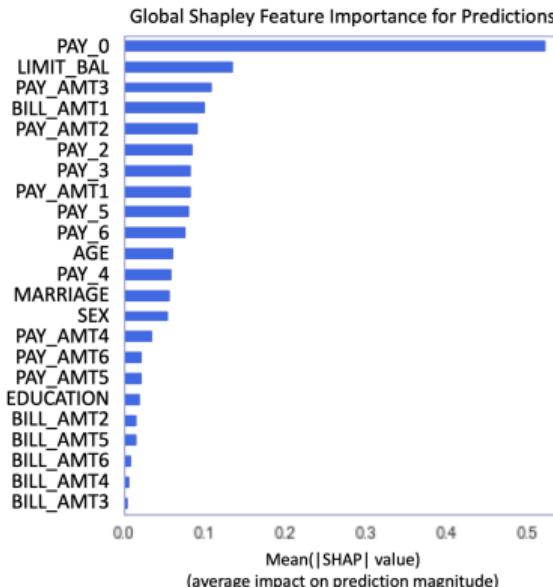
For SEX, accuracy and error characteristics vary little across individuals represented in the training data. Non-discrimination should be confirmed by disparate impact analysis.

# Residual Analysis: Mean Local Feature Contributions



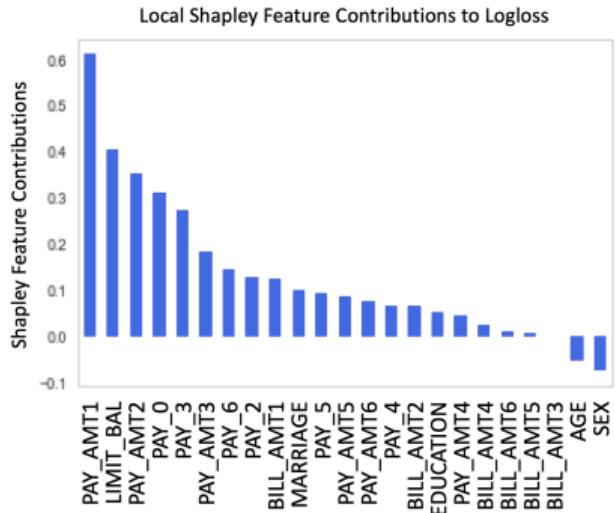
Exact local Shapley feature contributions [[shapley](#)], which are available at scoring time for unlabeled data, are noticeably different for low and high residual predictions. (Both monotonicity constraints and Shapley values are available in [h2o-3](#) and [XGBoost](#).)

## Residual Analysis: Non-Robust Features



Globally important features PAY\_3 and PAY\_2 are more important, on average, to the loss than to the predictions. (Shapley contributions to XGBoost logloss can be calculated using the [shap](#) package. This is a **time-consuming** calculation.)

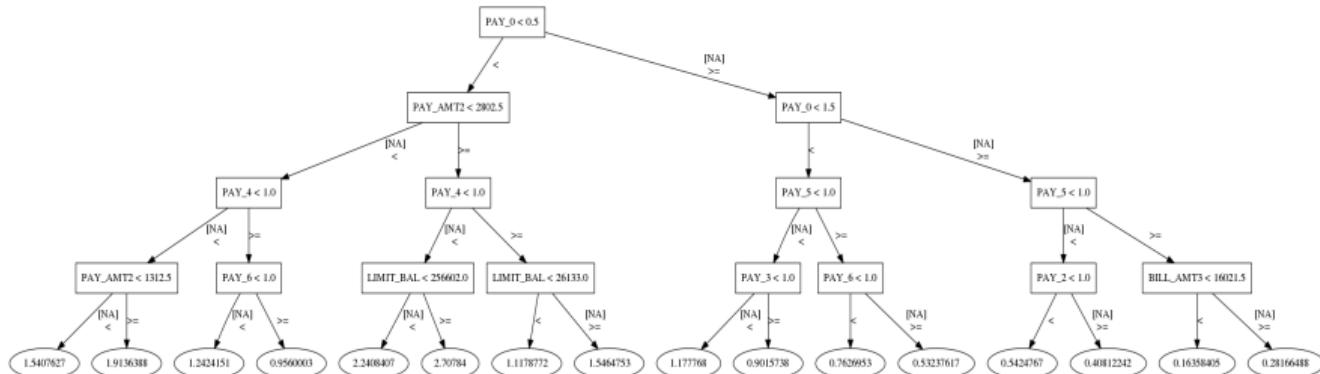
# Residual Analysis: Local Contributions to Logloss



Exact, local feature contributions to logloss can be calculated, enabling ranking of features contributing to logloss residuals for **each prediction**.

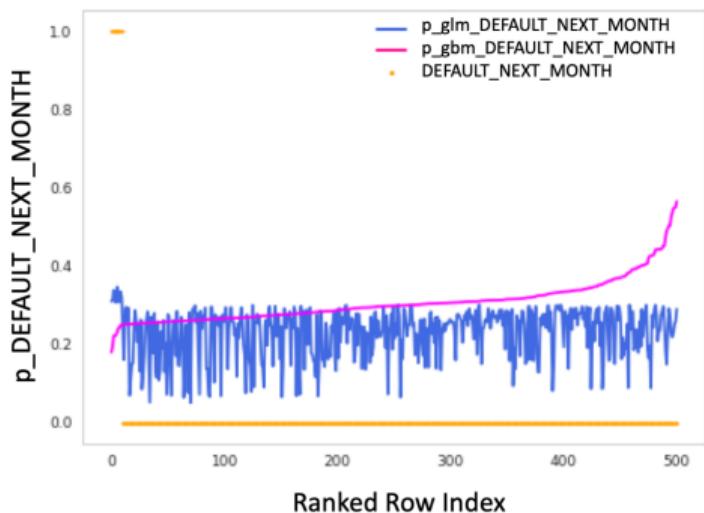
# Residual Analysis: Modeling Residuals

Decision tree model of  $g_{\text{mono}}$   $\text{DEFAULT\_NEXT\_MONTH} = 1$  logloss residuals with 3-fold CV MSE = 0.0070 and  $R^2 = 0.8871$ .



This tree encodes rules describing when  $g_{\text{mono}}$  is probably wrong.

## Benchmark Models: Compare to Linear Models



For a range of probabilities  $\in (\sim 0.2, \sim 0.6)$ ,  $g_{\text{mono}}$  displays exactly incorrect prediction behavior as compared to a benchmark GLM.

# Remediation: for $g_{\text{mono}}$

- **Over-emphasis of PAY\_0:**
  - Engineer features for payment trends or stability.
  - Missing value injection during training or scoring.
- **Sparsity of PAY\_0 > 2 training data:** Increase observation weights.
- **Payments  $\geq$  credit limit:** Scoring-time model assertion [[kangdebugging](#)].
- **Disparate impact:** Adversarial de-biasing [[zhang2018mitigating](#)] or model selection by minimal disparate impact.
- **Security vulnerabilities:** API throttling, authentication, real-time model monitoring.
- **Large logloss importance:** Evaluate dropping non-robust features.
- **Poor accuracy vs. benchmark GLM:** Blend  $g_{\text{mono}}$  and GLM for probabilities  $\in (\sim 0.2, \sim 0.6)$ .
- **Miscellaneous strategies:**
  - Local feature importance and decision tree rules can indicate additional scoring-time model assertions, e.g. alternate treatment for locally non-robust features in known high-residual ranges of the learned response function.
  - Incorporate local feature contributions to logloss into training or scoring processes.

## Remediation: General Strategies

- Calibration to past data.
- Data collection or simulation for model blindspots.
- Detection and elimination of non-robust features.
- Missing value injection during training or scoring.
- Model or model artifact editing.

## Acknowledgments

Thanks to Lisa Song for her continued assistance in developing these course materials.

Some materials ©Patrick Hall and the H2O.ai team 2017-2020.

## References

Advanced code examples for this presentation:

[https://www.github.com/jphall663/interpretable\\_machine\\_learning\\_with\\_python](https://www.github.com/jphall663/interpretable_machine_learning_with_python)

[https://www.github.com/jphall663/responsible\\_xai](https://www.github.com/jphall663/responsible_xai)