

SHAP: Interpreting ML Models with IML



March 7, 2020
SatRday



Who am I?



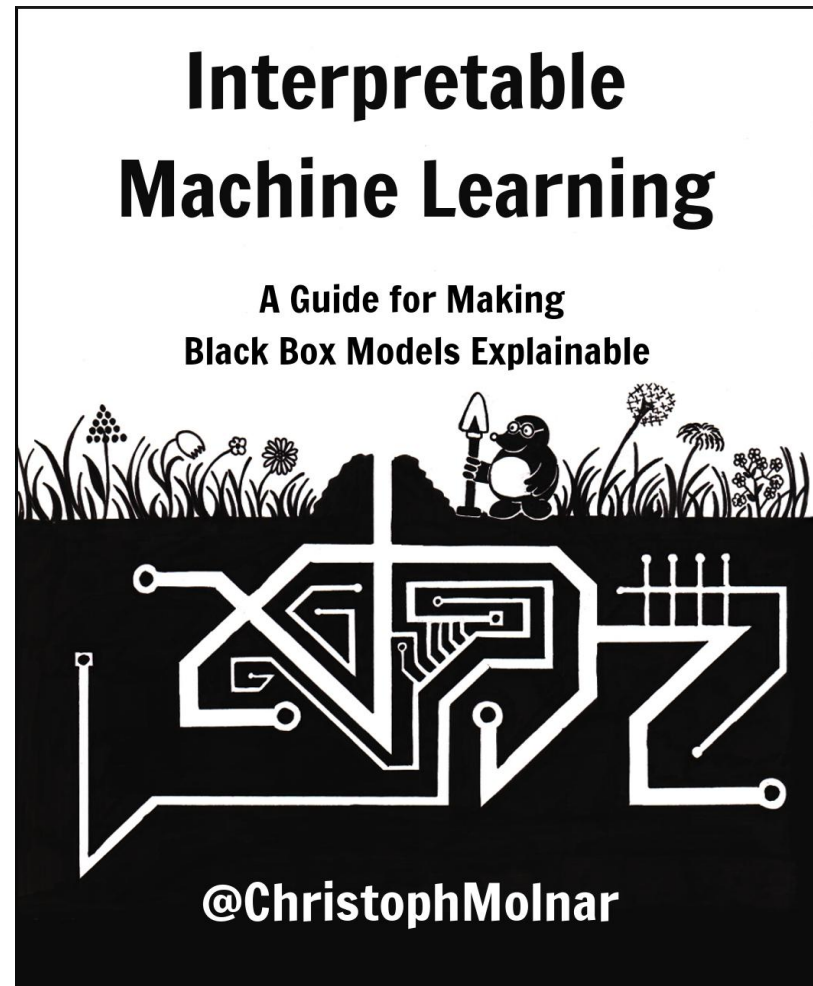
Drikus du Toit

Data Scientist
Decision Science
Capitec Bank



Resource

Molnar, Christoph. "Interpretable machine learning. A Guide for Making Black Box Models Explainable", 2019.
<https://christophm.github.io/interpretable-ml-book/>.



What is ML interpretability and why is it important?

The Client

Applies for loan, gets rejected.

Client Questions:

- Why did I not get the loan?
- What should I do to improve my credit score?

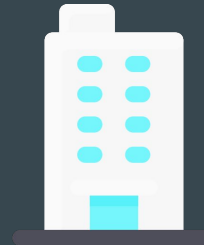


The Business

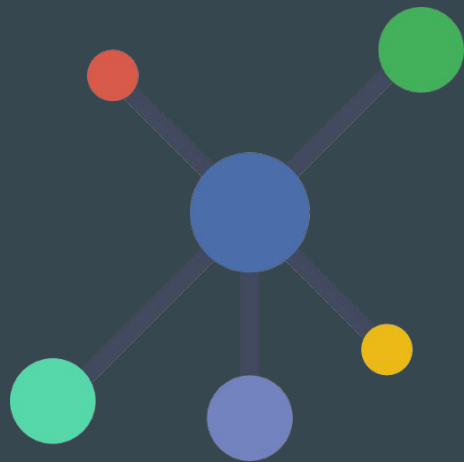
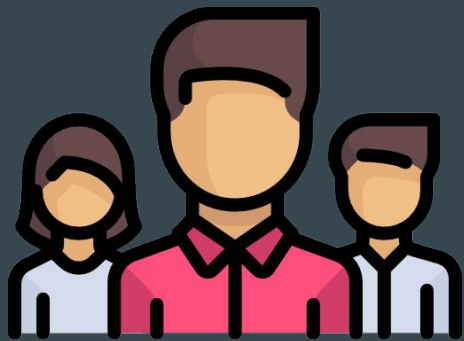
Build a credit default model. High accuracy usually goes with high complexity.

Business Questions:

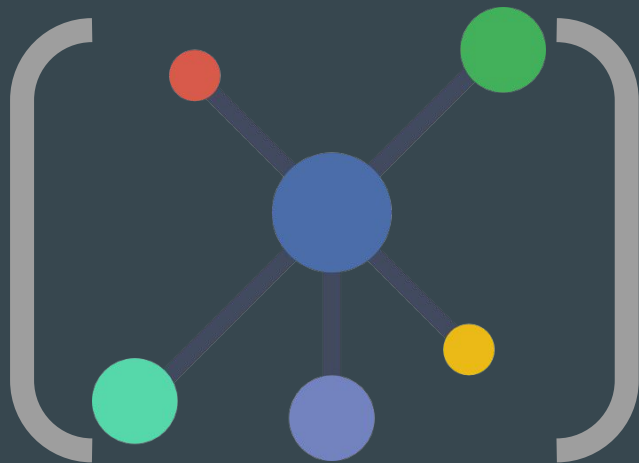
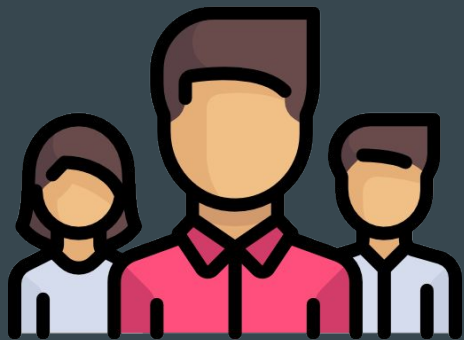
- Is there bias in our model?
- Do we understand our underlying data?
- What will cause the model to not perform as expected?
- Are we within regulatory framework?



Context..



Context..



Context..

Interpretable

Accurate

Complex Model

Simple Model

Context..

Interpretable

Accurate

Complex Model

Simple Model



Context..

Interpretable


Accurate

Complex Model






Simple Model

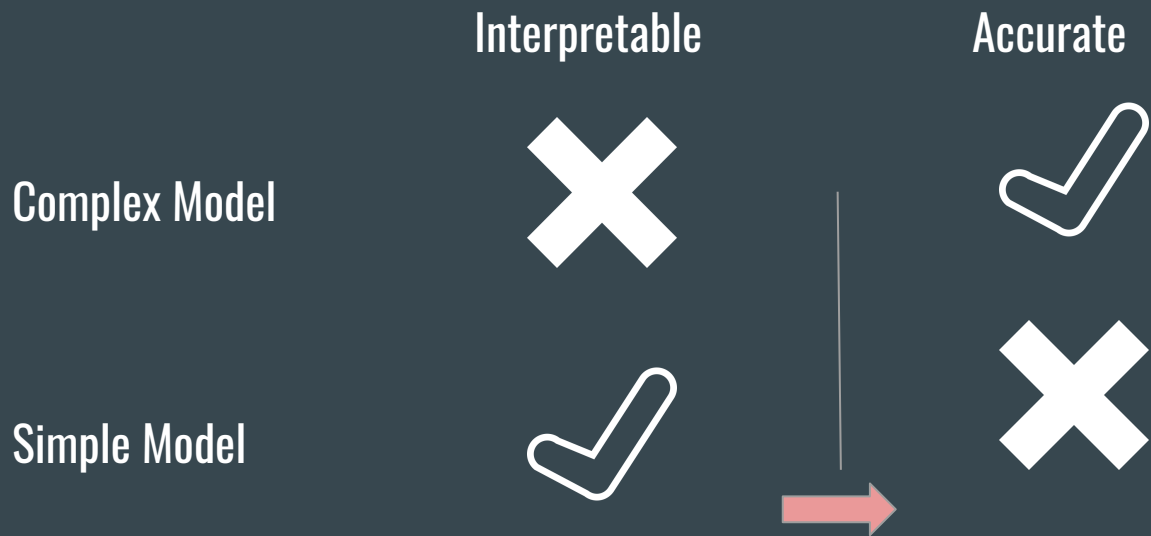
Context..

	Interpretable	Accurate
Complex Model		
Simple Model		

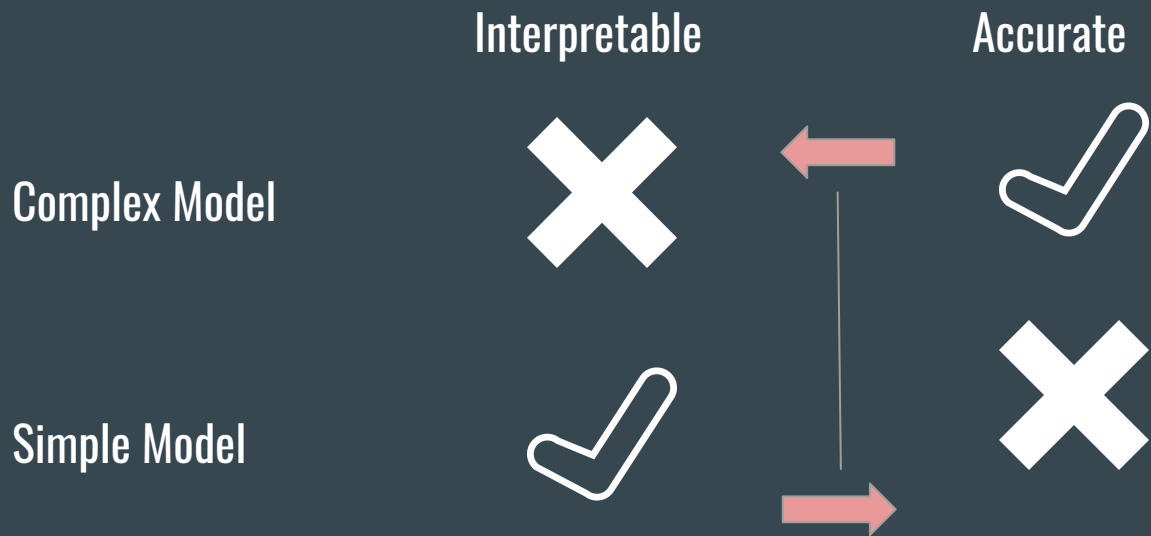
Context..

	Interpretable	Accurate
Complex Model		
Simple Model		

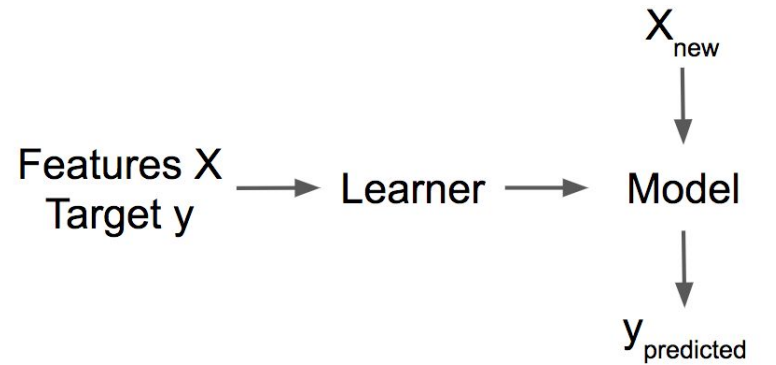
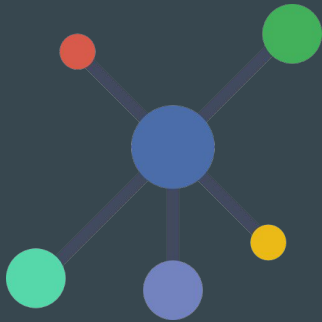
Context..



Context..

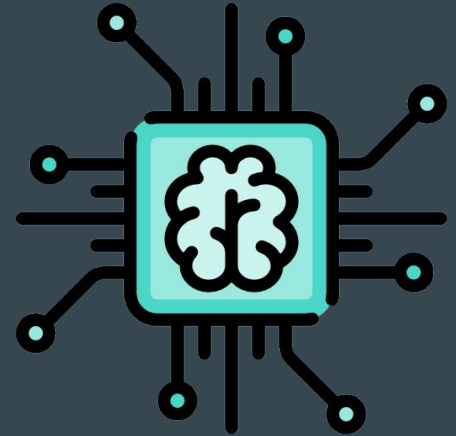


Machine Learning Model



Importance of Interpretability (what or why)

1. **Human curiosity and learning**
2. **Goal of science**
3. **Safety measures**
4. **Detecting bias**
5. **Manage social interactions**
6. **Debugged and audited**



Interpretability Techniques



- **Reveal its internal mechanisms**
- Fully understood by looking at their parameters
- Also called interpretable models



- **Does not reveal its internal mechanisms**
- Cannot be understood by looking at their parameters (e.g. a neural network)

Black Box Models (interpretability techniques)



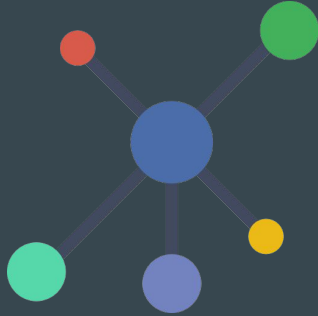
Model Agnostic Methods	Example Based Explanations
Partial Dependence Plot (PDP)	Counterfactual Explanations
Individual Conditional Expectation (ICE)	Adversarial Examples
Accumulated Local Effects (ALE) Plot	Prototypes and Criticisms
Feature Interaction	Influential Instances
Permutation Feature Importance	
Global Surrogate	
Local Surrogate (LIME)	
Scoped Rules (Anchors)	
Shapley Values	
SHAP (SHapley Additive exPlanations)	

Black Box Models (interpretability techniques)



Model Agnostic Methods	Example Based Explanations
Partial Dependence Plot (PDP)	Counterfactual Explanations
Individual Conditional Expectation (ICE)	Adversarial Examples
Accumulated Local Effects (ALE) Plot	Prototypes and Criticisms
Feature Interaction	Influential Instances
Permutation Feature Importance	
Global Surrogate	
Local Surrogate (LIME)	
Scoped Rules (Anchors)	
Shapley Values	
SHAP (SHapley Additive exPlanations)	

Shapley Values

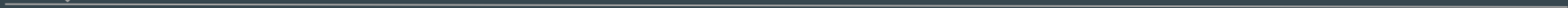


23 %





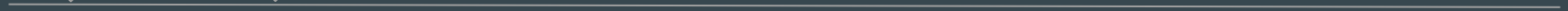
0





Average Pred
15%

0





Average Pred
15%
 $E[f(x)]$

Joe
23%
 $f(x)$

0



?



Average Pred
15%
 $E[f(x)]$

0



Base
 $\Phi(0)$



Average Pred
15%
 $E[f(x)]$

0

19%



Base
 $\Phi(0)$

Income not verified
 $\Phi(1)$



Average Pred
15%
 $E[f(x)]$

0

19%

21%



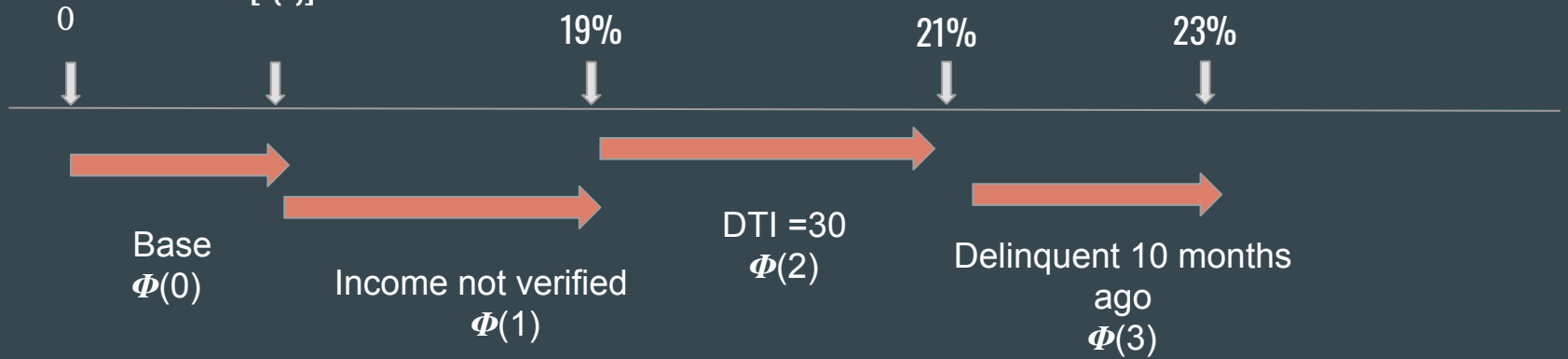
Base
 $\Phi(0)$

Income not verified
 $\Phi(1)$

DTI = 30
 $\Phi(2)$

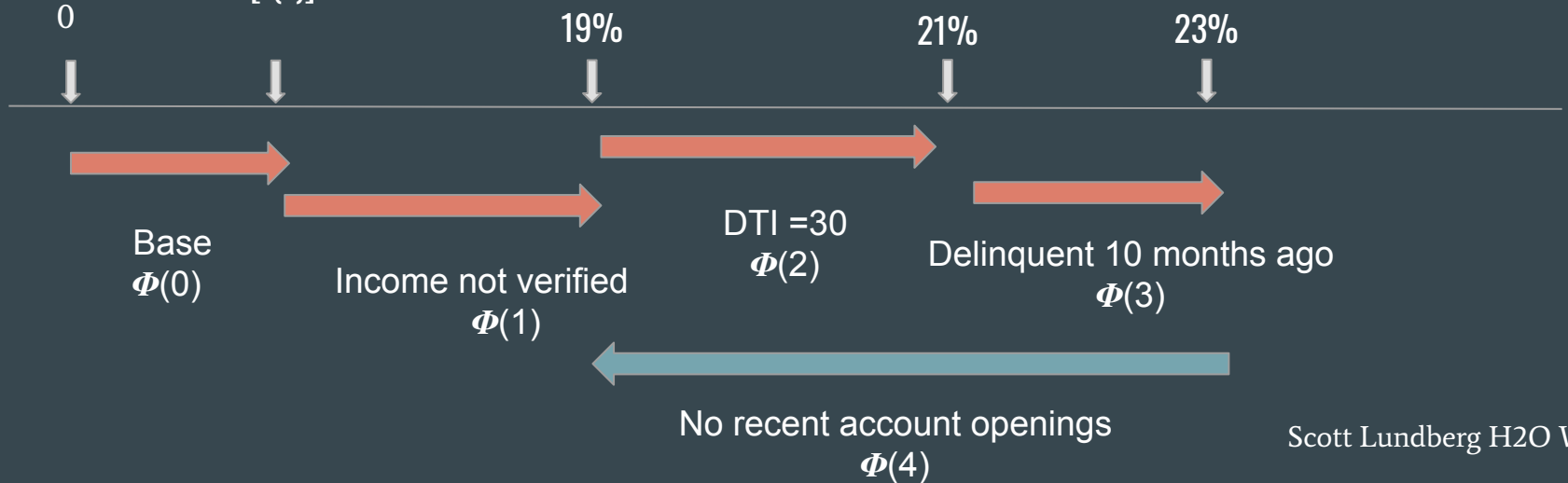


Average Pred
15%
 $E[f(x)]$



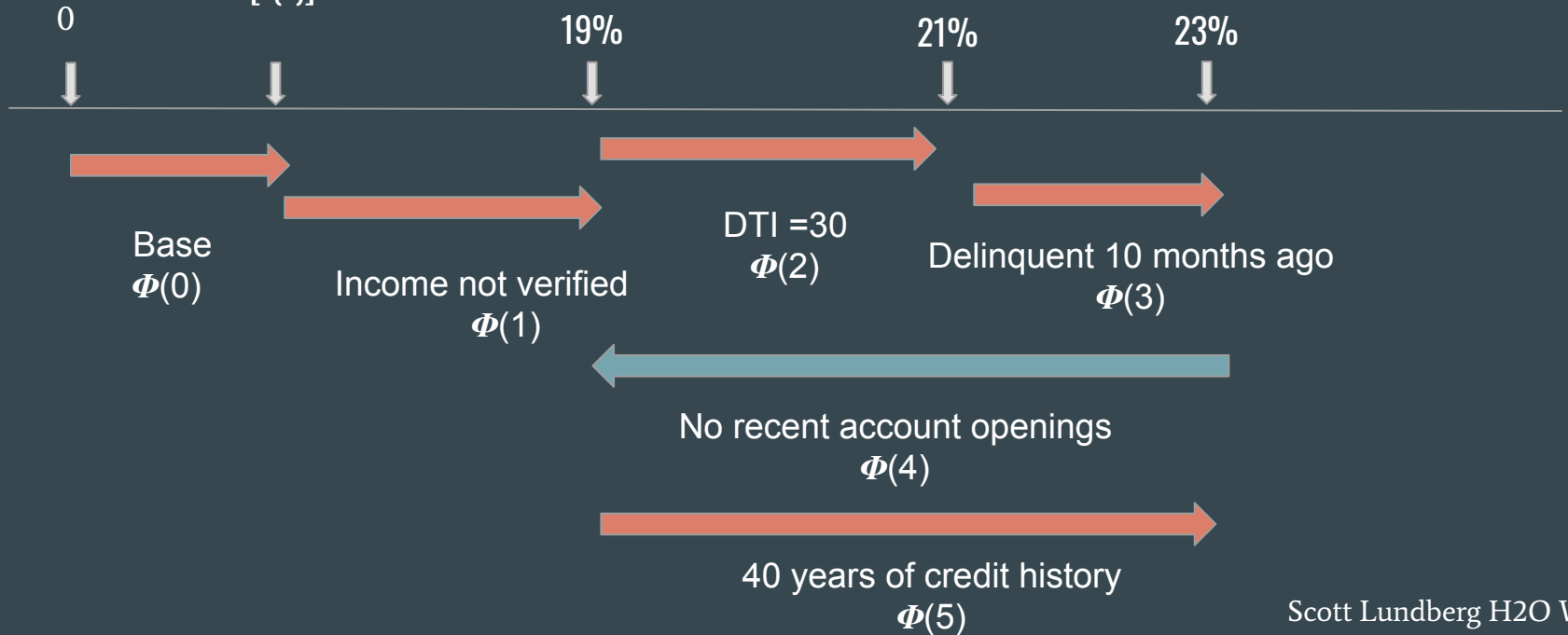


Average Pred
15%
 $E[f(x)]$



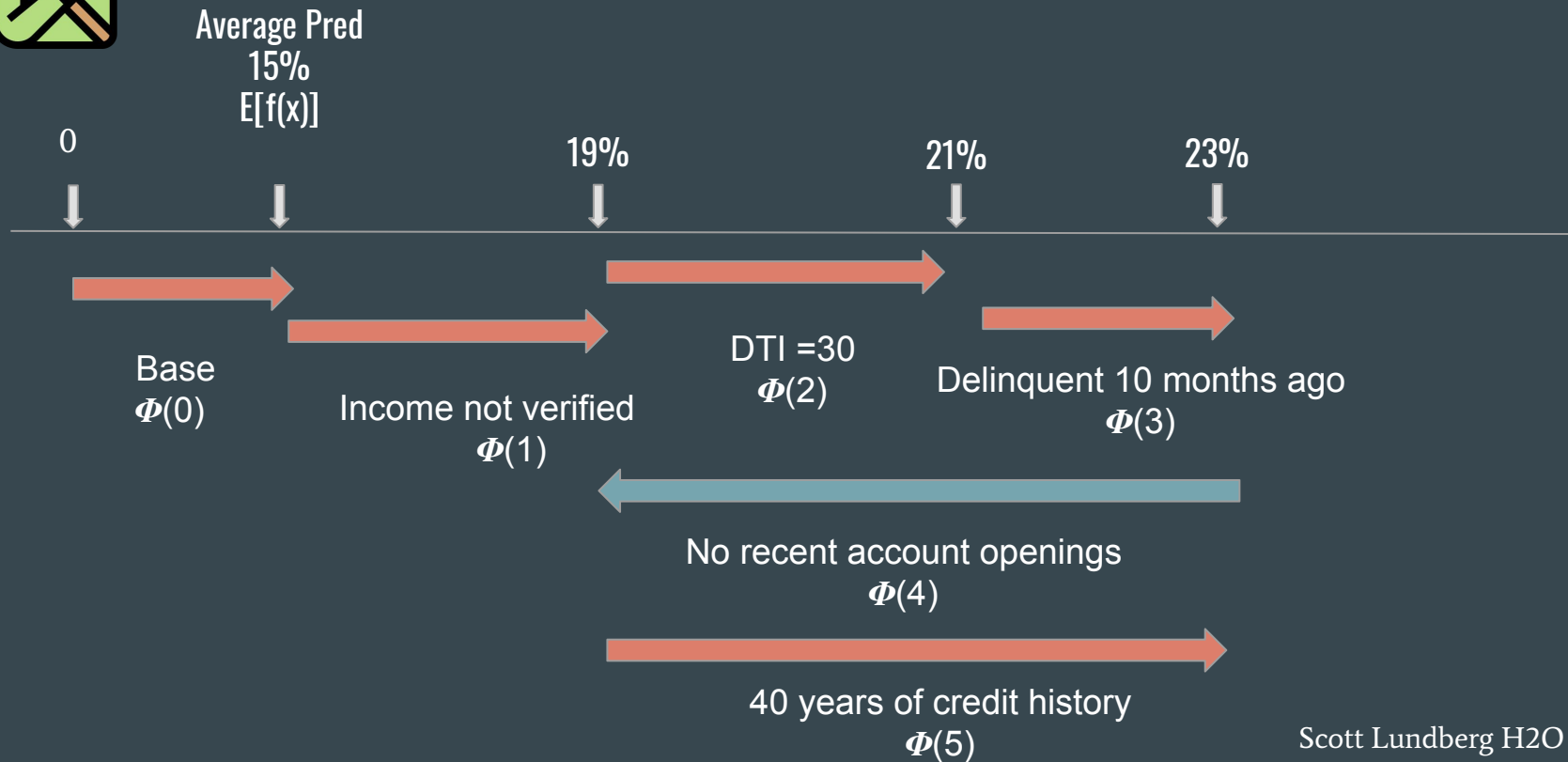


Average Pred
15%
 $E[f(x)]$



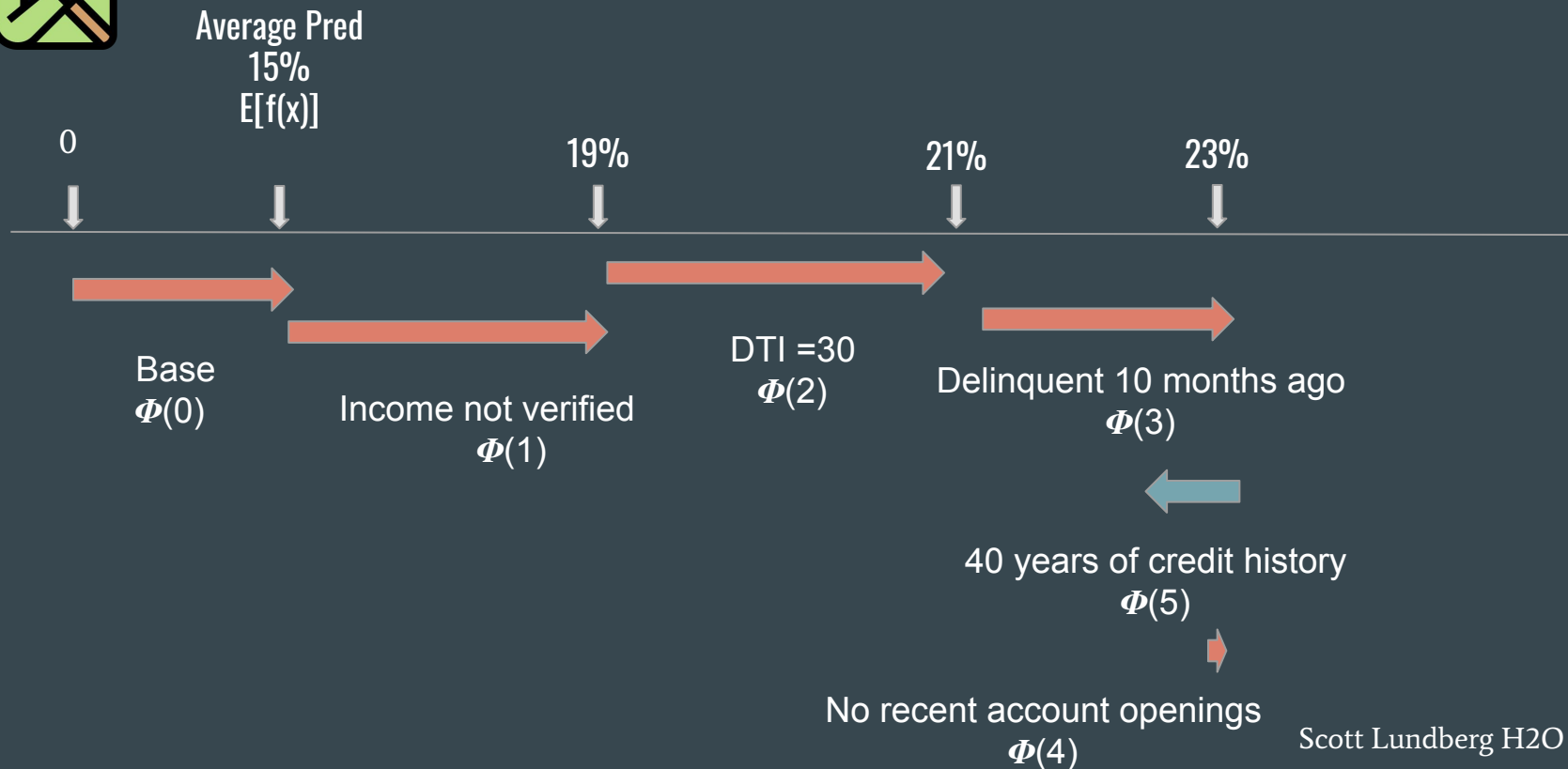


What about the order?



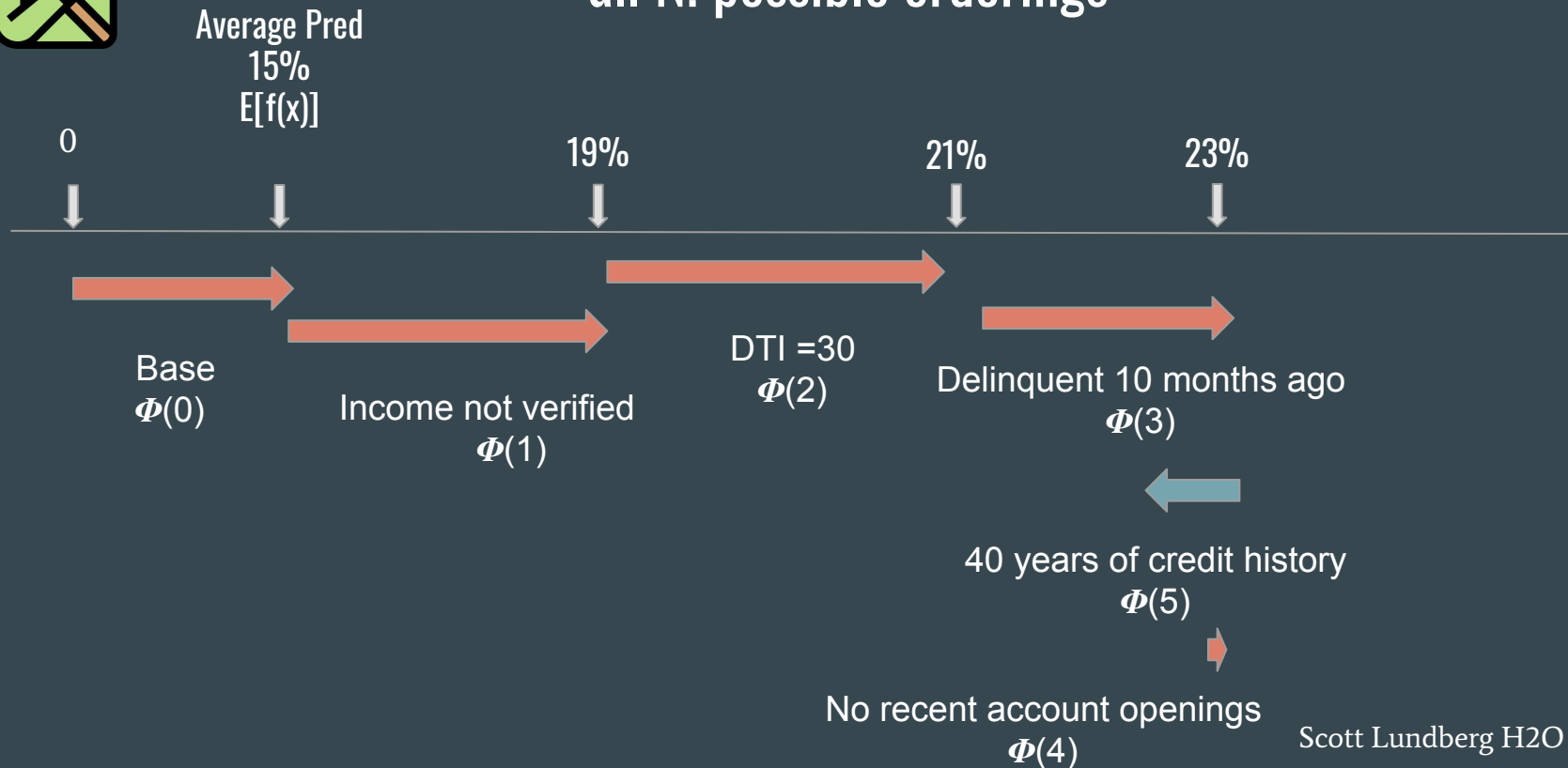


What about the order?





Shapley values results from averaging over all $N!$ possible orderings



Examples & Interpretation

```
iml_shapley.Rmd
1
2 # First we fit a machine learning model on the Boston housing data
3
4 set.seed(42)
5 library("iml")
6 library("randomForest")
7
8 data("Boston", package = "MASS")
9 rf = randomForest(medv ~ ., data = Boston, ntree = 50)
10 X = Boston[-which(names(Boston) == "medv")]
11 mod = Predictor$new(rf, data = X)
12
13 # Then we explain the first instance of the dataset with the Shapley method:
14
15 x.interest = X[1,]
16 shapley = Shapley$new(mod, x.interest = x.interest)
17
18 # plot
19
20 plot(shapley)
21
```

Environment History Connections

Files Plots Packages Help Viewer

Zoom Export Publish

Actual prediction: 25.75
Average prediction: 22.56

feature.value	phi
lstat=4.98	3.5
indus=2.31	1.0
ptratio=15.3	0.8
dis=4.09	0.4
black=396.9	0.0
chas=0	0.0
zn=18	0.05
tax=296	-0.2
nox=0.538	-0.3
age=65.2	-0.5
rad=1	-0.7
crim=0.00632	-0.8
rm=6.575	-1.2

21:1 (Top Level) R Markdown

Console

SHAP (Shapley Additive Explanations)

KernelSHAP

An alternative, kernel-based estimation approach for Shapley values inspired by local surrogate models

TreeSHAP

An efficient estimation approach for tree-based models

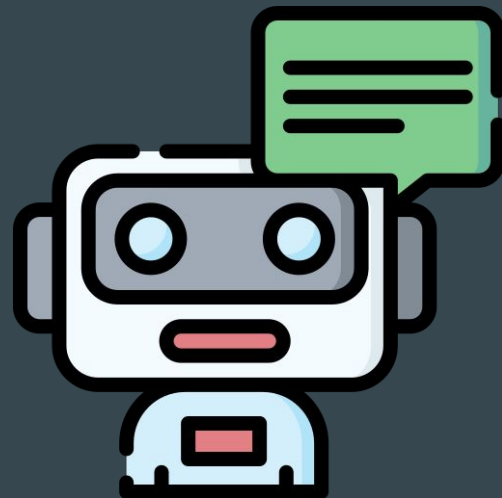
SHAP comes with many global interpretation methods based on aggregations of Shapley values

The future of interpretability

The focus will be on model-agnostic interpretability tools

Robots and programs will explain themselves

Ethical Issues



Resource

Miller, Tim. "Explanation in artificial intelligence: Insights from the social sciences." arXiv Preprint arXiv:1706.07269. (2017). [2]

Kim, Been, Rajiv Khanna, and Oluwasanmi O. Koyejo. "Examples are not enough, learn to criticize! Criticism for interpretability." Advances in Neural Information Processing Systems (2016). [3]

Doshi-Velez, Finale, and Been Kim. "Towards a rigorous science of interpretable machine learning," no. MI: 1–13. <http://arxiv.org/abs/1702.08608> (2017)

Package iml

<https://www.youtube.com/watch?v=ngOBhhINWb8>

Molnar, Christoph. "Interpretable machine learning. A Guide for Making Black Box Models Explainable", 2019. <https://christophm.github.io/interpretable-ml-book/>.

Lundberg, Scott M., and Su-In Lee. "A unified approach to interpreting model predictions." Advances in Neural Information Processing Systems. 2017.

Resource

Icons:

Icons made by [Smashicons](https://www.flaticon.com/authors/smashicons "Smashicons") from [www.flaticon.com](https://www.flaticon.com/ "Flaticon")

Icons made by [Freepik](https://www.flaticon.com/authors/freepik "Freepik") from [www.flaticon.com](https://www.flaticon.com/ "Flaticon")

Icons made by [monkik](https://www.flaticon.com/authors/monkik "monkik") from [www.flaticon.com](https://www.flaticon.com/ "Flaticon")

Icons made by [DinosoftLabs](https://www.flaticon.com/authors/dinosoflabs "DinosoftLabs") from [www.flaticon.com](https://www.flaticon.com/ "Flaticon")

Icons made by [Vitaly Gorbachev](https://www.flaticon.com/authors/vitaly-gorbachev "Vitaly Gorbachev") from [www.flaticon.com](https://www.flaticon.com/ "Flaticon")

Icons made by [Flat Icons](https://www.flaticon.com/authors/flat-icons "Flat Icons") from [www.flaticon.com](https://www.flaticon.com/ "Flaticon")

Icons made by [photo3idea_studio](https://www.flaticon.com/authors/photo3idea-studio "photo3idea_studio") from [www.flaticon.com](https://www.flaticon.com/ "Flaticon")

Icons made by [srip](https://www.flaticon.com/authors/srip "srip") from [www.flaticon.com](https://www.flaticon.com/ "Flaticon")

Icons made by [Eucalyp](https://www.flaticon.com/authors/eucalyp "Eucalyp") from [www.flaticon.com](https://www.flaticon.com/ "Flaticon")

Icons made by [ultimatearm](https://www.flaticon.com/authors/ultimatearm "ultimatearm") from [www.flaticon.com](https://www.flaticon.com/ "Flaticon")



Questions?