



Machine Learning Model Validation

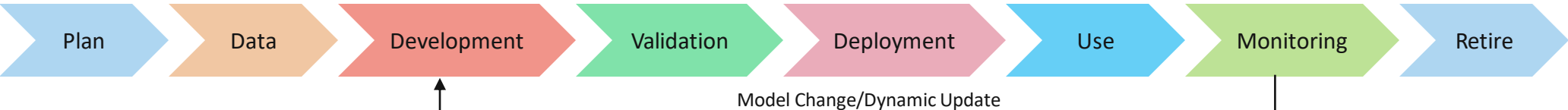
Hands-on Training with PiML Toolbox

Aijun Zhang
Head of Validation Engineering
Wells Fargo

The 13th Annual Risk Americas – AI in Banking Forum | May 20, 2024 | NYC

Disclaimer: This material represents the views of the presenter and does not necessarily reflect those of Wells Fargo.

Machine Learning Lifecycle



Data

Collection, Labelling,
Loading, Preprocessing,
Quality Check,
Data Visualization,
Feature Engineering,
Variable Selection



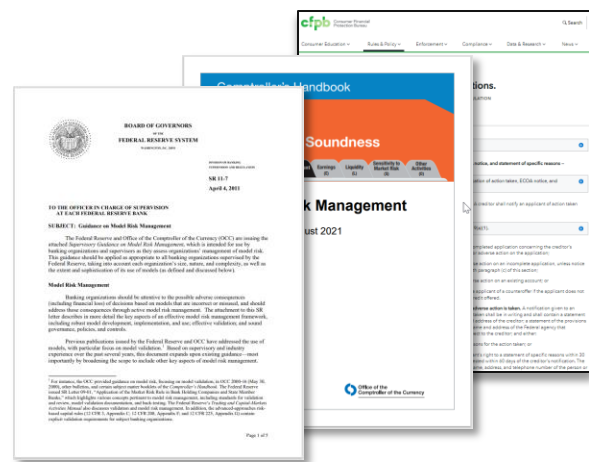
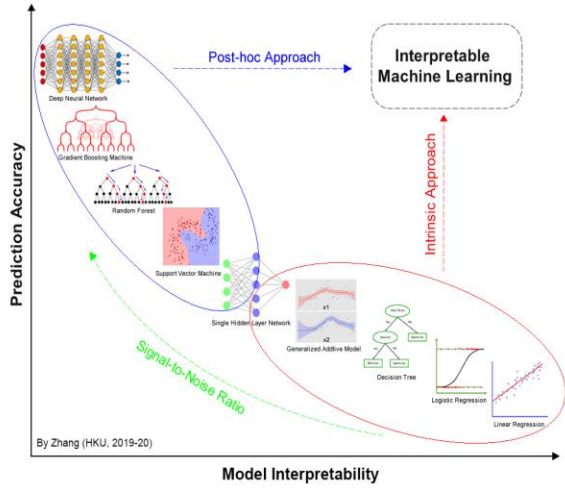
Model Development

Model Design and Assumptions,
Model Training,
Hyperparameter Tuning,
Model Calibration,
Developmental Testing,
Developmental Benchmarking

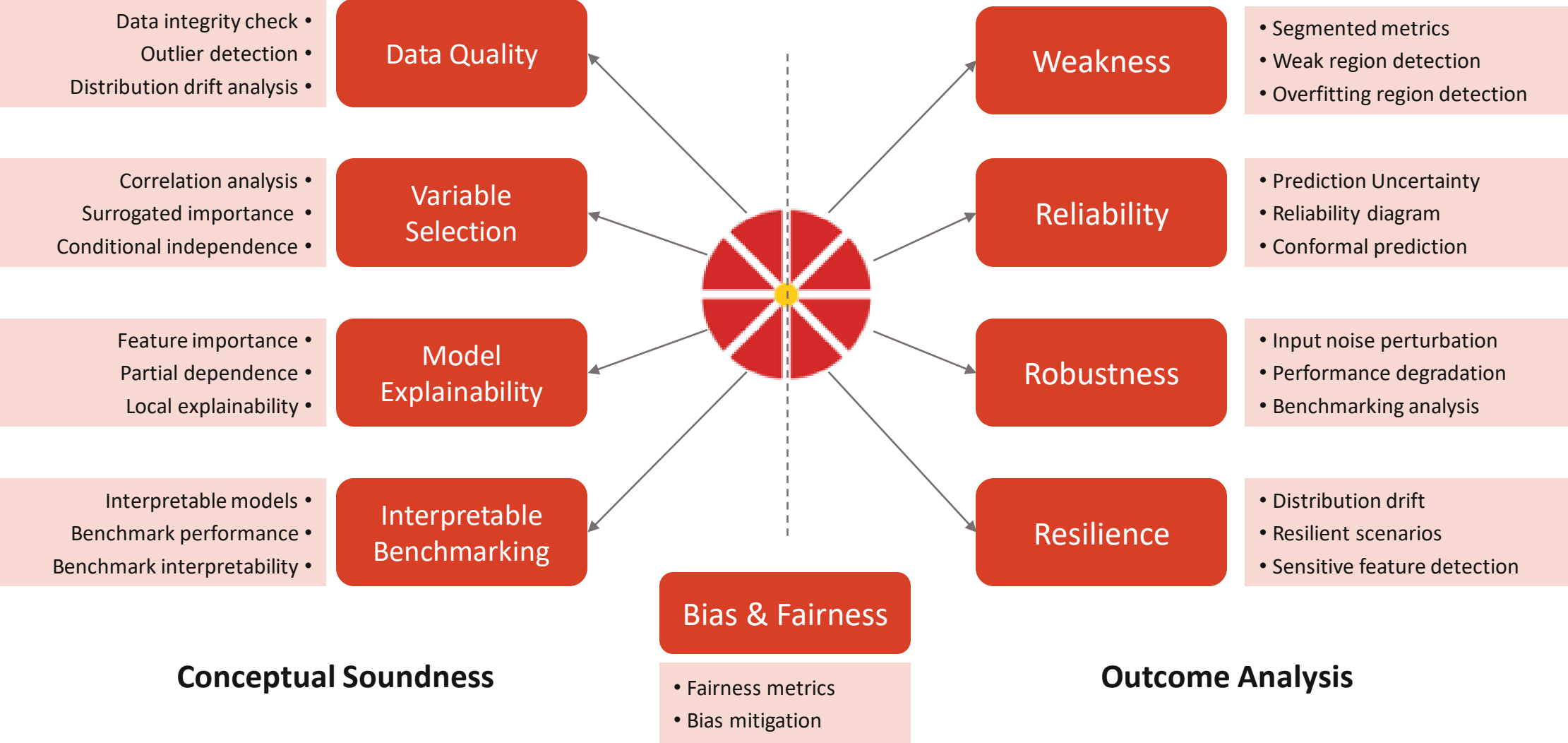


Model Validation

Independent Testing,
Independent Benchmarking,
Data Quality Check,
Conceptual Soundness,
Outcome Analysis



Machine Learning Model Validation - Key Elements



PiML Toolbox Overview



An integrated Python toolbox for interpretable machine learning

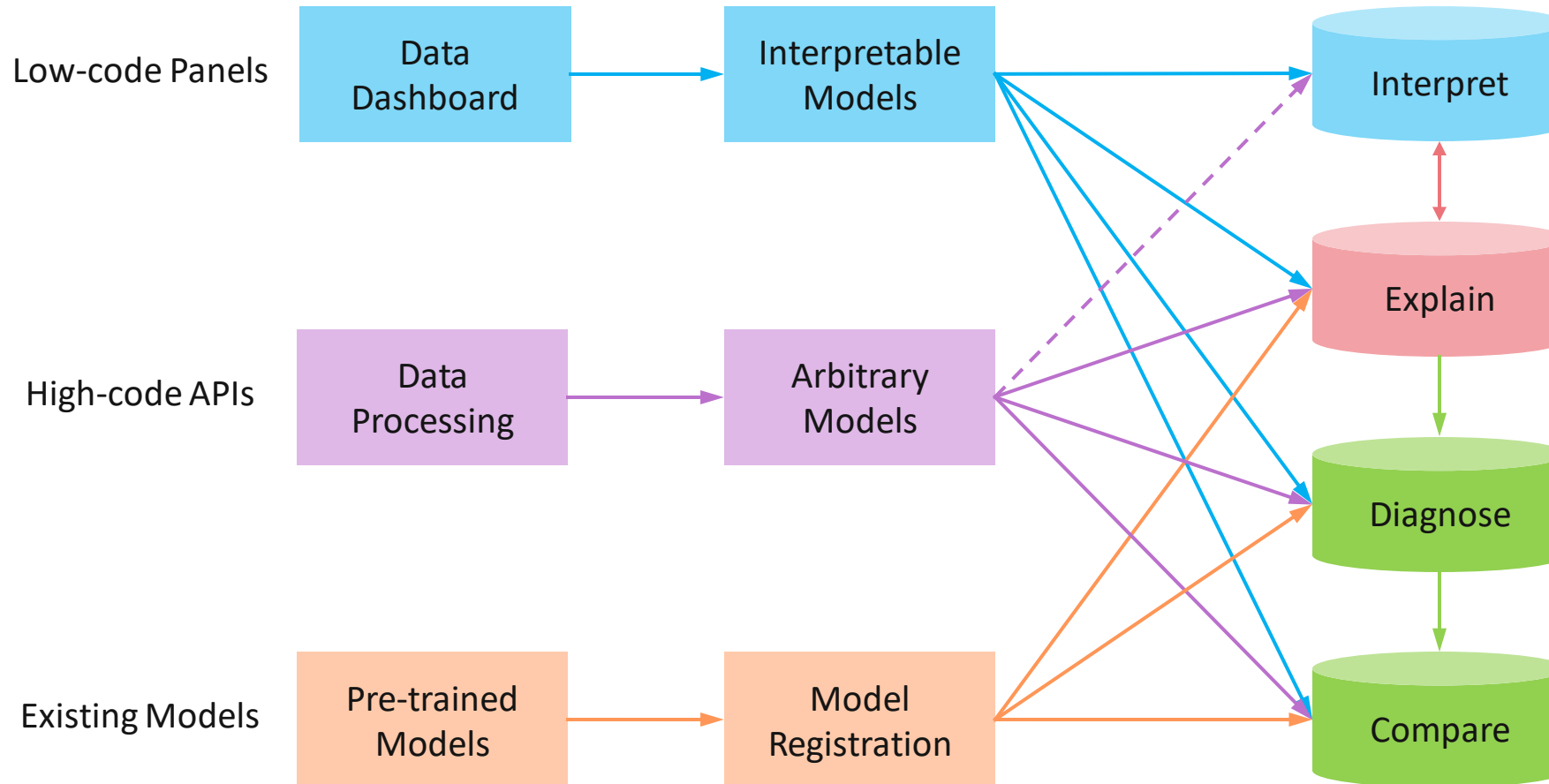
Model Development

- Data Exploration and Quality Check
- Inherently Interpretable ML Models
 - GLM, GAM, XGB1
 - XGB2, EBM, GAMI-Net, GAMI-Lin-Tree
- Locally Interpretable ML Models
 - Decision Tree
 - Sparse ReLU Neural Networks
- Model-specific Interpretability
- Model-agnostic Explainability

Model Testing

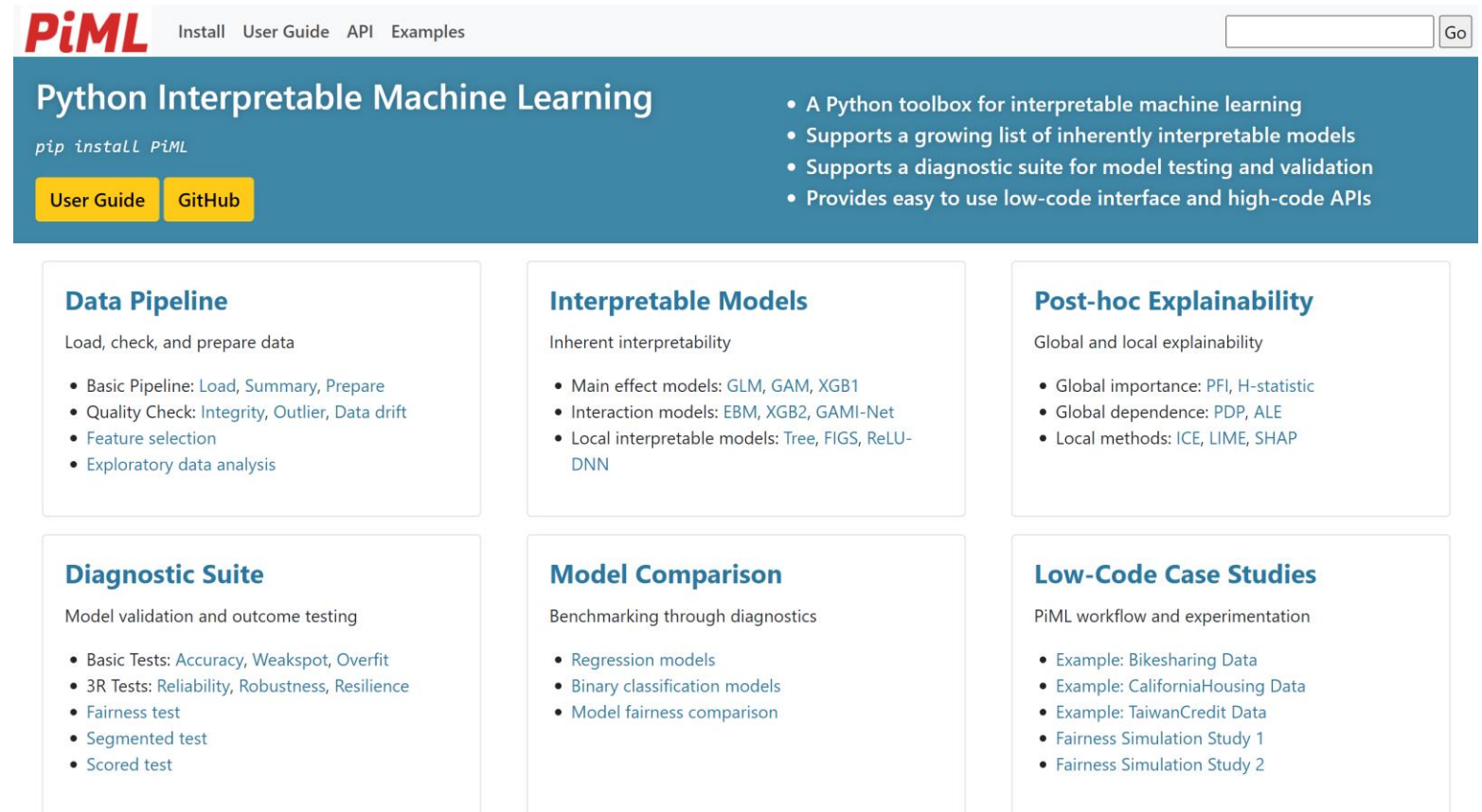
- Model Diagnostics and Outcome Analysis
 - Prediction Accuracy
 - Hyperparameter Turning
 - Weakness Detection
 - Reliability Test (Prediction Uncertainty)
 - Robustness Test (Covariate perturbation)
 - Resilience Test (Extreme scenarios)
 - Bias and Fairness
- Model Comparison and Benchmarking

PiML Pipelines



PiML User Guide

-  **May 4, 2022:** V0.1.0 is launched with low-code UI/UX.
-  **June 26, 2022:** V0.2.0 is released with high-code APIs.
-  **July 26, 2022:** V0.3.0 is released with classic statistical models.
-  **October 31, 2022:** V0.4.0 is released with enriched models and enhanced diagnostics.
-  **May 4, 2023:** V0.5.0 is released with [PiML user guide](#) with many examples.
-  **December 1, 2023:** V0.6.0 is released with enhanced data handling and model analytics.



The screenshot shows the PiML website with a navigation bar containing 'PiML', 'Install', 'User Guide', 'API', and 'Examples'. The main header area has the title 'Python Interpretable Machine Learning' and the command 'pip install PiML'. Below this are two buttons: 'User Guide' and 'GitHub'. To the right of the header is a list of features: 'A Python toolbox for interpretable machine learning', 'Supports a growing list of inherently interpretable models', 'Supports a diagnostic suite for model testing and validation', and 'Provides easy to use low-code interface and high-code APIs'. The main content area is divided into six boxes: 'Data Pipeline' (Load, check, and prepare data; Basic Pipeline: Load, Summary, Prepare; Quality Check: Integrity, Outlier, Data drift; Feature selection; Exploratory data analysis), 'Interpretable Models' (Inherent interpretability; Main effect models: GLM, GAM, XGB1; Interaction models: EBM, XGB2, GAMI-Net; Local interpretable models: Tree, FIGS, ReLU-DNN), 'Post-hoc Explainability' (Global and local explainability; Global importance: PFI, H-statistic; Global dependence: PDP, ALE; Local methods: ICE, LIME, SHAP), 'Diagnostic Suite' (Model validation and outcome testing; Basic Tests: Accuracy, Weakspot, Overfit; 3R Tests: Reliability, Robustness, Resilience; Fairness test; Segmented test; Scored test), 'Model Comparison' (Benchmarking through diagnostics; Regression models; Binary classification models; Model fairness comparison), and 'Low-Code Case Studies' (PiML workflow and experimentation; Example: Bikesharing Data; Example: CaliforniaHousing Data; Example: TaiwanCredit Data; Fairness Simulation Study 1; Fairness Simulation Study 2).

<https://selfexplainml.github.io/PiML-Toolbox/>

See also: <https://github.com/SelfExplainML/PiML-Toolbox> (workshop materials) and <https://piml.medium.com/> (blog tutorials)

Hands-on PiML Training

Installation: pip install piml (224K+ PyPI downloads in 2 years)

Session 1: Model Interpretability

Jupyter Notebook: [Github link](#), [Open in Google Colab](#)

- Data Pipeline
- ML Models: XGB and MLP
- Post-hoc Explainability
- Inherently Interpretable Benchmark Models
 - FANOVA models
 - Inherent interpretability
 - Monotonic constraints

Session 2: Outcome Analysis

Jupyter Notebook: [Github link](#), [Open in Google Colab](#)

- Prediction Accuracy
- Hyperparameter Tuning
- Weakness Detection
- Overfitting Analysis
- Prediction Uncertainty
- Robustness and Resilience
- Bias and Fairness

Explainability Test

- **Post-hoc explainability test** is model-agnostic, i.e., it works for any pre-trained model.
 - Useful for explaining black-box models; but need to use with caution (there is no free lunch).
 - Post-hoc explainability tools sometimes have pitfalls, challenges and potential risks.
- **Local explainability tools** for explaining an individual prediction
 - **ICE** (Individual Conditional Expectation) plot
 - **LIME** (Local Interpretable Model-agnostic Explanations)
 - **SHAP** (SHapley Additive exPlanations)
- **Global explainability tools** for explaining the overall impact of features on model predictions
 - **Examine relative importance of variables**: **VI** (Variable Importance), **PFI** (Permutation Feature Importance), **SHAP-FI** (SHAP Feature Importance), **H-statistic** (Importance of two-factor interactions), etc.
 - **Understand input-output relationships**: 1D and 2D **PDP** (Partial Dependence Plot) and **ALE** (Accumulated Local Effects).

Post-hoc Explainability vs. Inherent Interpretability

- **Post-hoc explainability** is model agnostic, but there is no free lunch. According to Cynthia Rudin, use of auxiliary post-hoc explainers creates “**double trouble**” for black-box models.
- Various post-hoc explanation methods, including VI/FI, PDP, ALE, ... (for global explainability) and LIME, SHAP, ... (for local explainability), **often produce results with disagreements**.
- Lots of academic discussions about pitfalls, challenges and potential risks of using post-hoc explainers.
- This echoes CFPB Circular 2022-03 (May 26, 2022): Adverse action notification requirements in connection with credit decisions based on complex algorithms¹.

- **Inherent interpretability** is intrinsic to a model. It facilitates gist and intuitiveness for human insightful interpretation. It is important for evaluating a model’s **conceptual soundness**.
- Model interpretability is a loosely defined concept and can be hardly quantified. Sudjianto and Zhang (2021)² proposed a **qualitative rating assessment** framework for ML model interpretability.
- **Interpretable model design**: a) interpretable feature selection and b) interpretable architecture constraints³ such as additivity, sparsity, linearity, smoothness, monotonicity, visualizability, projection orthogonality, and segmentation degree.

¹ CFPB Circular 2022-03 Footnote 1: “While some creditors may rely upon various **post-hoc explanation methods**, such explanations approximate models and creditors must still be able to **validate the accuracy** of those approximations, which **may not be possible with less interpretable models**.” [consumerfinance.gov](https://www.consumerfinance.gov)

² Sudjianto and Zhang (2021): Designing Inherently Interpretable Machine Learning Models. [arXiv: 2111.01743](https://arxiv.org/abs/2111.01743)

³ Yang, Zhang and Sudjianto (2021, IEEE TNNLS): Enhancing Explainability of Neural Networks through Architecture Constraints. [arXiv: 1901.03838](https://arxiv.org/abs/1901.03838)

Inherently Interpretable FANOVA Models

- One effective way is to design inherently interpretable models by the functional ANOVA representation

$$g(\mathbb{E}(y|\mathbf{x})) = g_0 + \sum_j g_j(x_j) + \sum_{j < k} g_{jk}(x_j, x_k) + \sum_{j < k < l} g_{jkl}(x_j, x_k, x_l) + \dots$$

It additively decomposes into the overall mean (i.e., intercept) g_0 , main effects $g_j(x_j)$, two-factor interactions $g_{jk}(x_j, x_k)$, and higher-order interactions ...

- GAM main-effect models: Binning Logistic, XGB1, GAM (estimated using Splines, etc.)
- GAMI main-effect plus two-factor-interaction models:
 - **EBM** (Nori, et al. 2019) → explainable boosting machine with shallow trees
 - **XGB2** (Lengerich, et al. 2020) → boosted trees of depth 2 with effect purification
 - **GAMI-Net** (Yang, Zhang and Sudjianto, 2021) → specialized neural nets
 - **GAMI-Lin-Tree** (Hu, et al. 2023) → specialized boosted linear model-based trees
- **PiML Toolbox** integrates GLM, GAM, XGB1, XGB2, EBM and GAMI-Net, and provides their inherent interpretability.

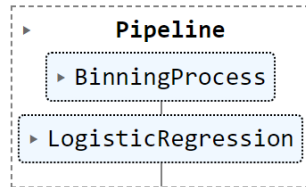
Binning Logistic vs. XGB1

```
from sklearn.pipeline import Pipeline
from optbinning import BinningProcess
from sklearn.linear_model import LogisticRegression

feature_names = exp.get_feature_names()
train_x, train_y, _ = exp.get_data(train=True)

lr = Pipeline(steps=[('Step 1', BinningProcess(feature_names)),
                     ('Step 2', LogisticRegression())])

lr.fit(train_x, train_y.ravel())
```



```
# Register it as PiML pipeline
tmp = exp.make_pipeline(model=lr)
exp.register(tmp, "BinningLogistic")
exp.model_diagnose(model="BinningLogistic", show='accuracy_table')
```

	ACC	AUC	Recall	Precision	F1
Train	0.6787	0.7374	0.7144	0.6716	0.6923
Test	0.6760	0.7341	0.7142	0.6728	0.6929
Gap	-0.0027	-0.0034	-0.0002	0.0012	0.0006

```
from piml.models import XGB1Classifier

exp.model_train(XGB1Classifier(), name='XGBoostDepth1')

exp.model_diagnose(model="XGBoostDepth1", show='accuracy_table')
```

	ACC	AUC	Recall	Precision	F1
Train	0.6940	0.7531	0.7313	0.6851	0.7075
Test	0.6883	0.7465	0.7298	0.6828	0.7055
Gap	-0.0057	-0.0066	-0.0015	-0.0023	-0.0019

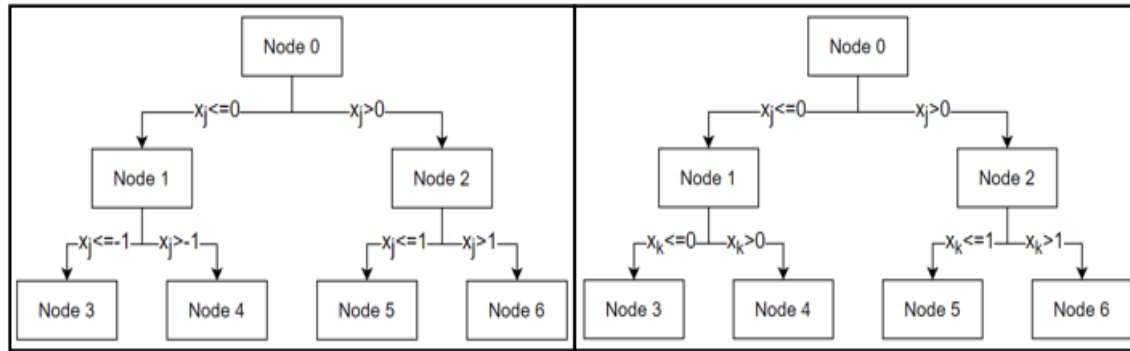
- Binning Logistic is a GAM main effect model with piecewise constant basis functions (feature engineering). It performs manual binning one variable at a time.
- XGB1 is also a GAM main effect model of the same type. It performs automated binning jointly for all variables.
- Both GAM models are inherently interpretable, easy to quantify feature importance and draw main effect plots.

XGB1, XGB2 and Beyond

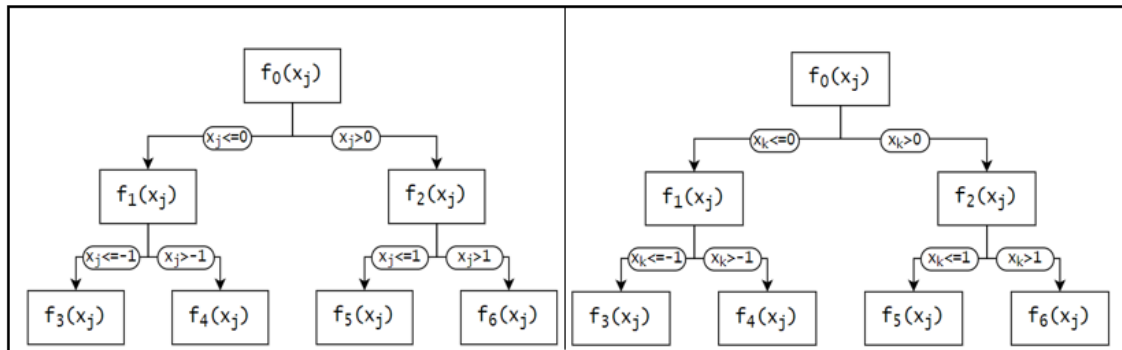
- **Proposition:** A depth- K tree-ensemble can be reformulated to an FANOVA model with main effects and k -way interactions with $k \leq K$.
- Examples: XGB1 is GAM with main effects; XGB2 is GAMI with main effects plus two-factor interactions.
- Three-step unwrapping technique for tree ensembles (e.g., RF, GBDT, XGBoost, LightGBM, CatBoost):
 1. **Aggregation:** all leaf nodes with the same set of k distinct split variables sum up to a raw k -way interaction.
 2. **Purification:** recursively cascade effects from high-order interactions to lower-order ones to obtain a unique FANOVA representation subject to hierarchical orthogonality constraints (Lengerich, et al., 2020).
 3. **Attribution:** quantify the importance of purified effects either locally (for a sample) or globally (for a dataset).
- Strategies to enhance model (e.g., XGBoost) interpretability without sacrificing model performance
 - XGB hyperparameters: max_tree_depth, max_bins, candidate interactions, monotonicity, L1/L2 regularization, etc.
 - Pruning of purified effects: effect selection by L1 regularization, forward and backward selection with early stopping
 - Other strategies such as post-hoc smoothing of purified effects, local flattening, and boundary effect adjustment.

EBM, GAMI-Lin-Tree, GAMI-Net

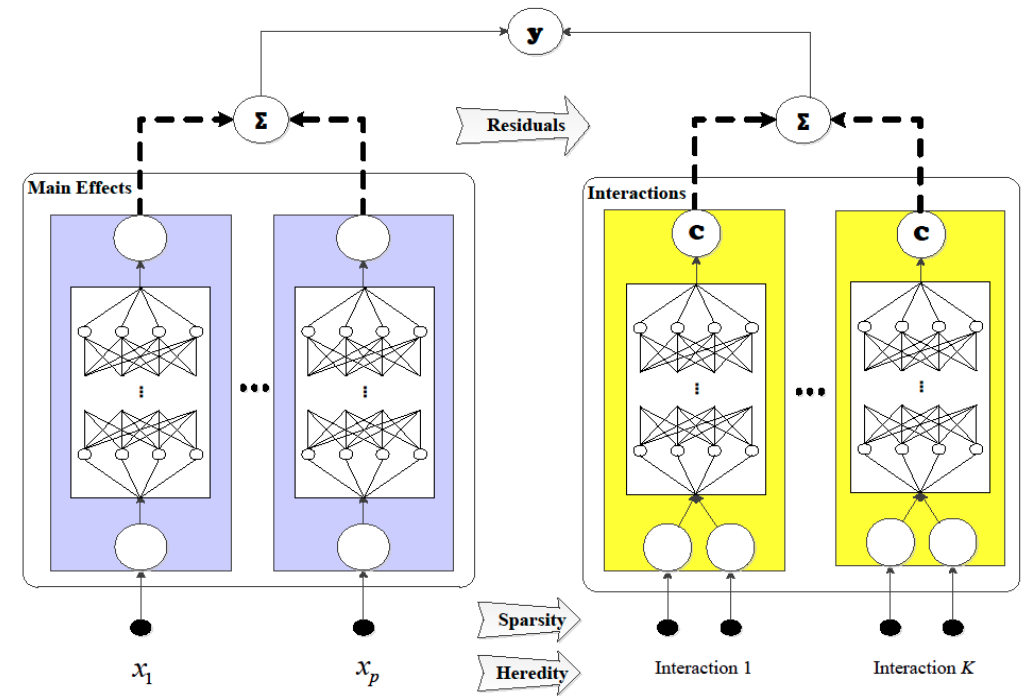
$$\text{GAMI: } g(E(y|x)) = \mu + \sum h_j(x_j) + \sum f_{jk}(x_j, x_k)$$



[EBM \(Nori, et al. 2019\)](#)



[GAMI-Lin-Tree \(Hu, et al. 2023\)](#)



GAMI-Net: An explainable neural network based on generalized additive models with structured interactions

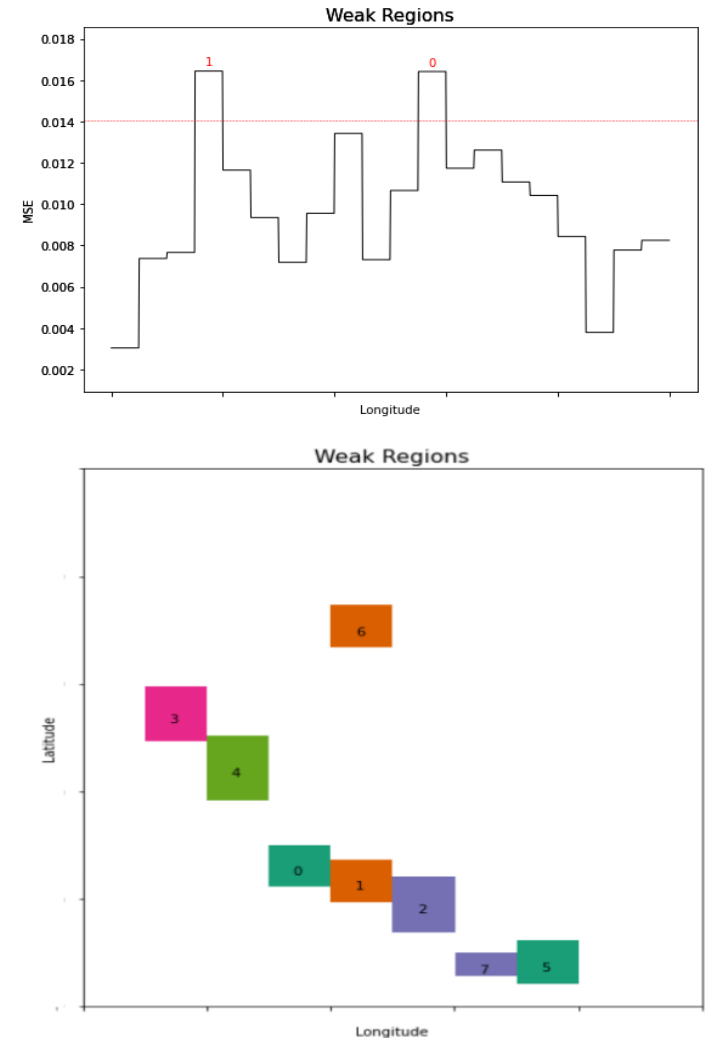
[Z Yang, A Zhang, A Sudjianto](#) - Pattern Recognition, 2021 - Elsevier

... models with structured interactions (**GAMI-Net**) is proposed to pursue a good balance between prediction accuracy and model interpretability. **GAMI-Net** is a disentangled feedforward ...

☆ Save 📄 Cite Cited by 112 Related articles All 4 versions

Weakness Detection by Error Slicing

1. **Specify an appropriate metric** based on individual prediction residuals: e.g., MSE for regression, ACC/AUC for classification, train-test performance gap (for checking overfit), prediction interval bandwidth, ...
2. Specify 1 or 2 slicing features of interest;
3. Evaluate the metric for each sample in the target data (training or testing) as pseudo responses;
4. **Segment the target data** along the slicing features, by
 - a) [Unsupervised] Histogram slicing with equal-space binning, or
 - b) [Supervised] fitting a decision tree to generate the sub-regions
5. **Identify the sub-regions** with average metric exceeding the pre-specified threshold, subject to minimum sample condition.



Prediction Uncertainty by Reliability Test

- Prediction uncertainty is important to understand where the model produces less reliable prediction:

Wider prediction interval \rightarrow Less reliable prediction

- Quantification of prediction uncertainty can be done through **Split Conformal Prediction** under the exchangeability assumption:

Given a pre-trained model $\hat{f}(x)$, a hold-out calibration data $\mathcal{X}_{\text{calib}}$, a pre-defined conformal score $S(x, y, \hat{f})$ and the error rate α (say 0.1)

- Calculate the score $S_i = S(x, y, \hat{f})$ for each sample in $\mathcal{X}_{\text{calib}}$;
- Compute the calibrated score quantile

$$\hat{q} = \text{Quantile}\left(\{S_1, \dots, S_n\}; \frac{\lceil (n+1)(1-\alpha) \rceil}{n+1}\right);$$

- Construct the prediction set for the test sample x_{test} by

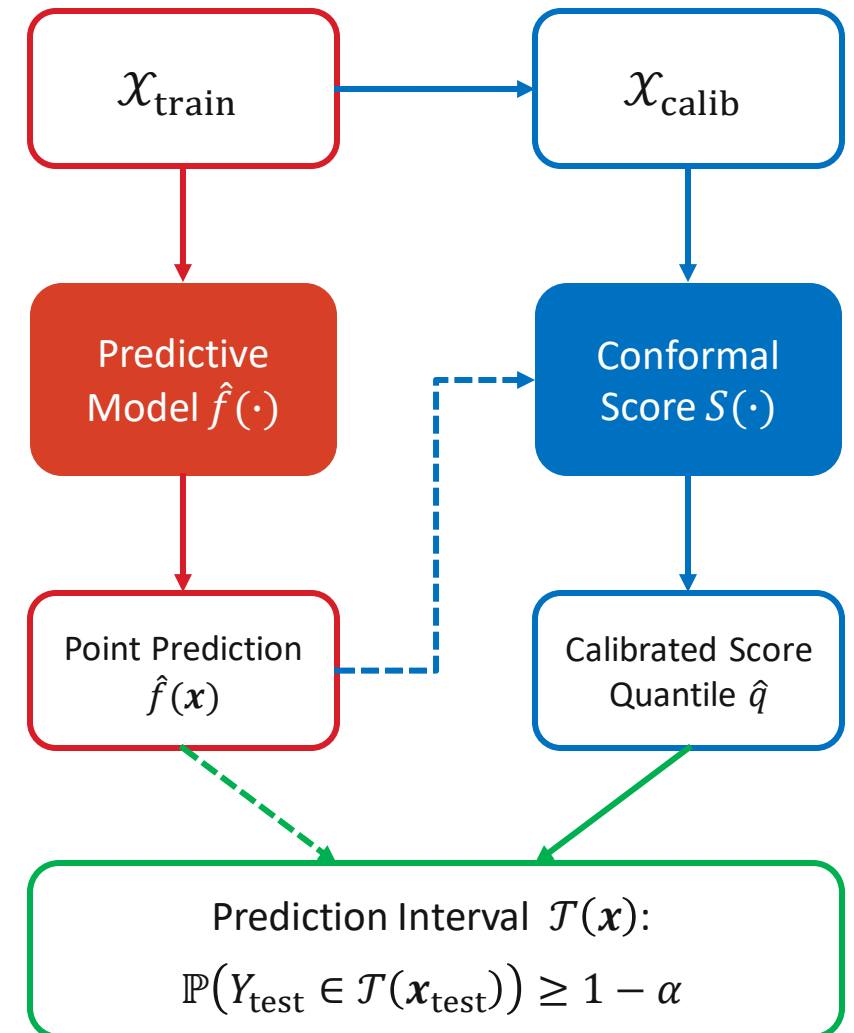
$$\mathcal{T}(x_{\text{test}}) = \{y: S(x_{\text{test}}, y, \hat{f}(x_{\text{test}})) \leq \hat{q}\}.$$

Under the exchangeability condition of conformal scores, we have that

$$1 - \alpha \leq \mathbb{P}(Y_{\text{test}} \in \mathcal{T}(x_{\text{test}})) \leq 1 - \alpha + \frac{1}{n+1}.$$

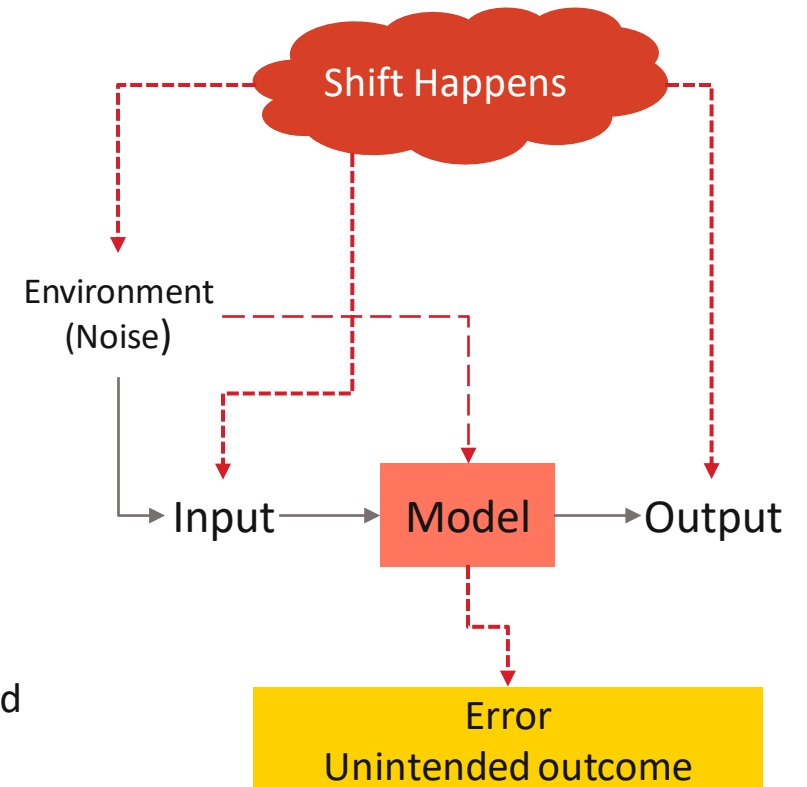
This provides the prediction bounds with α -level acceptable error.

- PiML team implements a sophisticated residual-quantile conformal method for regression models. See details in [this tutorial](#).



Robustness and Resilience Tests

- Train-test data split (i.i.d.) leads to over-optimism of model performance, since model in production will be exposed to data distribution shift.
- **Robustness test:** evaluate the performance degradation under covariate noise perturbation:
 - Perturb testing data covariates with small random noise;
 - Assess model performance of perturbed testing data.
 - Overfitting models often perform poorly in changing environments.
- **Resilience test:** evaluate the performance degradation under distribution drift scenarios
 - Scenarios: worst-sample, worst-cluster, outer-sample, hard-sample
 - Measure distribution drift (e.g., PSI) of variables between worst performing sample and the remaining sample.
 - Variables with notable drift are deemed to be sensitive in the resilience test.



Streamlined Validation of AI/ML Models

- Developing an effective AI/ML model risk management program: ValOps platform
- **Key objective:** streamline validation process to reduce cycle time and enable automated validation/ monitoring for AI/ML models (including dynamically updating models).
- **Standard model wrapping** - provides a standardized model management protocol for managing data and model complexity and diversity.
- **Standard validation tests** - centralizes test codes and validation suites for data quality check, evaluation of conceptual soundness, and outcome analysis.





Thank you

Aijun Zhang, Ph.D.

Email: Aijun.Zhang@wellsfargo.com

LinkedIn: <https://www.linkedin.com/in/ajzhang/>