# Monitor a model for quality and fairness

# Overview

AI is used in everyday life to support human decision-making.
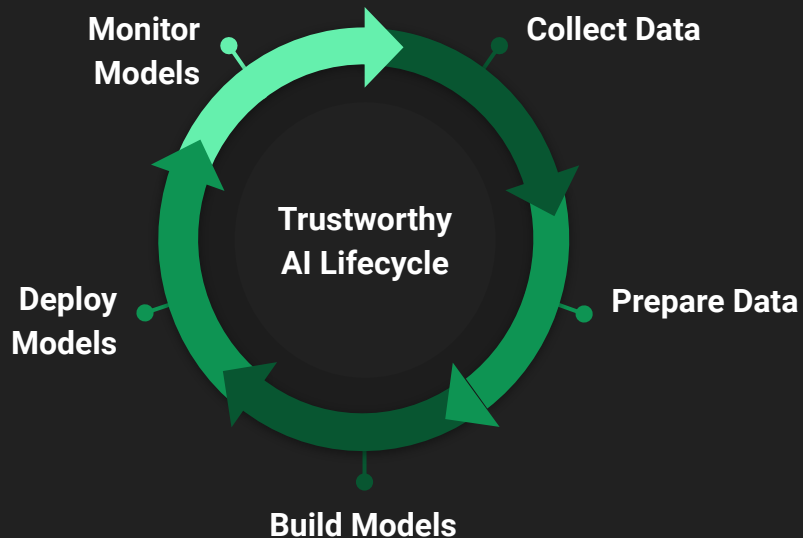
AI is computer algorithm. For many people, it's a block box.

How do we operationalize AI with confidence?
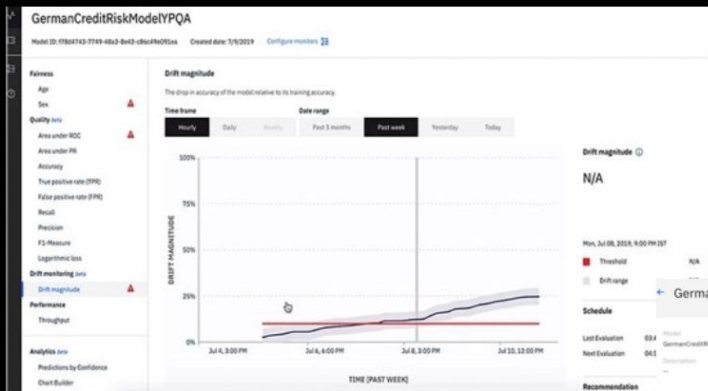
What is Fairness?

How to measure bias?

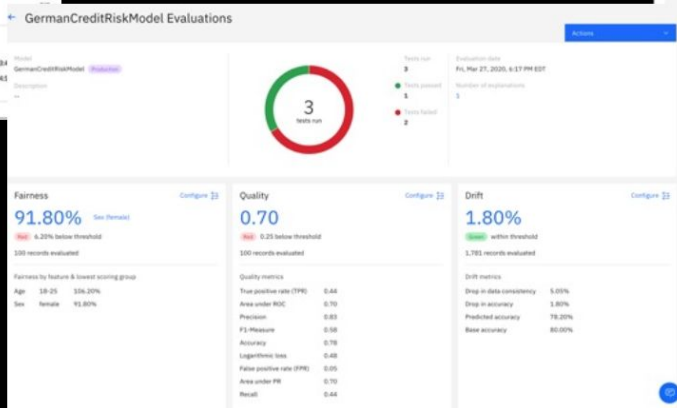How to enable responsible use of AI?

**Monitor Models**

**Collect Data**

**Trustworthy AI Lifecycle**

**Prepare Data**

**Deploy Models**

**Build Models**

# IBM Watson OpenScale

# Implement responsible, explainable AI
*Mitigate drift, bias, and model risk*



Drift: monitor model drift by hourly, daily or weekly

Model risk evaluation: fairness, quality and drift metrics to share model insights

Explain transactions: Determine what features reach different outcomes

# Demo

1. Tutorial scenario: https://github.com/IBM/ai-data-workshop/tree/main/monitor-model-with-openscale
2. Notebook
3. Evaluation
4. Insight Dashboard
   a. Fairness monitor
   b. Quality monitor
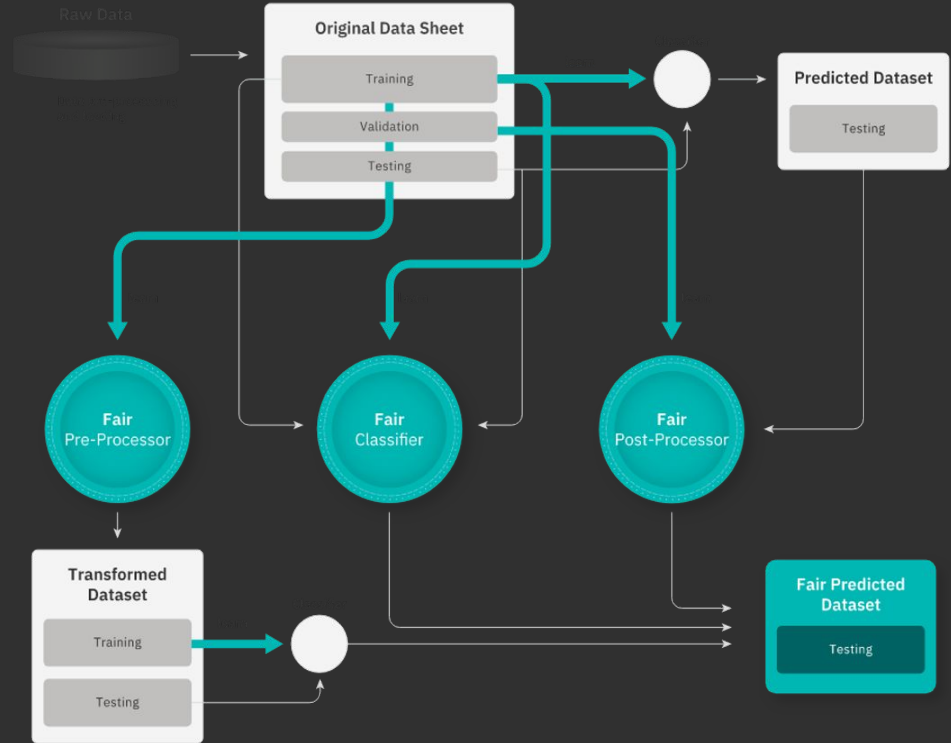   c. Explain transaction

# AI Fairness 360

# AI Fairness 360 (Open Source)

AIF360 toolkit is an open-source library to help detect and remove bias in machine learning models. AIF360 translates algorithmic research from the lab into practice. Applicable domains include finance, human capital management, healthcare, and education.

The AI Fairness 360 Python package includes a comprehensive set of metrics for datasets and models to test for biases, explanations for these metrics, and algorithms to mitigate bias in datasets and models.

Toolbox
Fairness metrics (70+)
Fairness metric explanations
Bias mitigation algorithms (10+)

http://aif360.mybluemix.net/

Raw Data

**Original Data Sheet**
- Training
- Validation
- Testing

**Predicted Dataset**
- Testing

**Fair Pre-Processor**

**Fair Classifier**

**Fair Post-Processor**

**Transformed Dataset**
- Training
- Testing

**Fair Predicted Dataset**
- Testing

# AIF360 Demo

1. Import Notebook
2. Run notebook