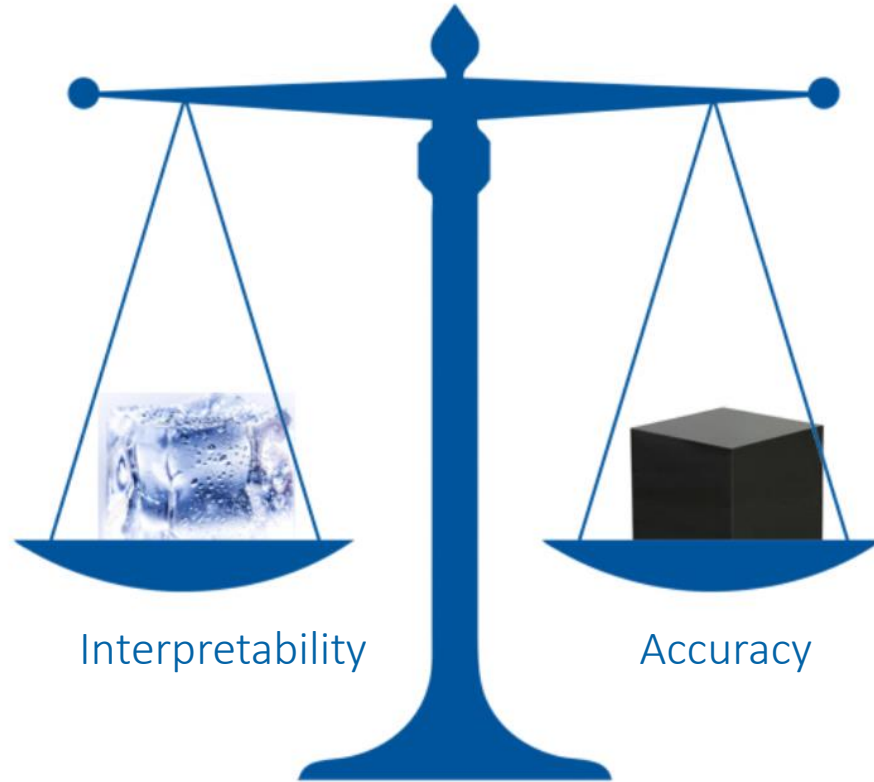# Model Interpretability

FANS Network Meeting | Data Science |June 2, 2021
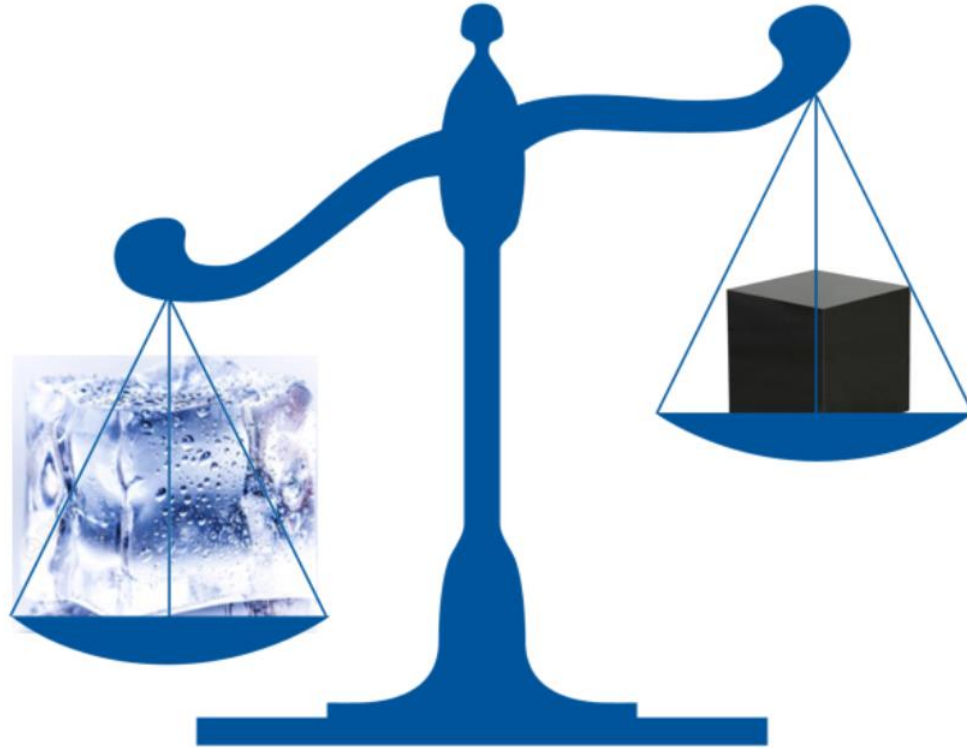
Mathias Lanner SAS Institute
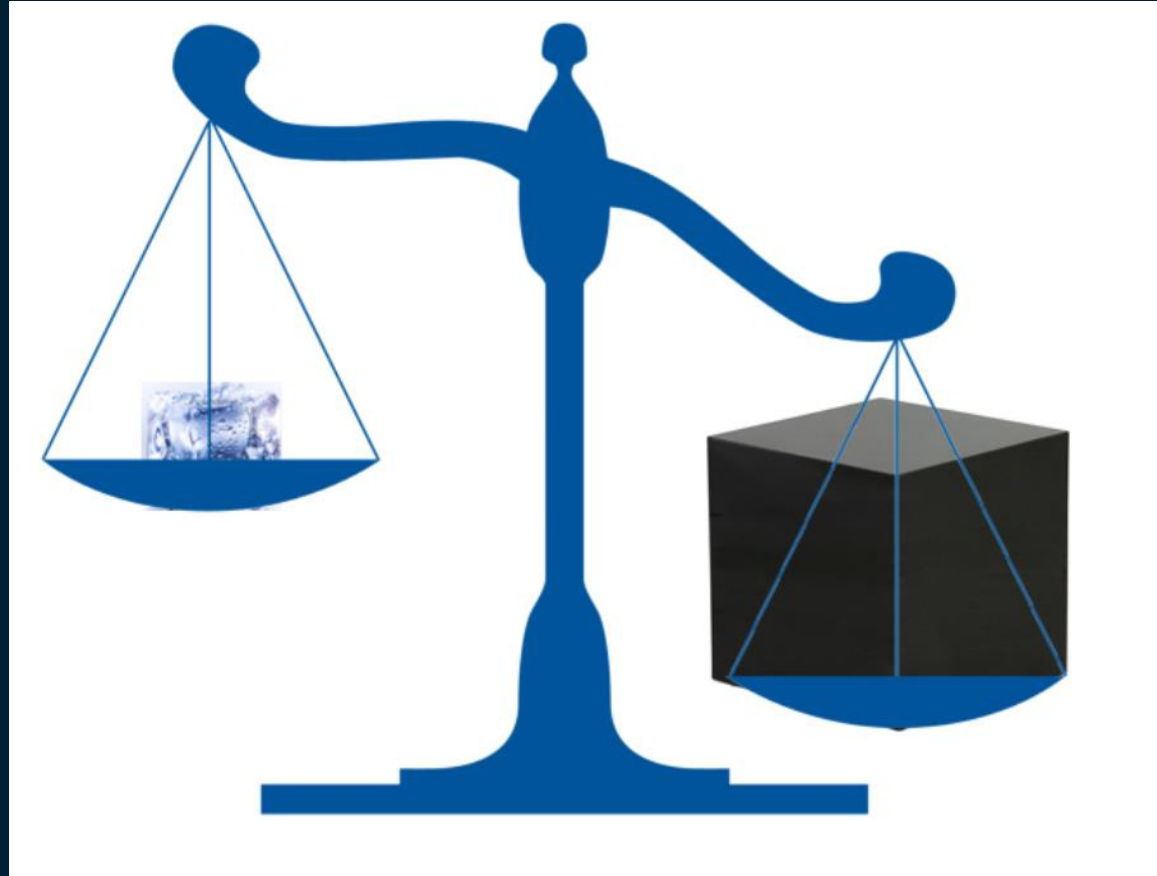
§sas

# White Box- vs Black Box-models



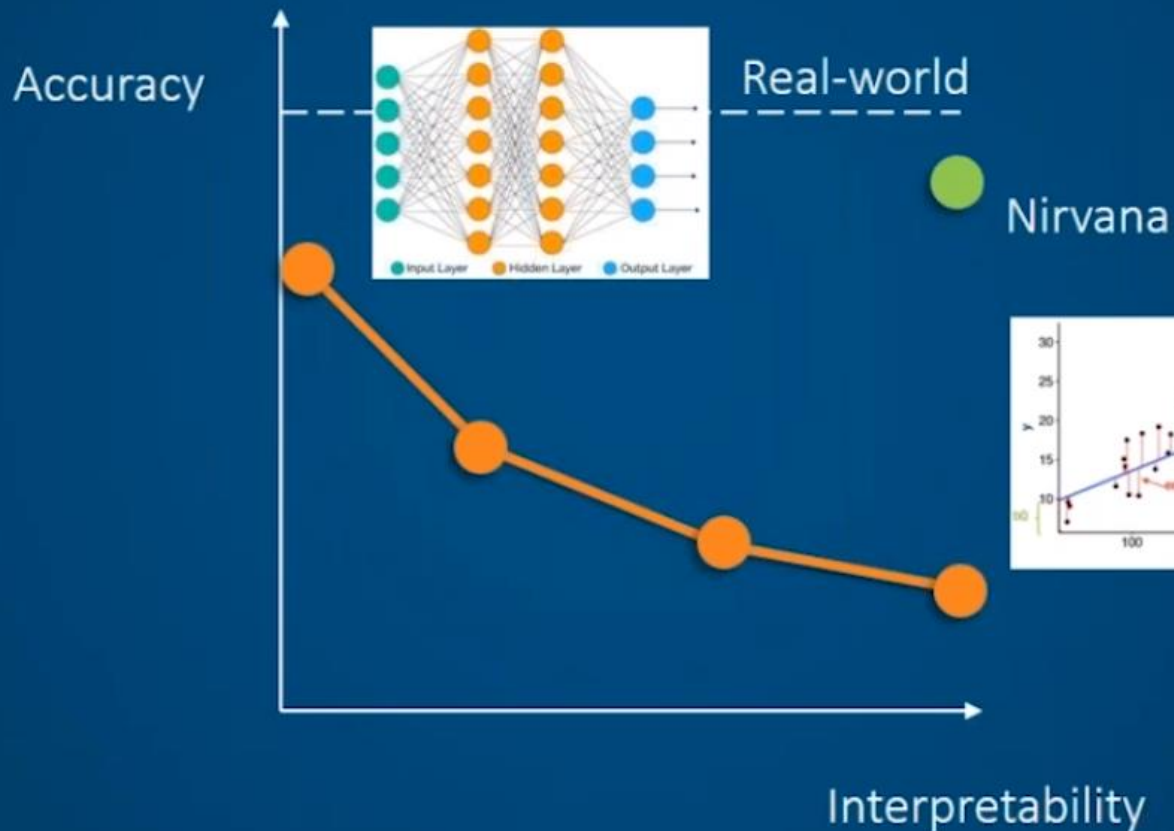Interpretability

Accuracy

§sas

# What do you need?

# What do you need?

# Accuracy vs Interpretability

# Can you taste the difference?

How does the ingrediencies impact the taste experience?



- Gin & Tonic

- Boeuf Bourguignon

§sas

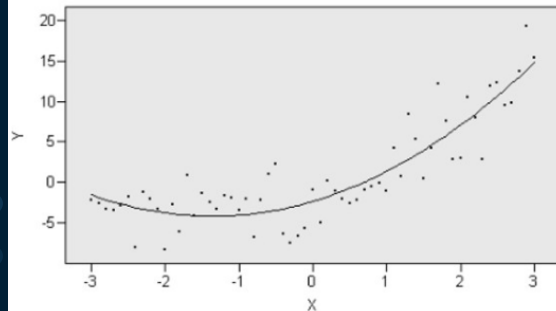| Transparent "White Box" | Opaque "Black Box" |
|---|---|
| Regressions | Neural Networks |
| Decision Trees | Random Forests |
| Rules-Based | Gradient Boosting |



$$y = \beta_0 + \beta_1 x + \beta_2 x^2$$
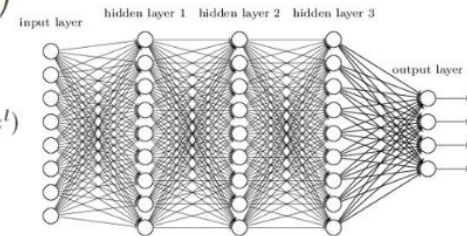


$$a_j^l = \sigma\left(\sum_k w_{jk}^l a_k^{l-1} + b_j^l\right)$$

$$\delta^L = \nabla_a C \odot \sigma'(z^L)$$

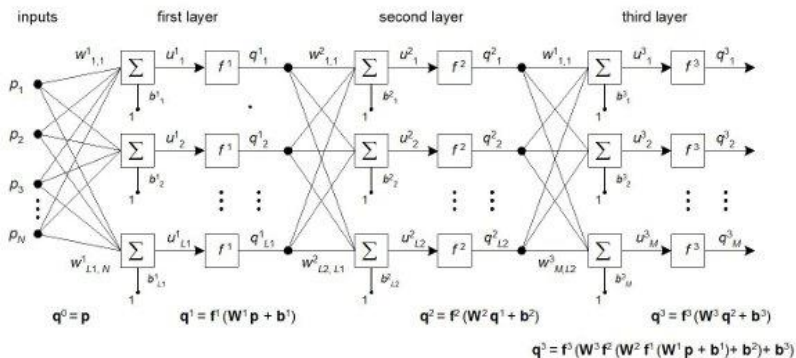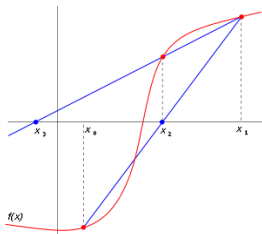$$\delta^l = ((w^{l+1})^T \delta^{l+1}) \odot \sigma'(z^l)$$

$$\frac{\partial C}{\partial b_j^l} = \delta_j^l$$

$$\frac{\partial C}{\partial w_{jk}^l} = a_k^{l-1} \delta_j^l$$

input layer   hidden layer 1   hidden layer 2   hidden layer 3   output layer



§.sas

# The black box



$$s = \sqrt{\frac{\sum(x - \bar{x})^2}{n - 1}}$$



inputs | first layer | second layer | third layer

$$q^0 = p \qquad q^1 = f^1(W^1 p + b^1) \qquad q^2 = f^2(W^2 q^1 + b^2) \qquad q^3 = f^3(W^3 q^2 + b^3)$$
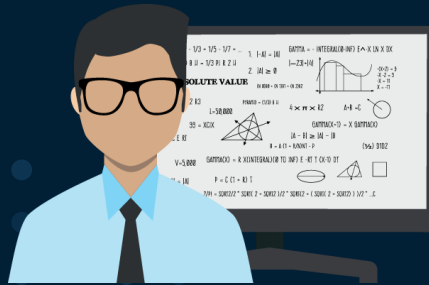
$$q^3 = f^3(W^3 f^2(W^2 f^1(W^1 p + b^1) + b^2) + b^3)$$

How can we generate models which are not only accurate, but
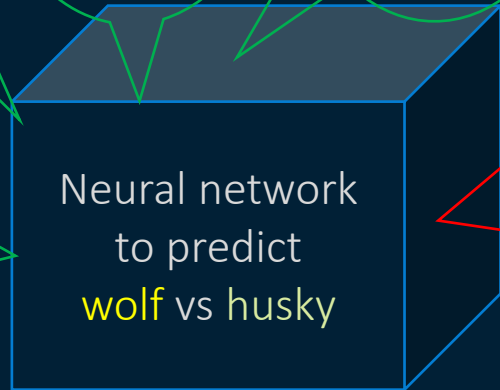
Fair
Accountable
Transparent
Trustworthy
Explainable
?

Ssas

# Figure out when NOT to trust a model



Prediction accuracy is very high. It is time to put this system online

You are detecting snow, not wolves! I can't trust you

Data scientist

Neural network to predict wolf vs husky

Predicted: wolf
True: wolf

Predicted: wolf
True: wolf

Predicted: husky
True: husky

Predicted: husky
True: husky

Predicted: wolf
True: husky

§.sas

# Shine some light on the black box

## Road to Trustworthy AI



Being able to interpret and explain machine learning models is key to trustworthy AI

§sas

# Proxy Methods and Diagnostics

## Proxy Methods:

- **Surrogate model approach**
  Fit a black box machine learning model (deep learning, gradient boosting, random forest, etc.) to your training data. Then use those predicted outcomes as the targets for a more interpretable model (decision tree, regression)

- **Machine learning as benchmark**
  Use a complex model to set the goal for potential accuracy metrics that could be achieved, then use that as the standard against which you compare the outputs of more interpretable model types

- **Machine learning for feature creation**
  Use ML/Deep Learning to extract the features, then use those features as inputs to a more explainable model type

## Post-Modeling diagnostics:

- Variable importance(VI)
- Partial Dependence (PD)
- Individual Conditional Expectation (ICE)
- Local Interpretable Model-agnostic Explanations (LIME)
- **SHapley Additive exPlanations (SHAP)**

§.sas

# Post-Modeling Diagnostics

## Input-Output Relationship

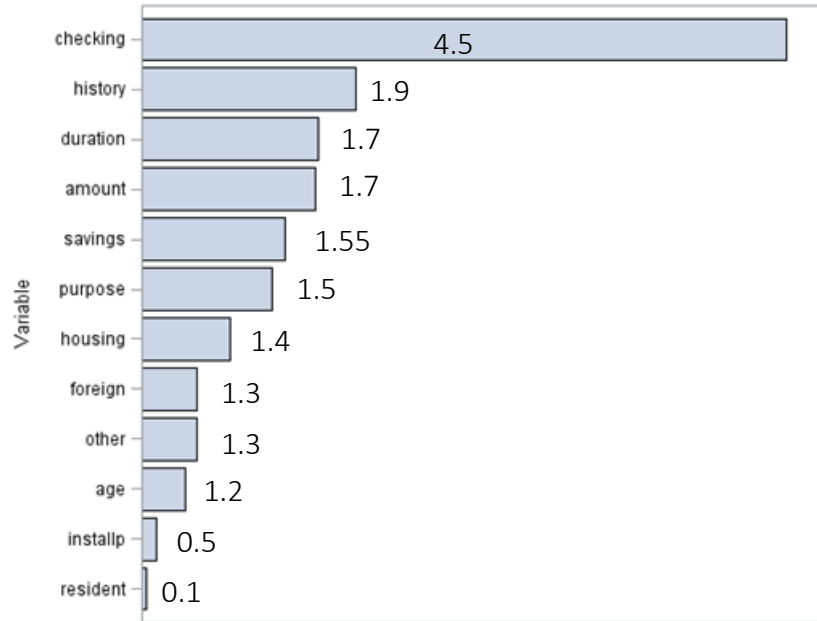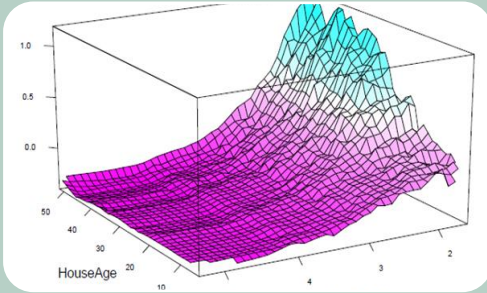| Question | Technique |
|---|---|
| What are the top inputs? | **Variable Importance(VI)/Relative VI** |
| How do the drivers work? | **Partial Dependence (PD)**<br>**Individual Conditional Expectation (ICE)** |
| What is the explanation for a particular prediction? | **Local Interpretable Model-agnostic Explanations (LIME)**<br>**SHapley Additive exPlanations (SHAP)** |

Ssas

# Variable Importance

## Variable Importance

| Variable Name | Train Importance |
|---|---|
| curr_days_susp | 81.3108 |
| handset_age_grp | 46.5221 |
| ever_days_over_plan | 30.2300 |
| pymts_late_ltd | 24.8962 |
| billing_cycle | 16.2053 |
| avg_days_susp | 15.8511 |
| calls_care_ltd | 13.3913 |
| call_category_1 | 13.1900 |

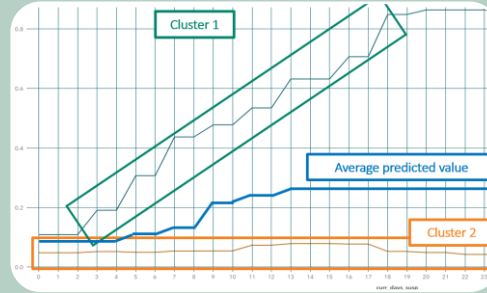Train Importance is calculated as sum of the decrease in error when split by a variable.

### Selected Variable Importance

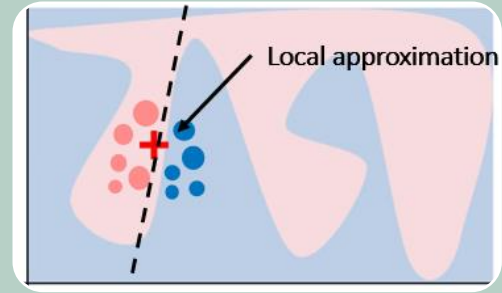| Variable | Importance |
|---|---|
| checking | 4.5 |
| history | 1.9 |
| duration | 1.7 |
| amount | 1.7 |
| savings | 1.55 |
| purpose | 1.5 |
| housing | 1.4 |
| foreign | 1.3 |
| other | 1.3 |
| age | 1.2 |
| installp | 0.5 |
| resident | 0.1 |

§.sas

# Partial Dependence Plots

depicts relationship between the value of an input variable and the value of the model predictions after the influence of all other variables has been averaged out

# Individual Conditional Expectation (ICE)
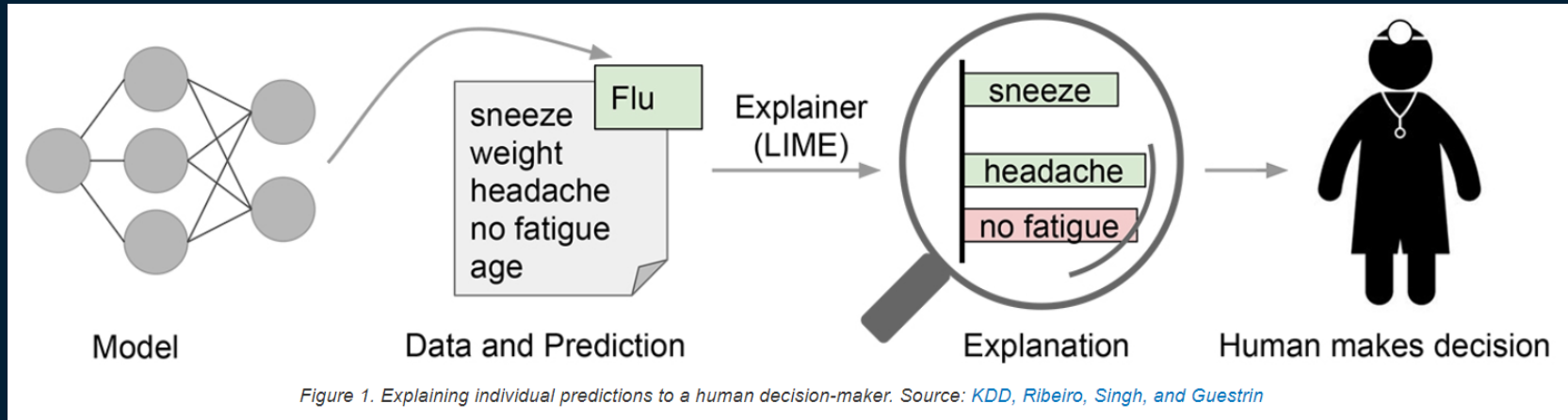
helps identify subgroups and interactions

# Local Interpretable Model-agnostic Explanations (LIME)

builds an interpretable model of explanatory data samples at local areas in the analyzed data

§sas

# LIME (Local Interpretable Model-agnostic Explanations)
## Flu Prediction (Neural Network)



Figure 1. Explaining individual predictions to a human decision-maker. Source: KDD, Ribeiro, Singh, and Guestrin

# Some examples…….

sas.com