# Feature Importance Analysis Report
## Performance Evaluation and Recommendations

Omer, Ahad

June 5, 2025

# 1 Executive Summary

This report provides a comprehensive evaluation of four feature importance techniques implemented for Task 3. Our analysis identifies critical methodological issues requiring immediate attention, while highlighting the reliability and effectiveness of each approach. The evaluation is based on mathematical rigor, statistical validity, and practical interpretability of results.

## 1.1 Key Findings

- **Linear Regression Coefficient Analysis** demonstrates optimal performance with mathematically sound implementation

- **Tree-Based Feature Importance** provides reliable baseline measurements suitable for production use

- **Correlation Analysis** requires enhancement due to high-dimensional data challenges

- **Permutation Importance Analysis** contains critical mathematical errors requiring immediate correction

# 2 Methodology Assessment

## 2.1 1. Tree-Based Feature Importance (XGBoost)

### 2.1.1 Performance Status: SATISFACTORY

**Implementation Quality:** The XGBoost feature importance implementation produces interpretable results. The methodology correctly leverages the model's internal feature splitting statistics to quantify variable importance.

**Technical Observations:**

- Feature importance scores are well-distributed across the feature space

- Clear hierarchical ranking facilitates feature selection decisions

- Computational efficiency suitable for large-scale applications

- No mathematical inconsistencies detected in the implementation

**Limitations:** The technique exhibits known bias toward high-cardinality features and may not generalize beyond tree-based model architectures.

## 2.2 2. Permutation Importance Analysis

### 2.2.1 Performance Status: CRITICAL ISSUES IDENTIFIED

**Mathematical Error Analysis:** Our evaluation has identified a fundamental computational error in the percentage calculation methodology, rendering current results unreliable for decision-making purposes.

**Specific Issues Detected:**

1. **Percentage Calculation Error:** Importance percentages sum to 882% rather than the mathematically required 100%

2. **Negative Value Handling:** The current implementation fails to properly process negative permutation importance scores

3. **Statistical Anomalies:** Top-ranked features exhibit importance values exceeding 200%, which is mathematically impossible

4. **Data Quality Concerns:** High importance scores for features with low statistical significance weights suggest potential overfitting

   **Root Cause:** The error originates from the percentage calculation formula:

$$\text{percentage}_i = \frac{\text{importance}_i}{\sum_{j=1}^{n} \text{importance}_j} \times 100 \tag{1}$$

When $\sum_{j=1}^{n} \text{importance}_j$ contains negative values or approaches zero, the calculation becomes unstable.

## 2.3 3. Linear Regression Coefficient Analysis

### 2.3.1 Performance Status: OPTIMAL PERFORMANCE

**Implementation Excellence:** This technique demonstrates exemplary implementation quality with proper feature standardization, mathematically consistent calculations, and interpretable results.
   **Technical Strengths:**

- Coefficient magnitudes within expected range (0.052 to 0.085)

- Balanced distribution of positive and negative coefficients

- Percentage calculations correctly sum to 100%

- Feature standardization ensures fair comparison across variables

   **Statistical Characteristics:**

- Individual feature contributions remain below 0.5%, indicating collective feature behavior

- Well-regularized model characteristics observed

- Linear relationship patterns clearly quantified

## 2.4 4. Feature-Target Correlation Analysis

### 2.4.1 Performance Status: REQUIRES ENHANCEMENT

**Technical Assessment:** While mathematically correct, the analysis reveals concerning patterns indicating high-dimensional data challenges that limit practical utility.
   **Observed Patterns:**

- Correlation coefficients range from 0.29 to 0.37 (moderate strength)

- Extremely diluted percentage distribution (0.014% to 0.017%)

- Evidence of high-dimensional feature space with distributed importance

**Implications:**

- Large feature space may require dimensionality reduction techniques

- Individual feature selection becomes challenging due to distributed importance

- Statistical significance testing recommended for robust interpretation

# 3  Comparative Performance Analysis

Table 1: Feature Importance Technique Performance Summary

| Technique | Reliability | Mathematical Rigor |
|---|---|---|
| Linear Regression | Excellent | Excellent |
| Tree-Based (XGBoost) | Good | Good |
| Correlation Analysis | Moderate | Good |
| Permutation Importance | Poor | Poor |

# 4  Critical Recommendations

## 4.1  Immediate Actions Required

1. **Priority 1 - Permutation Importance Correction:**

   - Implement robust handling of negative importance values
   - Correct percentage calculation methodology
   - Add statistical significance testing framework

2. **Priority 2 - Correlation Analysis Enhancement:**

   - Integrate p-value calculations for statistical significance
   - Implement correlation threshold filtering ($|\rho| > 0.1$)
   - Add multicollinearity detection capabilities

3. **Priority 3 - Data Quality Investigation:**

   - Examine features with high importance but low statistical significance
   - Implement overfitting detection mechanisms
   - Consider feature reduction strategies for high-dimensional data

## 4.2  Strategic Enhancements

- **Cross-Validation Framework:** Implement k-fold validation to assess feature importance stability

- **Ensemble Methodology:** Combine insights from multiple reliable techniques for robust feature selection

- **Statistical Rigor:** Establish significance thresholds and confidence intervals for all methods

- **Automated Quality Assurance:** Implement validation checks to prevent mathematical errors in future analyses

# 5 Conclusion

The feature importance analysis demonstrates mixed performance across implemented techniques. Linear Regression Coefficient Analysis exhibits optimal performance characteristics suitable for immediate deployment, while Tree-Based methods provide reliable baseline measurements. Critical mathematical errors in Permutation Importance analysis require immediate attention before the technique can be considered production-ready.

The high-dimensional nature of the dataset presents challenges that may benefit from dimensionality reduction strategies and enhanced statistical rigor across all methodologies.

**Prepared by:** Omer, Ahad
**Date:** June 6, 2025