

# Explainable, interpretable, and trustworthy AI for an intelligent digital twin: A case study on remaining useful life

Kazuma Kobayashi, SB Alam  
*Engineering Applications of Artificial Intelligence*. 2024



# Abstract



- ✓ Artificial intelligence (AI) and Machine learning (ML) are increasingly used for digital twin development in energy and engineering systems.
- ✓ Need to be unbiased, interpretable and explainable.
- ✓ Explains XAI and IML.
- ✓ Justify the role of XAI/IML for digital twin systems for the prediction of RUL (Remaining useful Life), both locally and globally.

# Problem Statement



- ✓ Prognostics and health management (PHM) use DTs to monitor and manage system health utilizing statistical algorithms and models.
- ✓ DTs help in:
  - ✓ Conducting condition based maintenance.
  - ✓ Making maintenance decisions.
  - ✓ Anticipate future failures in advance.
- ✓ XAI can help increase trust the RUL prediction of an Intelligent DT system (physical asset such machinery or infrastructure).

# Digital Twin System



- ✓ virtual representations of physical systems.
- ✓ replicates a physical asset in the virtual environment, including its functionality, features, and behavior.
- ✓ DT's technology can reflect, mimic, and forecast the status of the operating physical system in real-time and is one of the best tools for RUL prediction.



- ✓ Incorporating AI/ML in DTs can facilitate in
  - ✓ risk informed decision making related to RUL of a component.
  - ✓ Streamlining high performing simulations.
- ✓ With ML comes uncertainty and lack of trustworthiness especially in complex non-linear models.
- ✓ Need explainability in ML ( XAI/IML).
- ✓ the use of explainable, interpretable, and trustworthy AI is crucial for accurately predicting the remaining useful life in an intelligent DT system

# Explainable AI



- ✓ Explainable [AI](#) is used to describe an AI model, its expected impact and potential biases.
- ✓ XAI is important for several reasons:
  - ✓ Builds trust by explaining the rational behind a decision.
  - ✓ Transparency.
  - ✓ International regulations requires to explain decision making process.
  - ✓ Improves user's experience.
  - ✓ Debug and improve.
  - ✓ Increase reliability and accuracy of prediction of a ML model.

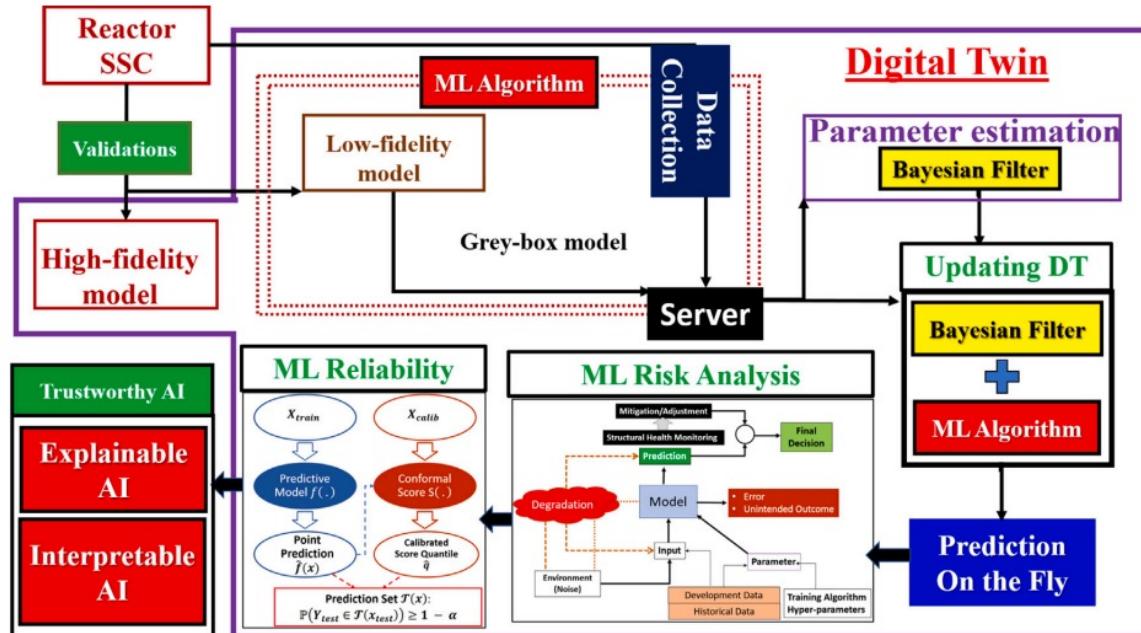


- ✓ Explainability allows users to understand the factors most important in the model's prediction.
- ✓ Interpretability allows the model's predictions to be easily understood by non-technical users.
- ✓ Trustworthiness ensures that the model is not making predictions based on irrelevant or misleading features and is unbiased.

# Intelligent Digital Twin framework



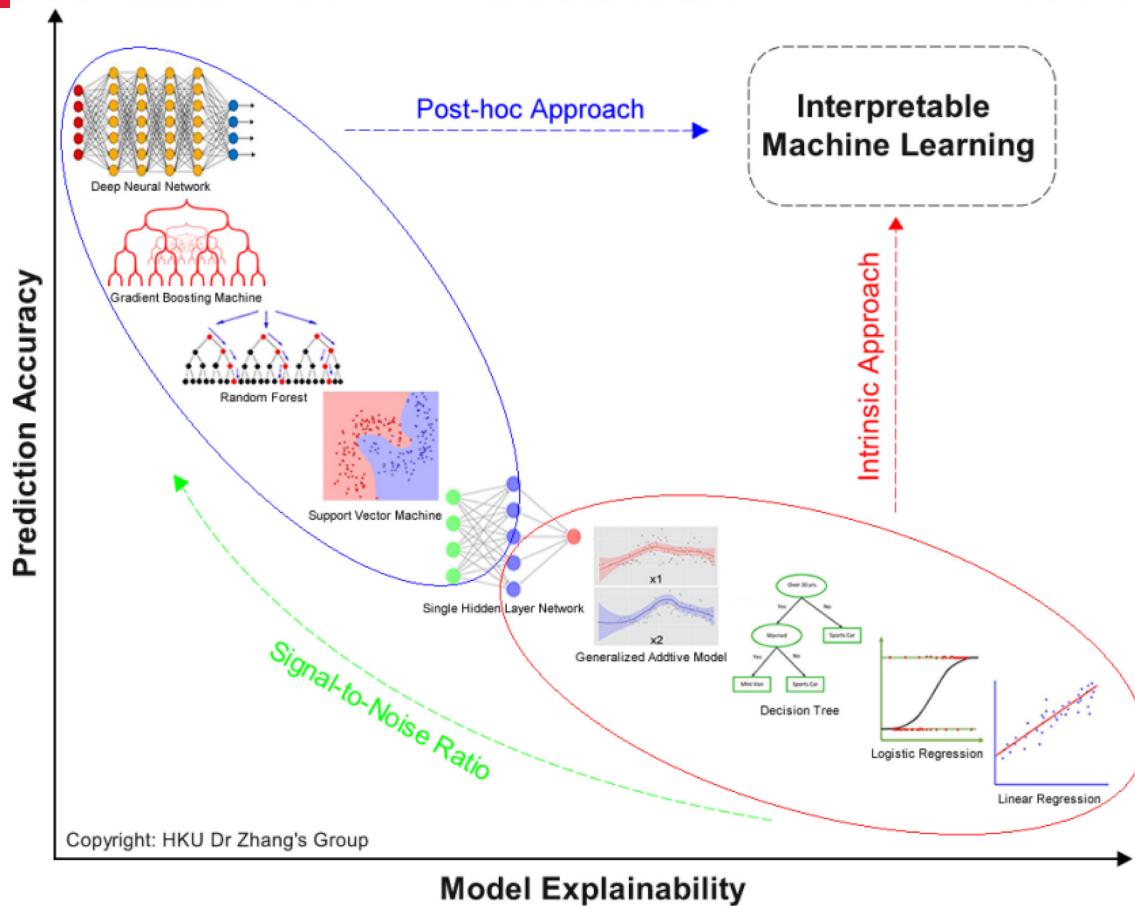
- ✓ An Intelligent DT framework has five essential components.
  1. Prediction module
  2. System Update Module.
  3. Data processing module
  4. Visualization module
  5. Decision Making module.
- ✓ System update and prediction module require sophisticated ML algorithms.



# Interpretable AI



- Interpretable AI refers to the AI system's capacity to justify its judgement and predictions.
- Enables to comprehend the decision-making process for the user.
- Two strategies to develop such systems
  - Use transparent models (decision trees, linear models).
  - Post hoc explanations (feature importance).



# XAI and Interpretable AI



- XAI means creating AI systems that can explain their actions (transparent models or post hoc explanations).
- Interpretable AI refers to interpreting and comprehending a model's outcomes.
- XAI is a broader concept of developing explainable and transparent AI system whereas Interpretable explains the relationship between input and output variables related to decision making.

# Interpretable methods



- ReLu-DNN
- Explainable boosting machine (EBM).
- Fast interpretable greedy-tree sums (FIGS).
- Partial dependence plot (PDP).
- Explainable neural networks(XNN).



- This study uses inherently interpretable models.
- Dataset used : PHM08 (Generated by C-MAPSS simulation tool developed by NASA to mimic turbofan engine degradation under different operating condition)

# Feature selection



- Reduces a large set of features to a manageable subset.
- It can improve model's precision and interpretability.
- Removes noise and reduce overfitting.
- Methods:
  - Pearson correlation
  - Distance correlation
  - Feature importance

# Pearson correlation



- Statistical technique to determine the linear relationship between two continuous variables.
- Computed for each feature wrt target variable.
- Simple and straightforward method to implement.
- Most frequently used metric.
- Not good with non-linear interactions in the data.

# Distance correlation

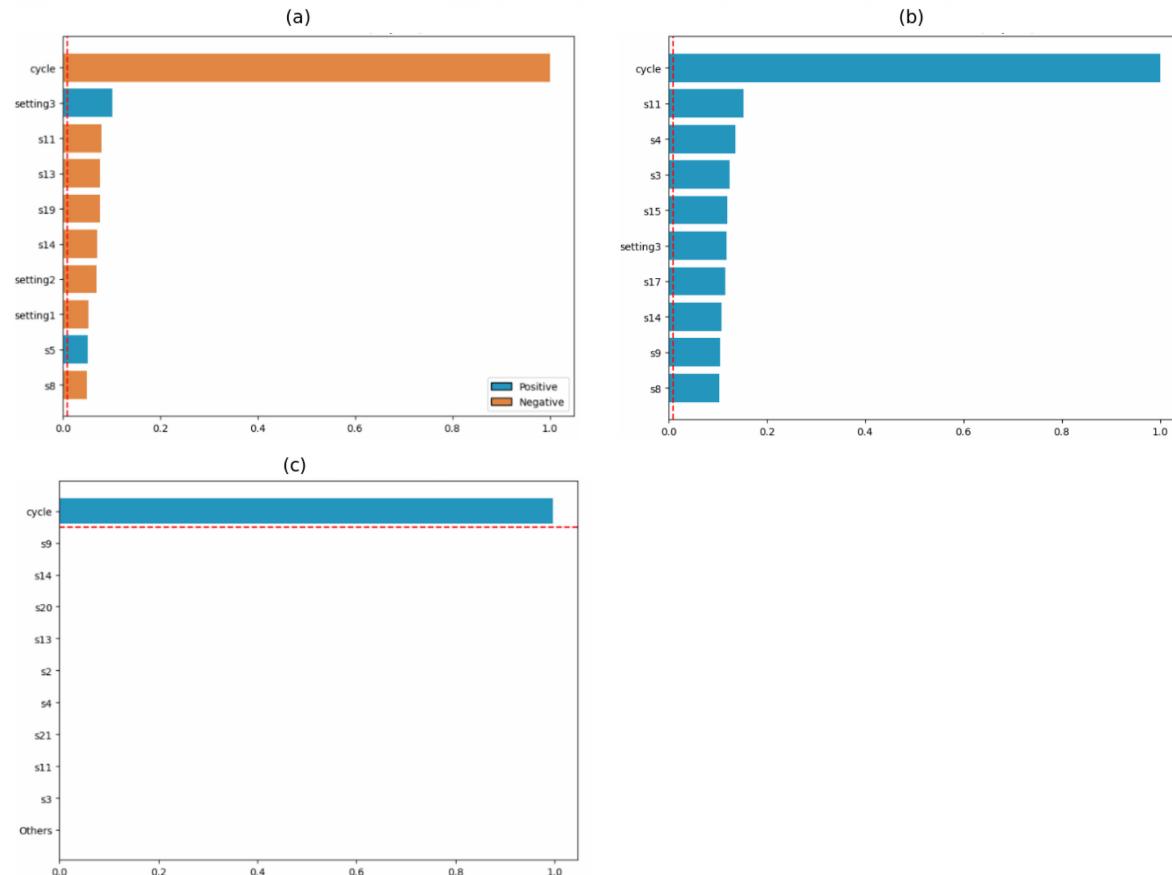


- Quantify relationship between two variables.
- Determine the most significant parameter for predicting the target variable.
- Can reveal non-linear relationships between variables.
- Computationally demanding.

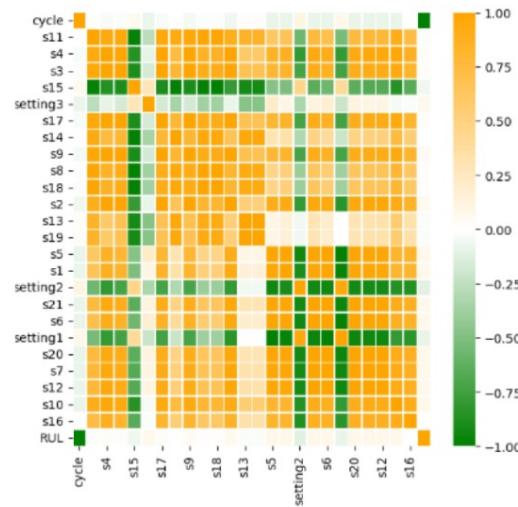
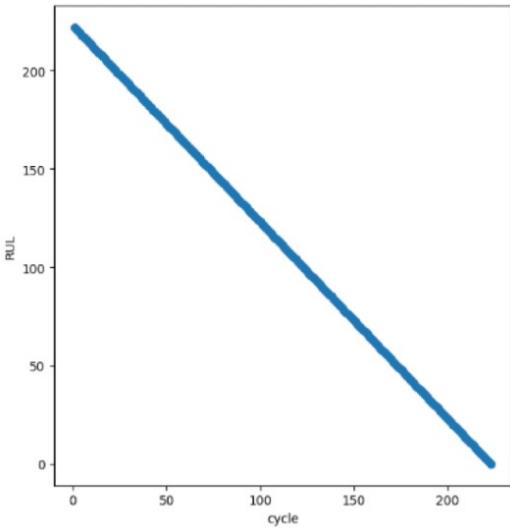
# LGBM based feature importance



- Light gradient boosting machine.
- LGBM utilizes decision trees as its base model and train multiple trees through an iterative process.
- The algorithm assigns higher weights to observations poor predicted by previous trees, which helps improve accuracy.
- Capable of handling large datasets, missing values and categorical features.



# Multivariate correlation Heatmap



# Global explainability

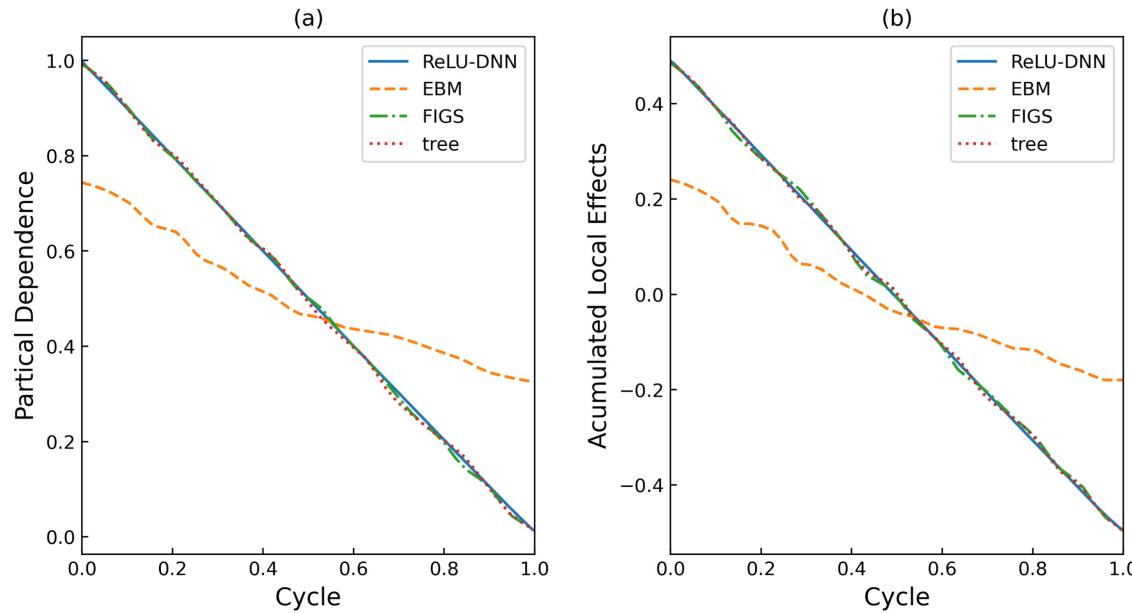


## Partial Dependence plot

- PDPs are an interpretability ML method to understand the relationship between a single feature and the target variable.
- Capable of finding linear or more complex relationships in the data.
- Identify confounding factors that might obscure the connection between an intriguing quality and the desired result.
- Easy to understand (works well for non-technical audiences).
- Provide only partial picture if a feature is linked with other features.



- Accumulated Local effects
- Plot the model's average prediction as a function of the feature of interest while maintaining the observed values of all other features.
- Capable to reveal linear and more complex relationship between features and target variable.
- Rely on observed values rather than mean and median to fix the values of other features, hence less sensitive to this decision.



# Local Explainability



Local explainability model agnostic explanations (LIME).

- Model agnostic.
- Provides local explanations for a single data point or a small set of data points.

- Toy example to present intuition for LIME.
- The black box model's complex decision function  $f$  is represented by the blue/pink background.
- The bold red cross is the instance being explained.
- LIME samples instances, gets predictions using  $f$  and weighs them by the proximity to the instance being explained (represented by size).
- The dashed line is the learned explanation.

