

Assessing Global Sensitivity Analysis for Credibility in Machine Learning Explainability



SAMSI Presentation
Numerical Analysis for Data Sciences Program
Global Sensitivity Analysis Working Group
2020-09-17



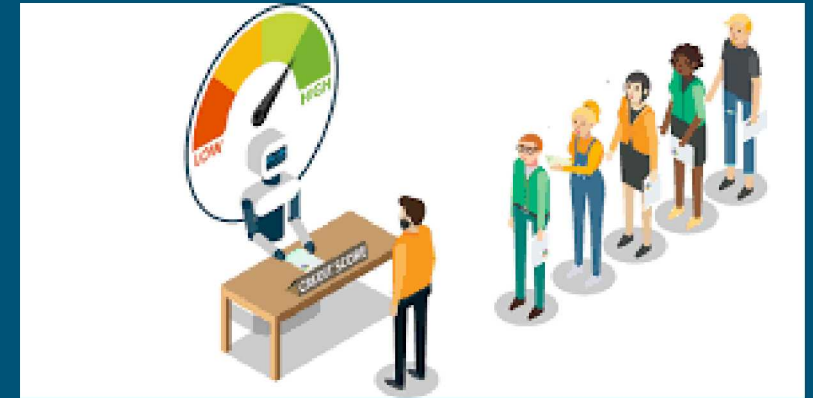
PRESENTED BY

Erin Acquesta and Mike Smith

Rich Field, Trevor Maxfield, Ahmad Rushdi,
and many others...

The Need for Credible Explainability

ML is being used in an increasingly number of high-consequence applications.
ML explainability has emerged as field that seeks to build trust.



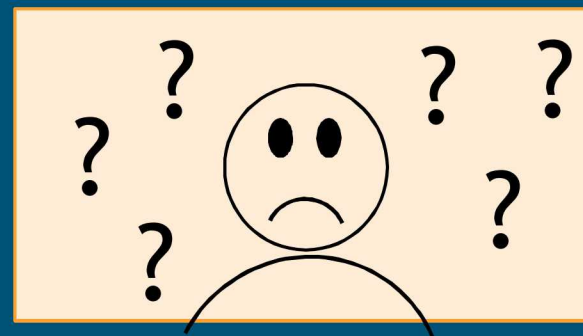
The Need for Credible Explainability

ML is being used in an increasingly number of high-consequence applications.
ML explainability has emerged as field that seeks to build trust.



Can we trust the explanation?

- Computational shortcuts
- Assumes some understanding of machine learning
- Lack verifiable foundations



An aerial photograph of a university campus, likely the University of Washington, showing various academic buildings, green spaces, and a surrounding landscape of rolling green hills under a clear sky. The image is used as a background for the slide, with a semi-transparent blue overlay on the right side.

Current ML Explainability Methods

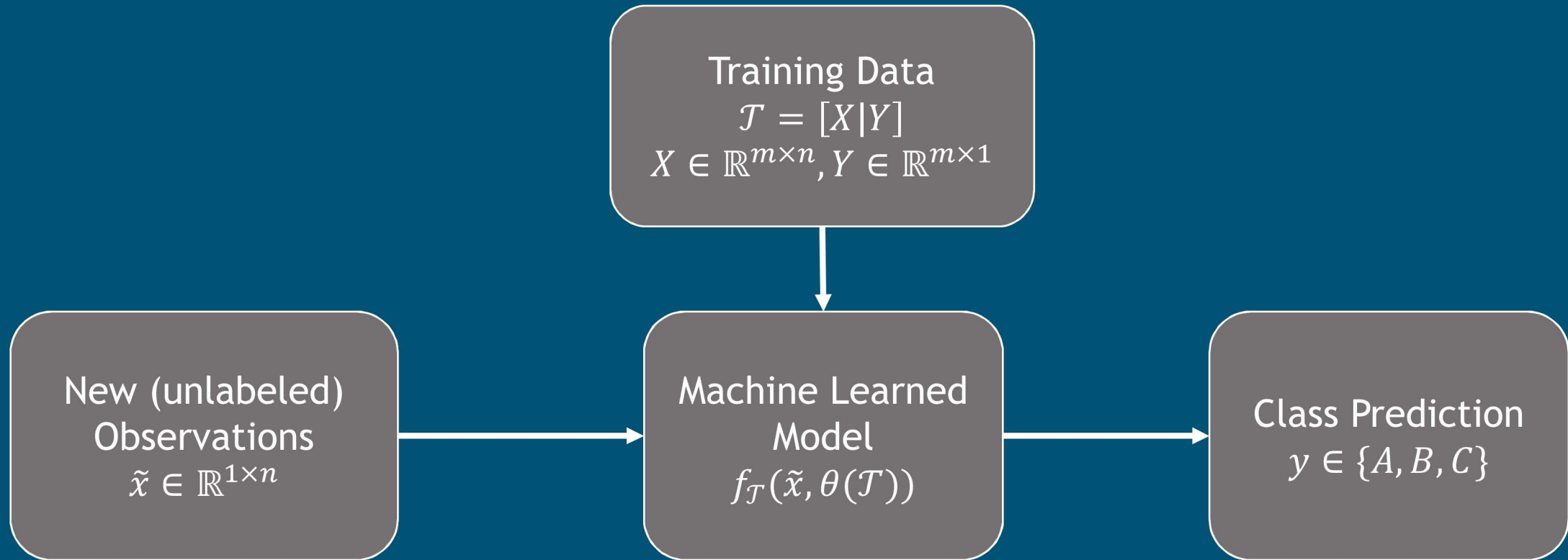
Sensitivity Analysis Guided Explainability

Correlated Feature

Training Data Statistics Preserving Sampling

Gaps and Limitations

Machine Learning (Supervised) Classification Model Diagram



ML Explainability

Attempt to describe the decision process that a machine learned model uses to make a prediction

Interpretable models

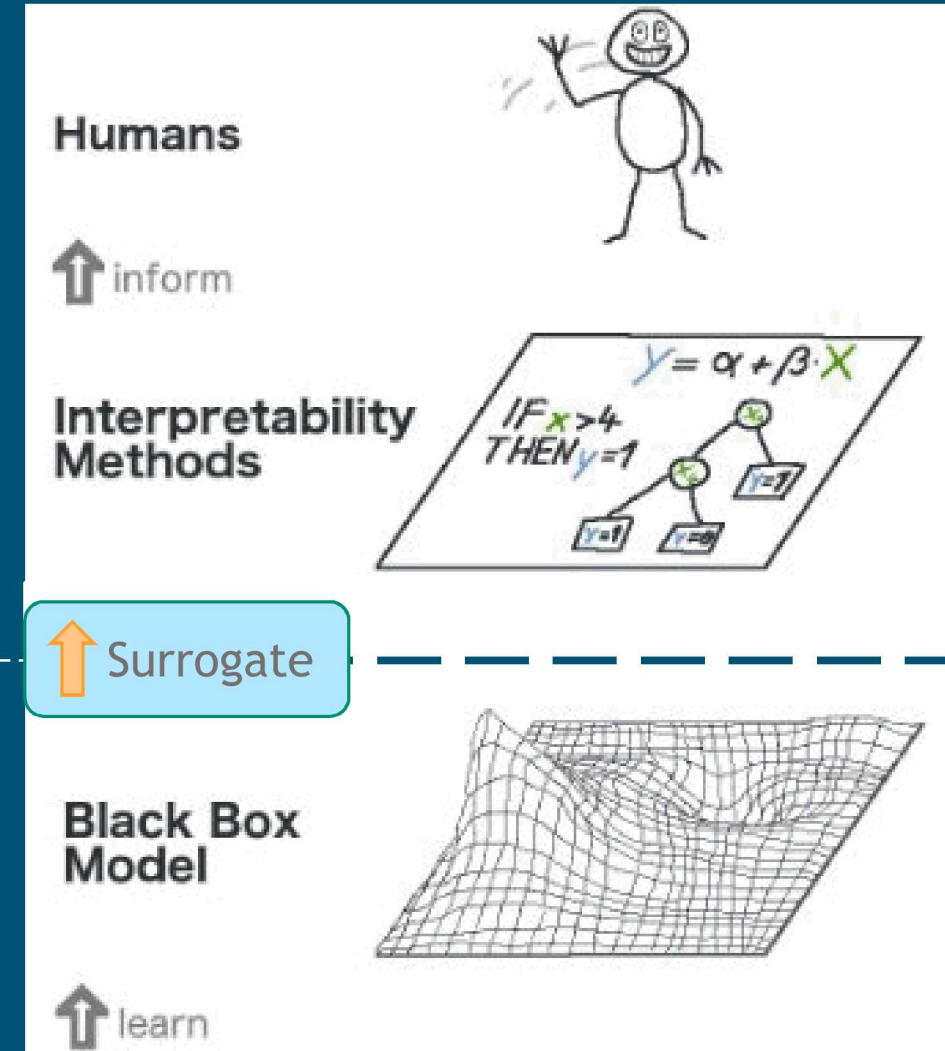
- Can inspect the model
- Models that are relatively easy to interpret (linear regression models, shallow decision trees)

White-box/Integrated

- Can inspect the model, but the model is sufficiently complex
- Gini-importance for decision trees
- Gradient-based methods for deep learning models (Saliency maps)

Black-box/Post-hoc

- Do not inspect the model, instead evaluates feature importance
- Create a surrogate model that is interpretable
- Often perturb the data and observe how the output changes



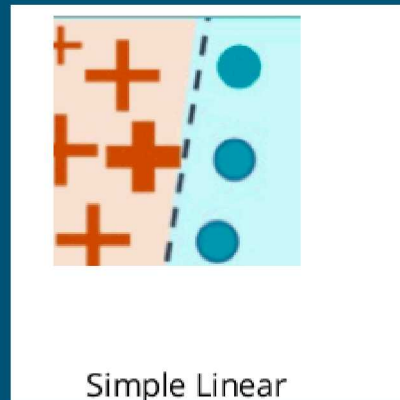
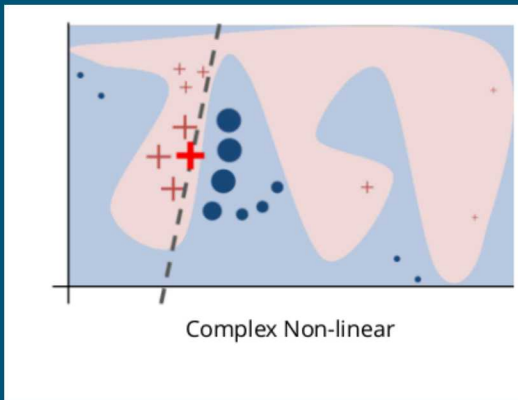
LIME: Local Interpretable Model-agnostic Explanation [1]

Perturbation Based Method

For $\tilde{x} \in \mathbb{R}^{1 \times n}$ (a new realization)

1. Sample \tilde{m} observation in a neighborhood of \tilde{x} , assuming that each feature of \tilde{x} is independent and normally distributed. Resulting in: $\{\tilde{z}^i\}$ for $i = 1 \dots \tilde{m}$
2. For each of the samples, get the original ML model prediction $f_{\mathcal{T}}(\tilde{z}^i) = \tilde{y}^i$
3. With the derived sample dataset $\tilde{\mathcal{T}} = \{\tilde{z}^i | \tilde{y}^i\}$, learn a linear regression model,

$$g_{\tilde{\mathcal{T}}}(\tilde{x}, \theta(\tilde{\mathcal{T}})) = \sum_{j \in \{1, \dots, \tilde{m}\}} \alpha_j \tilde{x}_j$$



LIME feature importance index =

$$\phi_j = \alpha_j * \tilde{x}_j$$

SHAP: SHapley Additive exPlanations [2]

Based on Cooperative Game Theory

N : the set of features

$j \in N$

$S \subseteq N$

$v(S)$: the total value of the S features

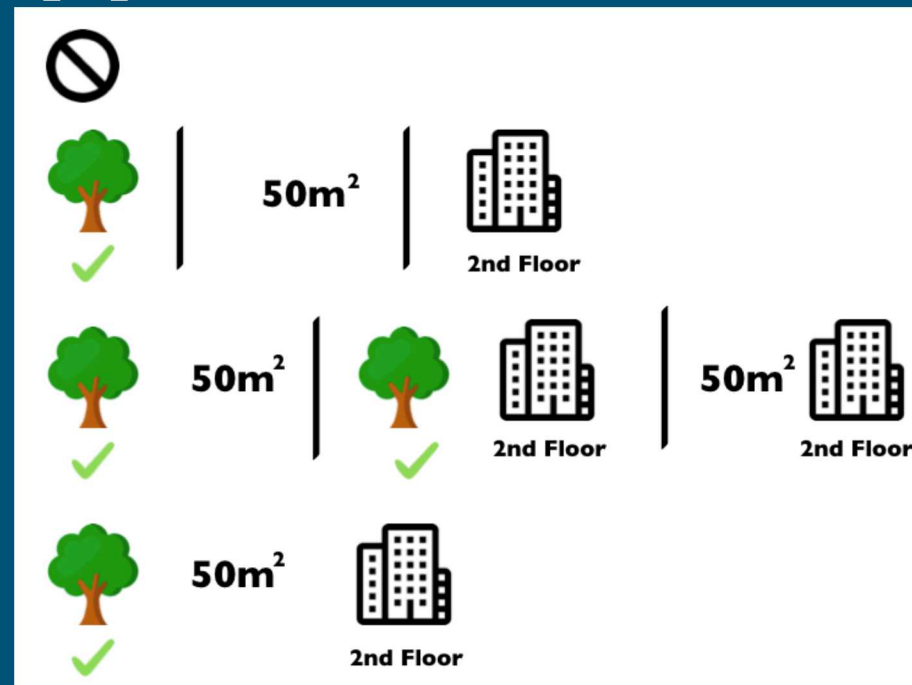
$(v(S \cup \{j\}) - v(S))$ for $S \subseteq N \setminus \{j\}$:
marginal contribution of feature j to set S

SHAP feature importance index:

$$\phi_j = \sum_{S \subseteq N \setminus \{j\}} \frac{|S|! (N - |S| - 1)!}{N!} (v(S \cup \{j\}) - v(S))$$

One Computational Approximation:

- To approximate $(v(S \cup \{j\}) - v(S))$ without retraining the model, “noise” is used as a surrogate for feature removal.
- “noise” is replicated as independent bootstrap samples from the training data



Summary of Black-Box Explanation Methods

Dependencies and Assumptions

- Dependent on a process for sampling the data
- Require distance on output—how much it changes
- Assumes independence and linearity

Observed Deficiencies [3]

- Descriptive Accuracy: Match when features are removed
- Instability: Produces different explanations on the same input
- Completeness: Generate explanations for all possible input vectors
- Efficiency: Can be slow to calculate especially as the dimensionality increases

No agreed upon definition of explainability or what constitutes an explanation

Can We Trust the ML Explanation??

Goal of explainability: provide credibility evidence by describing the decision process that a model considers for a prediction.

BUT....

LIME

- Introduces a surrogate model without assessing the fidelity of that surrogate in relation to the black-box model.
- Implements perturbation of the input uncertainty assuming independence and normally distributed random variates...which is uncharacteristic of the true statistics from the training dataset.

SHAP

- Has the most theoretical guarantees available to us today, but computational heuristics for the approximations break those theoretical guarantees.
- Bootstrapping from the training dataset to replicate "noise"...but is "noise" an appropriate replicate for removing a feature??

Can we use Sensitivity Analysis Guided Explainability to mitigate these challenges??

An aerial photograph of a university campus, showing various buildings, green spaces, and surrounding hills. The image is overlaid with a semi-transparent blue layer. On the left side, there is a vertical stack of five teal-colored rectangular boxes containing text. The second box from the top is highlighted with a white border.

Current ML Explainability Methods

Sensitivity Analysis Guided Explainability

Correlated Features

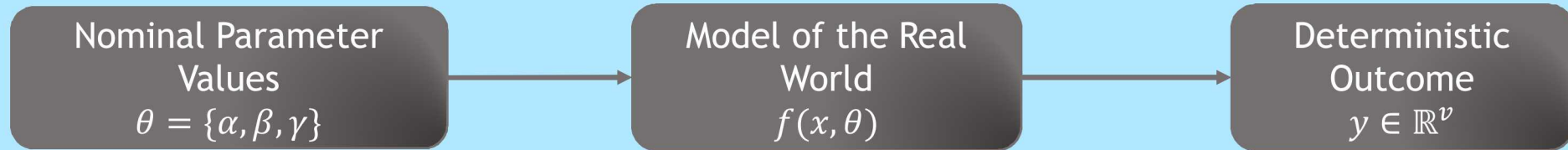
Training Data Statistics Preserving Sampling

Gaps and Limitations

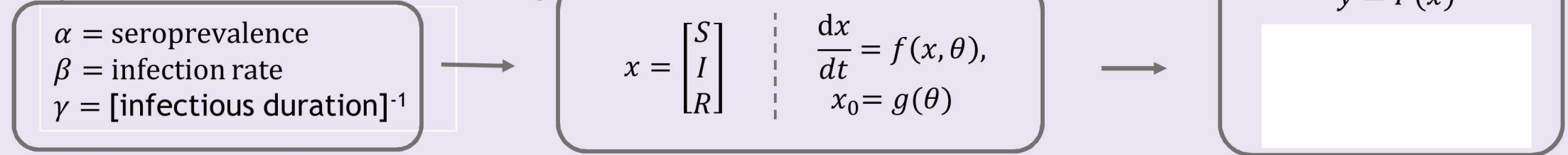
Global Sensitivity Analysis: Notational Example

The apportionment for the contributions of input uncertainties on output uncertainty. [4]

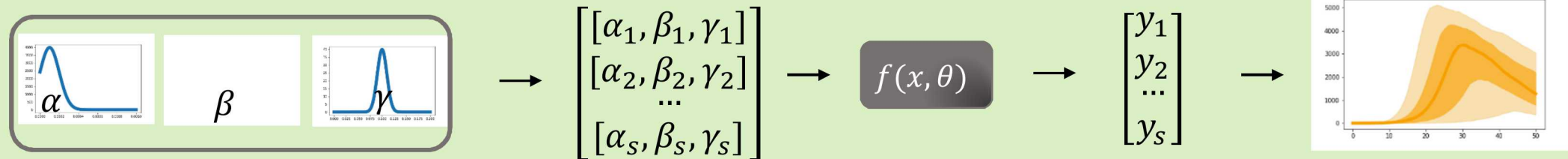
Modeling Flow Diagram



Example: Infectious Disease Modeling



Uncertainty Quantification of Model Forecast for the Infection Rate Curve



Sensitivity of the Infection Rate Curve with respect to the Parameters of the Model



Experimental Design for Sensitivity Analysis

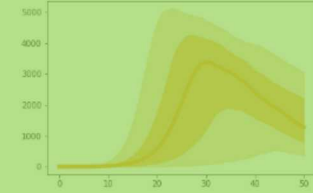
Experimental Design is a scientific approach for identifying the inputs to a *process* that are most influential to the outcome of that process; following particular design decisions.



$$\begin{bmatrix} [\alpha_1, \beta_1, \gamma_1] \\ [\alpha_2, \beta_2, \gamma_2] \\ \dots \\ [\alpha_s, \beta_s, \gamma_s] \end{bmatrix}$$

$$f(x, \theta)$$

$$\begin{bmatrix} y_1 \\ y_2 \\ \dots \\ y_s \end{bmatrix}$$



Inputs:

Uncertainty in Parameters

Design Decision I

Sampling

Sampling sufficient discrete realizations that preserve the statistics and introduces only marginal standard error.

Process:

Mathematical Model

Design Decision II

Controlled/Uncontrolled Random Behavior

Controlled: Sources of Variance
Uncontrolled: Random behaviors inherent to the model

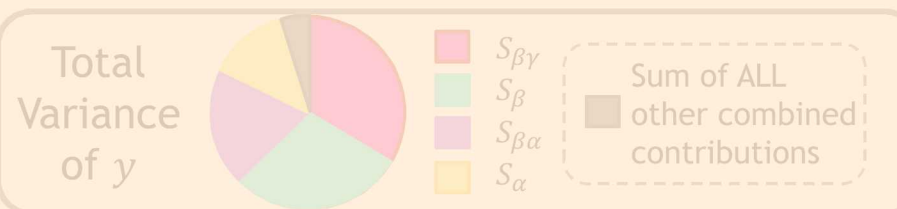
Outcome:

Uncertainty in Model Output

Design Decision III

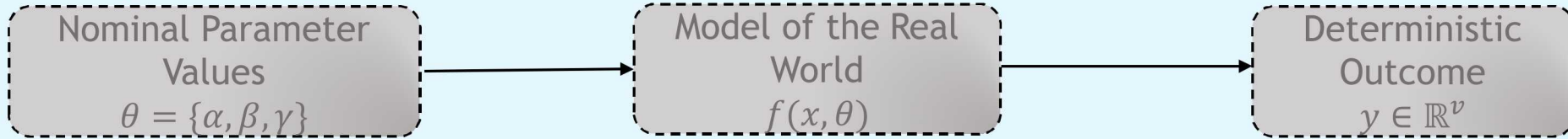
Quantity of Interest (QoI)

For the intended use case, what output from the model maps to quantitative metric for that intended purpose.

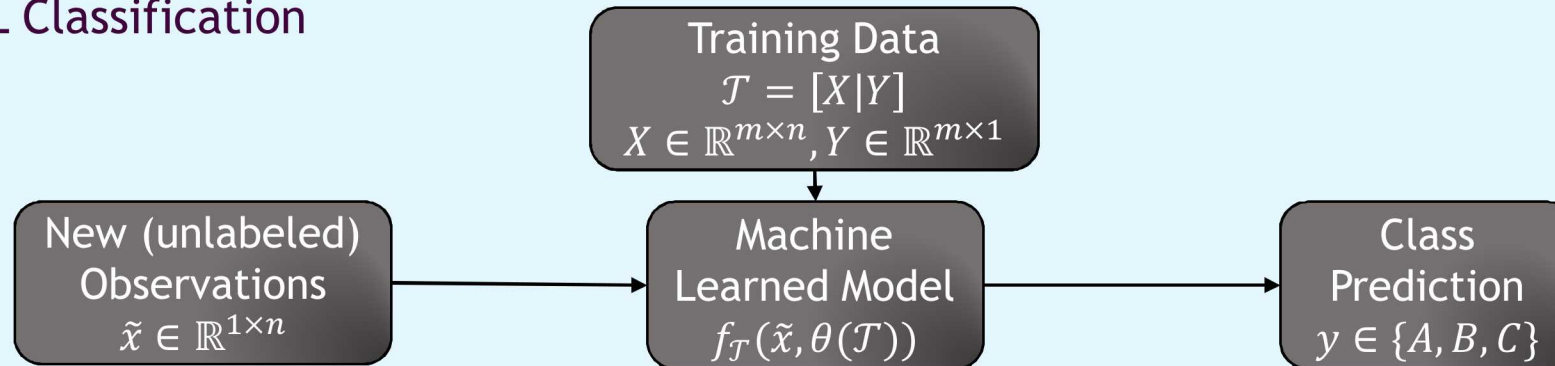


Experimental Design for ML Explainability

Original Modeling Flow Diagram



Translation to ML Classification



Inputs:

Uncertainty in Features

Sampling

Preserving the statistical properties of the training data: **non-Gaussian**, **discrete**, **correlated**, and **sparse**

Process:

Machine Learned Model

Controlled/Uncontrolled Random Behavior

Running sufficient replicates for the random behavior of stochastic machine learned models.

Outcome:

Uncertain Model Predictions

Quantity of Interest (QoI)

What is the appropriate QoI for which a sensitivity analysis will provide insight for ML explainability?

Methods to apportion the influence of sources of input uncertainty across output uncertainty, accounting for **higher-order interactions** in a model and **input correlations**.

An aerial photograph of a city, likely a university campus, with various buildings and green spaces. The image is overlaid with a semi-transparent blue layer. On the left side, there is a vertical stack of five teal-colored rectangular boxes containing text. The third box from the top is highlighted with a white border.

Current ML Explainability Methods

Sensitivity Analysis Guided Explainability

Correlated Features

Training Data Statistics Preserving Sampling

Gaps and Limitations

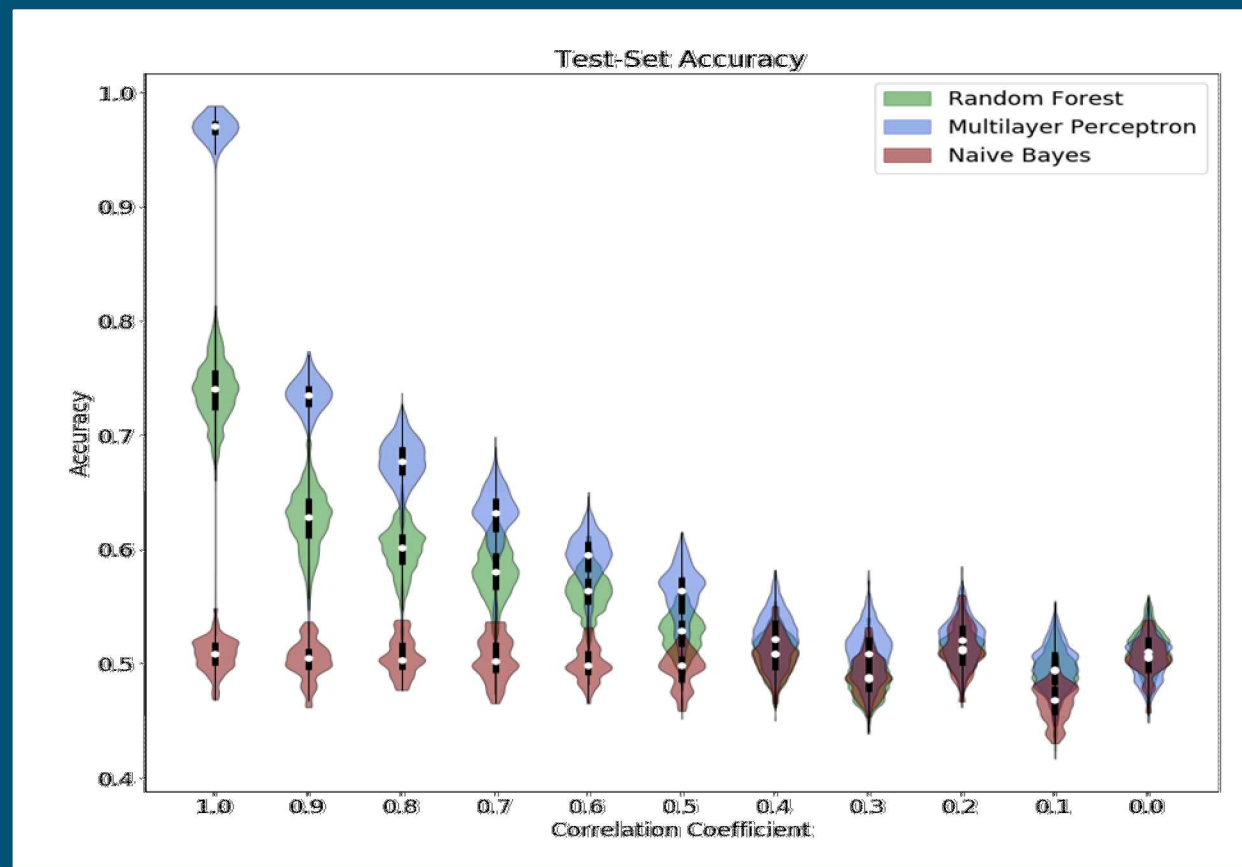
Correlated Features Provide Discriminative Power

Design of Experiments was used to debunk the myth that correlated variables only provide redundant information.

- Used synthetic data to control the amount of correlation as the distinguishing characteristic between classes
- Naïve Bayes is the baseline for linear relationships

Most explainability methods assume independence

- Incongruent explanations for the learned model
- LIME uses a linear model
- SHAP makes independence and linear assumptions
- Tested with quadratic regression



An aerial photograph of a city, likely Seattle, with a prominent blue overlay. The city features a mix of urban buildings, green spaces, and a waterfront area. The blue overlay is semi-transparent, allowing the city details to be visible while providing a cohesive background for the text.

Current ML Explainability Methods

Sensitivity Analysis Guided Explainability

Correlated Features

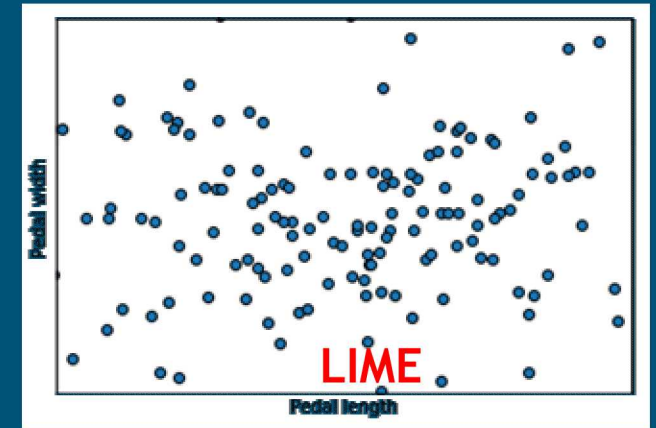
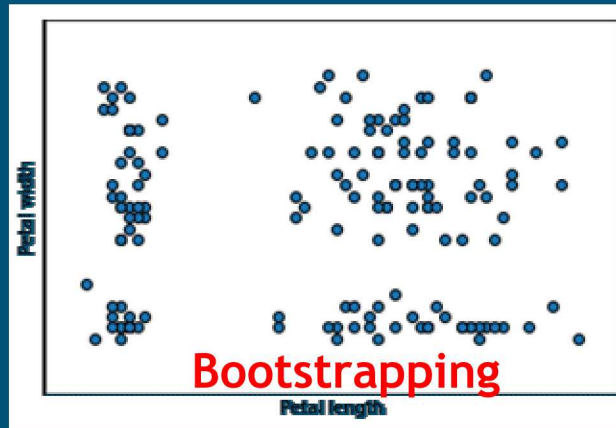
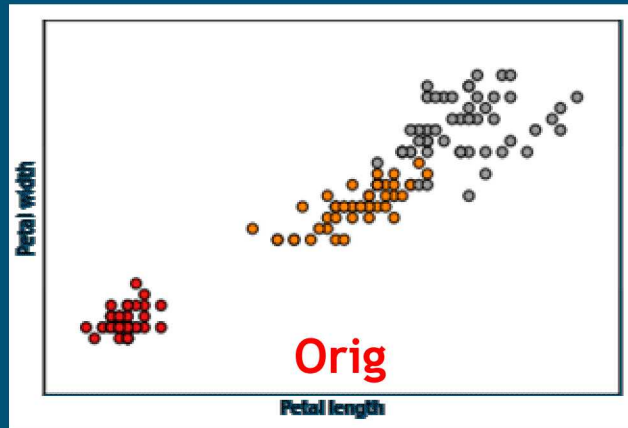
**Training Data Statistics Preserving
Sampling**

Gaps and Limitations

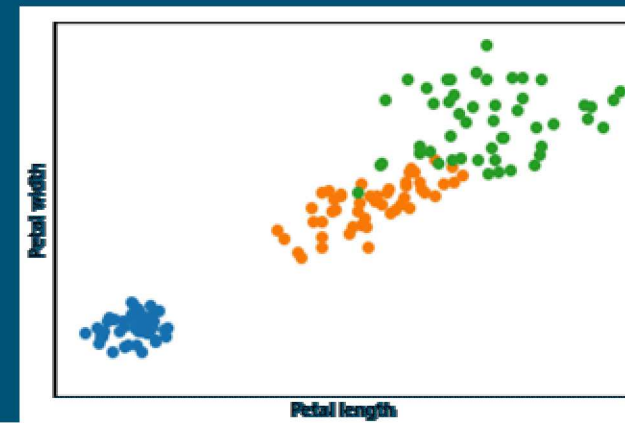
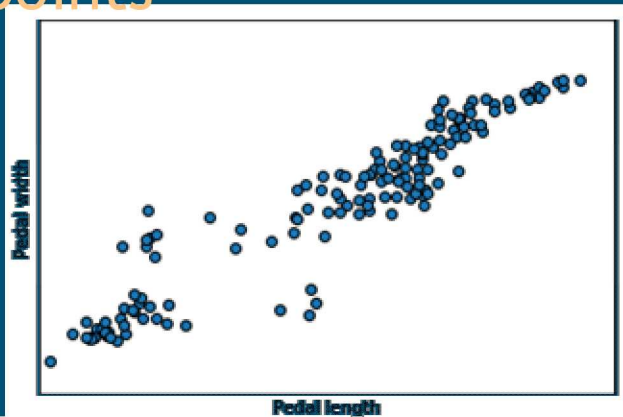
Training Data Statistics Preserving Sampling

Current sampling approaches can create unrealistic data points
(Do not preserve correlations or distributions)

Features 3 and 4 from the iris data set (correlated features)



Developed Sampling methods that preserve correlations and generate realistic data points



An aerial photograph of a city, likely a university campus, with various buildings and green spaces. The image is overlaid with a semi-transparent blue layer. On the left side, there is a vertical stack of five teal-colored rectangular boxes containing text. The bottom-most box is a darker green and has a white border.

Current ML Explainability Methods

Sensitivity Analysis Guided Explainability

Correlated Feature Influence

Training Data Statistics Preserving Sampling

Gaps and Limitations

Blockers to using GSA on ML Models for Explainability

- Dependence in our feature space, which defines our sources of uncertainty
- Appropriate Quantity of Interest that will map to an explanation.
- Expedient Computational Methods for our correlation and distribution preserving sampling technique.

Questions to the Working Group:

- Are there scalable methods for approximating the Shapley indices?
- How are dependent sources of uncertainty currently handled by the research community?

Continued Research Directions:

- Extensions of Sobol indices for dependent features. [5,6]
- Verification of computational estimators for SHAP feature importance indexes.

1. Ribeiro, M.T., Singh, S. and Guestrin, C., 2016, August. " Why should I trust you?" Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining* (pp. 1135-1144).
2. Lundberg, S.M. and Lee, S.I., 2017. A unified approach to interpreting model predictions. In *Advances in neural information processing systems* (pp. 4765-4774).
3. Warnecke, A., Arp, D., Wressnegger, C. and Rieck, K., 2019. Evaluating explanation methods for deep learning in security. *arXiv preprint arXiv:1906.02108*.
4. Saltelli, A., Ratto, M., Andres, T., Campolongo, F., Cariboni, J., Gatelli, D., Saisana, M. and Tarantola, S., 2008. *Global sensitivity analysis: the primer*. John Wiley & Sons.
5. Hart, J. and Gremaud, P.A., 2018. An approximation theoretic perspective of Sobol'indices with dependent variables. *International Journal for Uncertainty Quantification*, 8(6).
6. Chastaing, G., Gamboa, F. and Prieur, C., 2012. Generalized hoeffding-sobol decomposition for dependent variables-application to sensitivity analysis. *Electronic Journal of Statistics*, 6, pp.2420-2448.