

# Weakness Detection in ML Models

Using ModEva Framework for Identifying Underperforming Regions

Machine Learning Expert

July 5, 2025

# Outline

- 1 Introduction to Weakness Detection
- 2 Weakness Detection in ModEva
- 3 Error Slicing for Weakness Detection
- 4 Advanced Slicing Techniques
- 5 Model Comparison
- 6 Advanced Diagnostic Approaches
- 7 Acting on Weakness Insights

# Understanding Weakness Detection

## Definition

Weakness detection is the process of identifying areas in the input space where a machine learning model underperforms or makes incorrect predictions.

## Characteristics of Weak Regions:

- High residual errors
- Poor prediction accuracy
- Patterns of bias or inconsistency
- Unexpected behavior on certain data subsets

**Goal:** Identify and understand these regions to improve model reliability, robustness, and fairness across all data segments.

# Why Weakness Detection is Important

## Improve Model Performance

Identify struggling areas and implement targeted improvements

## Guide Data Collection

Pinpoint where additional data can enhance model performance

## Enhance Trustworthiness

Address systematic biases and recurring errors to build confidence

## Mitigate Risks

Detect and address weaknesses before deployment

# Key Approaches to Weakness Detection

## Residual Analysis

- Examine prediction errors
- Identify patterns in residuals
- Locate regions with systematic errors

## Data Slicing

- Divide dataset into smaller subsets
- Evaluate performance in each slice
- Compare across slices

## Feature Sensitivity

- Identify features linked to poor performance
- Analyze feature interactions
- Understand sensitivity to feature changes

## Visualization

- Plot performance across feature ranges
- Visualize error distributions
- Highlight underperforming regions

# Introduction to Error Slicing

## Concept

Error slicing involves dividing data into segments to assess model performance in specific regions of the feature space.

## Primary Binning Methods in ModEva:

- ❶ **Uniform Binning:** Equal-sized intervals across feature range
- ❷ **Quantile Binning:** Equal number of samples in each bin
- ❸ **Automatic Binning:** Using tree-based methods (XGBoost) to find optimal splits
- ❹ **User-defined Binning:** Custom-defined bin boundaries

**Goal:** Identify specific feature ranges or data segments where model performance drops below acceptable thresholds.

# Uniform Binning

## How It Works:

- Divides feature range into equal-sized intervals
- Simple and easy to interpret
- Works well with uniformly distributed features

## Disadvantages:

- Can result in empty or sparse bins for skewed distributions

```
1 # Analyze residual feature
  importance
2 results = ts.
    diagnose_residual_fi(
3     method="uniform")
4 results.plot()
5
6 # Uniform binning numerical
  feature
7 results = ts.
    diagnose_slicing_accuracy(
8     features=(("LIMIT_BAL", ),
9               ("PAY_1", )),
    method="uniform",
10    bins=10,
11    metric="AUC",
12    threshold=0.65)
13 results.plot()
14
```

# Quantile Binning

## How It Works:

- Divides data so each bin has equal number of samples
- Handles skewed distributions effectively
- Ensures equal representation in bins

## Disadvantages:

- Bin widths may vary, complicating interpretation
- Sensitive to outliers

```
1 # Quantile binning
2 results = ts.
   diagnose_slicing_accuracy(
3     features="LIMIT_BAL",
4     method="quantile",
5     bins=10,
6     metric="AUC",
7     threshold=0.65)
8 results.plot()
9
```

*This generates a visualization of performance across quantile bins*



# Automatic Binning with Tree-Based Models

## How It Works:

- Uses depth-1 or depth-2 XGBoost trees
- Automatically finds optimal split points
- Splits based on relationship with target

## Advantages:

- Optimized for target performance
- Captures meaningful feature-target relationship
- More intelligent than fixed-width binning

```
1 # Automatic binning
2 results = ts.
   diagnose_slicing_accuracy(
3     features="LIMIT_BAL",
4     method="auto-xgb1",
5     bins=10,
6     metric="AUC",
7     threshold=0.75)
8 results.plot() # Display
   results in plot
9 results.table # Display
   results in table
10
```

*This automatically identifies  
meaningful bins based on model  
performance*

# Custom Binning

```
1 # Custom binning
2 results = ts.diagnose_slicing_accuracy(
3     features="LIMIT_BAL",
4     method="precompute",
5     bins={"LIMIT_BAL": (0.0, 50000, 1000000)},
6     metric="AUC")
7 results.table # Display results in table
8
```

*This allows defining custom bin boundaries for specific feature ranges of interest*

## When to use custom binning:

- When specific feature thresholds are meaningful for the business context
- To focus on particular regions of interest (e.g., high-value customers)
- When domain expertise suggests particular breakpoints
- For comparing with established industry benchmarks

# Multiple Feature Slicing

```
1 # Slicing for a set of features
2 results = ts.diagnose_slicing_accuracy(
3     features=(("PAY_1", ), ("BILL_AMT1",), ("PAY_AMT1", )),
4     method="quantile",
5     metric="AUC",
6     threshold=0.6)
7 results.table
8
```

*This analyzes performance across multiple features independently*

## Benefits:

- Provides a comprehensive view of performance across multiple feature dimensions
- Identifies which features are most associated with weak performance areas
- Allows prioritization of feature improvement efforts

# Feature Interaction Slicing

```
1 # 2-Feature interaction slicing
2 results = ts.diagnose_slicing_accuracy(
3     features=("PAY_1", "PAY_AMT1"),
4     method="uniform",
5     bins=10,
6     metric="AUC",
7     threshold=0.5)
8 results.table
9
```

*This examines performance across combinations of feature values*

## Why interaction slicing matters:

- Models may perform well for individual feature ranges but struggle with specific combinations
- Reveals complex relationships that cause model weaknesses
- Identifies rare but important feature interaction scenarios
- Helps detect potential biases in specific feature combinations

# Analyzing Weak Regions

```
1 # Retrieving samples below threshold value
2 from modeva.testsuite.utils.slicing_utils import
   get_data_info
3 data_info = get_data_info(res_value=results.value)[("PAY_1",
   "PAY_AMT1")]
4 data_info
5
6 # Testing distribution difference between weak samples and
   the rest
7 data_results = ds.data_drift_test(
8     **data_info,
9     distance_metric="PSI",
10    psi_method="uniform",
11    psi_bins=10)
12 data_results.plot("summary")
13
```

*This identifies data points in weak regions and compares their distribution to the overall dataset*

# Visualizing Weak Regions

```
1 # To get the list of figure names in the "data_results"  
   object  
2 data_results.get_figure_names()  
3  
4 # Example of plotting from the list of figures  
5 data_results.plot(('density', 'PAY_1'))  
6
```

*This visualizes the distribution differences between weak samples and normal samples*

## Key visualizations for understanding weak regions:

- Density plots showing feature distributions
- PSI (Population Stability Index) summary plots
- Feature importance plots for weak regions
- Interaction heatmaps highlighting problematic combinations

# Comparing Weaknesses Across Models

```
1 # Compare models on numerical features
2 tsc = TestSuite(ds, models=[model_lgbm, model_xgb])
3 results = tsc.compare_slicing_accuracy(
4     features="PAY_AMT1",
5     method="quantile",
6     bins=10,
7     metric="AUC")
8 results.plot()
9
```

*This compares performance of different models across feature slices*

## Benefits of comparative weakness analysis:

- Identifies which model performs better in specific regions
- Reveals complementary strengths across models
- Informs potential ensemble strategies
- Guides targeted model improvement efforts

# Comparing Categorical Feature Performance

```
1 # Compare models on categorical features
2 tsc = TestSuite(ds, models=[model_lgbm, model_xgb])
3 results = tsc.compare_slicing_accuracy(
4     features="EDUCATION",
5     metric="AUC",
6     threshold=0.6)
7 results.plot()
8
```

*This compares model performance across categorical feature values*

## What to look for:

- Categories where models show significant performance differences
- Segments where all models struggle (potential data issues)
- Categories with inconsistent performance across models
- Low-frequency categories with high performance variance



# Beyond Basic Weakness Detection

## Robustness Testing

- Tests model sensitivity to input noise
- Identifies regions vulnerable to small perturbations
- Evaluates stability of predictions

## Reliability Testing

- Focuses on prediction uncertainty
- Identifies regions with low confidence
- Evaluates calibration of probability estimates

## Resilience Testing

- Tests performance under distribution shift
- Evaluates behavior across heterogeneous data
- Measures degradation under changing conditions

## Integration with Other Tests

- Combine with fairness analysis
- Link to explainability assessments
- Connect with feature importance

*For details, refer to the corresponding sections in ModEva documentation*

# From Weakness Detection to Model Improvement

## Data Strategies

- Collect more data in weak regions
- Balance representation of underperforming segments
- Engineer new features to address specific weaknesses
- Apply targeted transformations to problematic features

## Model Strategies

- Adjust hyperparameters to focus on weak areas
- Create specialized models for challenging segments
- Implement higher capacity approaches such as mixture of experts

# From Weakness Detection to Model Deployment

## Deployment Strategies

- Implement guardrails for detecting edge cases
- Add uncertainty estimates to flag low-confidence predictions
- Create monitoring dashboards focused on weak regions
- Design fallback mechanisms for known weak spots

## Business Integration

- Communicate limitations to stakeholders
- Align model capabilities with business risk tolerance
- Design workflows that accommodate model weaknesses
- Prioritize improvements based on business impact

# Summary: Weakness Detection with MoDeVa

- 1 **Comprehensive Approach:** MoDeVa provides multiple methods for identifying model weaknesses across feature spaces
- 2 **Flexible Binning:** Choose from uniform, quantile, automatic, or custom binning to effectively slice data
- 3 **Feature Interactions:** Identify weaknesses in specific feature value combinations
- 4 **Distribution Analysis:** Understand the characteristics of weak regions compared to the overall dataset
- 5 **Model Comparison:** Compare weakness patterns across different models to guide improvement strategies
- 6 **Advanced Diagnostics:** Link weakness detection to robustness, reliability, and resilience testing

## Key Takeaway

Weakness detection transforms model development from a metrics-focused process to one that addresses specific underperforming regions.