

Assignment Week 12

Critical Analysis of Microsoft Responsible AI Toolbox Tools

InterpretML

Mode of Usage: Refer to documentation or integrate via Python packages for model interpretability, fairness, or causal analysis.

Key Benefits:

- Provides interpretability to machine learning models by offering visual explanations.
- Supports both black-box and glass-box models for better transparency.
- Allows users to understand global and local feature importance.
- Helps identify potential bias or unfair behavior in models.
- Can be integrated into model development pipelines for continuous evaluation.

Fairlearn

Mode of Usage: Refer to documentation or integrate via Python packages for model interpretability, fairness, or causal analysis.

Key Benefits:

- Mitigates unfairness in classification and regression models.
- Provides disparity metrics and fairness dashboards.
- Enables decision-makers to assess trade-offs between fairness and accuracy.
- Useful in regulated industries like finance and healthcare.
- Encourages ethical practices in AI deployment.

DiCE

Mode of Usage: Refer to documentation or integrate via Python packages for model interpretability, fairness, or causal analysis.

Key Benefits:

- Generates counterfactual explanations for individual predictions.
- Helps understand what minimal changes can lead to different outcomes.
- Promotes transparency by showing actionable changes.
- Supports multiple models and platforms.
- Useful for debugging and improving model decisions.

Assignment Week 12

Error Analysis

Mode of Usage: Refer to documentation or integrate via Python packages for model interpretability, fairness, or causal analysis.

Key Benefits:

- Helps diagnose model failures and understand error distribution.
- Identifies subgroups where models perform poorly.
- Supports slicing and filtering for better analysis.
- Improves model performance by targeted retraining.
- Valuable for continuous improvement in industrial settings.

EconML

Mode of Usage: Refer to documentation or integrate via Python packages for model interpretability, fairness, or causal analysis.

Key Benefits:

- Estimates causal effects using machine learning.
- Designed for economic decision-making and policy modeling.
- Useful in scenarios where correlation does not imply causation.
- Supports treatment effect estimation and uplift modeling.
- Enables more accurate business decisions based on causal insights.