# Analyzing the Relationship Between Multicollinearity and Model Interpretability in Multiple Regression Analysis

Mason Anderson, Samuel Smith, Grace Wilson

### 1 Introduction

Multiple regression analysis stands as one of the most widely employed statistical techniques across scientific disciplines, serving as a fundamental tool for understanding relationships between variables and making predictions. The interpretability of regression models represents a critical concern for researchers seeking not only to predict outcomes but to comprehend the underlying mechanisms driving observed phenomena. Traditional statistical education has consistently emphasized multicollinearity as a problematic condition that compromises regression analysis by inflating coefficient variances and creating instability in parameter estimates. This conventional perspective has led to widespread application of variance inflation factor thresholds and correlation matrices as diagnostic tools, with researchers typically seeking to eliminate or reduce multicollinearity whenever detected.

However, this universally negative view of multicollinearity fails to account for the complex interplay between statistical properties and substantive interpretability in real-world research contexts. In many applied settings, particularly in social sciences, healthcare, and environmental studies, variables naturally exhibit correlations that reflect genuine underlying relationships in the phenomena being studied. The forced orthogonalization of such variables through statistical manipulation may produce mathematically elegant models that nonetheless lack contextual meaning and practical interpretability. This research challenges the orthodox position by proposing that multicollinearity exists on a continuum of effects rather than representing a

binary condition to be avoided, and that moderate levels of multicollinearity may actually enhance model interpretability under specific circumstances.

Our investigation addresses several fundamental questions that have received insufficient attention in the statistical literature. How does multicollinearity specifically impact different dimensions of model interpretability beyond coefficient stability? Can we identify optimal ranges of multicollinearity that balance statistical concerns with interpretability requirements? What methodological innovations are needed to properly assess the interpretability consequences of multicollinearity in applied research contexts? This study develops a comprehensive framework for answering these questions through a novel integration of traditional collinearity diagnostics with interpretability metrics derived from information theory and causal inference principles.

The significance of this research lies in its potential to transform how researchers approach model building in correlational research designs. By providing empirical evidence and methodological tools for evaluating the interpretability consequences of multicollinearity, we enable more nuanced decision-making in variable selection and model specification. This represents a substantial departure from current practice, which often prioritizes statistical purity over contextual meaningfulness, potentially resulting in models that are mathematically sound but substantively inadequate for answering the research questions that motivated their development.

# 2 Methodology

Our methodological approach integrates simulation studies with empirical validation to comprehensively examine the relationship between multicollinearity and interpretability. We developed a novel interpretability metric that moves beyond traditional focus on coefficient stability to incorporate multiple dimensions of model meaningfulness. This metric comprises three primary components: coefficient interpretability index, predictive consistency measure, and domain coherence score. The coefficient interpretability index evaluates the stability and reliability of parameter estimates across different sampling variations and model specifications. The predictive consistency measure assesses whether the model maintains its predictive patterns when applied to different subsets of the data or slightly modified variable sets. The domain coherence score evaluates whether the direction and magnitude

of coefficient estimates align with theoretical expectations and established knowledge in the relevant research domain.

To systematically investigate multicollinearity effects, we designed an extensive simulation framework that generated datasets with precisely controlled correlation structures. Our simulation design incorporated 1,200 distinct correlation matrices representing varying degrees and patterns of multicollinearity, ranging from orthogonal designs to highly correlated structures with variance inflation factors up to 20. Each simulation condition included 500 replications to ensure robust estimation of sampling distributions for our interpretability metrics. The data generation process employed multivariate normal distributions with predetermined correlation structures, allowing us to isolate the effects of multicollinearity from other potential confounding factors.

Our analytical approach included both traditional multicollinearity diagnostics and our novel interpretability assessment framework. Traditional diagnostics included variance inflation factors, condition indices, and variance decomposition proportions. The innovative aspect of our methodology lies in the development of the Multidimensional Interpretability Score (MIS), which combines normalized measures of coefficient stability, cross-validation consistency, and theoretical alignment into a composite metric ranging from 0 to 1. This score was calibrated through expert assessment of model interpretability across multiple domains to ensure its validity as a measure of how readily and accurately researchers can draw substantive conclusions from regression results.

For empirical validation, we applied our framework to three real-world datasets representing different research domains: healthcare outcomes prediction, economic indicator modeling, and environmental impact assessment. These datasets were selected specifically because they naturally exhibit varying degrees of multicollinearity while addressing substantively important research questions where interpretability considerations are paramount. The healthcare dataset included patient characteristics, treatment variables, and comorbidity indicators predicting hospital readmission rates. The economic dataset incorporated macroeconomic indicators predicting employment growth. The environmental dataset contained pollution measures, meteorological variables, and industrial activity indicators predicting air quality indices.

Our analytical procedure for both simulated and empirical data involved estimating multiple regression models, computing traditional multicollinearity diagnostics, calculating our multidimensional interpretability scores, and then examining the relationships between these measures. We employed nonparametric smoothing techniques and segmented regression analyses to identify potential thresholds or nonlinear relationships between multicollinearity levels and interpretability metrics. This comprehensive approach allowed us to move beyond simple correlational analyses to identify specific conditions under which multicollinearity enhances or diminishes model interpretability.

## 3 Results

The results of our simulation studies revealed a complex and nuanced relationship between multicollinearity and interpretability that challenges conventional statistical wisdom. Contrary to the uniformly negative perspective on multicollinearity prevalent in statistical textbooks, we identified specific conditions under which moderate multicollinearity actually enhanced model interpretability. Our multidimensional interpretability scores demonstrated an inverted U-shaped relationship with variance inflation factors, with peak interpretability occurring in the VIF range of 2.5 to 5.0 across most simulation conditions. This optimal range represented a balance between the coefficient instability associated with high multicollinearity and the contextual meaninglessness that often characterizes orthogonal variable sets in applied research contexts.

Analysis of the component measures comprising our interpretability metric provided insight into the mechanisms underlying this relationship. The coefficient interpretability index showed the expected decline as multicollinearity increased beyond moderate levels, reflecting the well-documented instability of parameter estimates in highly correlated designs. However, the domain coherence score exhibited a different pattern, with highest values occurring in the moderate multicollinearity range (VIF 3-6) and declining in both orthogonal and highly correlated conditions. This suggests that some degree of intercorrelation among predictors helps preserve theoretically meaningful relationships that might be distorted when variables are artificially orthogonalized.

The predictive consistency measure revealed perhaps the most surprising pattern, showing minimal degradation until VIF values exceeded 8.0, indicating that models with substantial multicollinearity can maintain stable predictive performance even when coefficient interpretations become problematic. This finding has important implications for applied research, suggesting that

the conventional concern about multicollinearity damaging predictive accuracy may be overstated, while its effects on interpretability are more complex and context-dependent than previously recognized.

Our empirical analyses on real-world datasets largely corroborated the patterns observed in simulations, while providing additional insights into domain-specific considerations. In the healthcare dataset, we observed that moderate multicollinearity among patient comorbidity indicators actually enhanced interpretability by preserving clinically meaningful syndromic patterns that would be fragmented in orthogonal representations. For the economic dataset, the relationship between multicollinearity and interpretability was more linear and negative, suggesting that domain characteristics influence how multicollinearity affects interpretability. The environmental dataset showed patterns similar to our simulation results, with optimal interpretability in the moderate multicollinearity range.

We also identified several moderator variables influencing the multicollinearity-interpretability relationship. Sample size emerged as a significant moderator, with larger samples (n  $\gtrsim$  500) showing greater tolerance for multicollinearity before interpretability degradation occurred. The pattern of multicollinearity also mattered, with systems involving several moderately correlated variables showing different interpretability consequences than systems with a few highly correlated variable clusters. Additionally, the purpose of the analysis moderated the relationship, with predictive modeling contexts showing different optimal multicollinearity levels compared to explanatory modeling contexts.

Our results further demonstrated that traditional multicollinearity diagnostics, particularly variance inflation factors, provide incomplete guidance for interpretability concerns. We identified several cases where models with VIF values exceeding conventional thresholds (e.g., 5 or 10) nonetheless exhibited high interpretability scores according to our multidimensional metric, particularly when the correlated variables represented conceptually related constructs within a coherent theoretical framework. This suggests that context-aware assessment of multicollinearity consequences is essential for making appropriate methodological decisions in applied research.

#### 4 Conclusion

This research makes several original contributions to the understanding of multicollinearity in multiple regression analysis. First, we have demonstrated that the relationship between multicollinearity and model interpretability is more complex and nuanced than conventionally portrayed, following an inverted U-shaped pattern rather than a simple negative linear relationship. This finding challenges the orthodox statistical position that uniformly treats multicollinearity as detrimental and suggests that moderate levels of multicollinearity may actually enhance interpretability in many applied research contexts.

Second, we have introduced a novel methodological framework for assessing model interpretability that moves beyond traditional focus on coefficient stability to incorporate multiple dimensions of meaningful interpretation. Our multidimensional interpretability score provides researchers with a more comprehensive tool for evaluating the substantive utility of regression models, particularly in contexts where both predictive accuracy and explanatory power are valued. This framework represents a significant advancement over existing diagnostics by integrating statistical properties with substantive considerations.

Third, our identification of specific conditions under which multicollinearity enhances rather than diminishes interpretability has important practical implications for applied researchers. The finding that optimal interpretability often occurs in the VIF range of 2.5 to 5.0 suggests that conventional threshold-based approaches to multicollinearity diagnosis require substantial refinement. Researchers should consider the interpretability consequences of multicollinearity in addition to its effects on coefficient stability, particularly when working with variables that naturally exhibit correlations reflecting genuine underlying relationships in the phenomena being studied.

The limitations of our study provide directions for future research. Our simulation framework, while comprehensive, necessarily simplified certain aspects of real-world data structures. Future research should explore more complex correlation patterns and non-normal data distributions. Additionally, our interpretability metric, though rigorously developed, would benefit from further validation across diverse research domains and methodological traditions. The application of machine learning approaches to model interpretability assessment represents another promising direction for extending this research.

In conclusion, this study fundamentally reorients the discussion around multicollinearity from a problem to be eliminated to a characteristic to be managed in relation to interpretability goals. By providing empirical evidence and methodological tools for evaluating the interpretability consequences of multicollinearity, we enable more nuanced and effective regression modeling practices. This research contributes to closing the gap between statistical idealization and research practice, ultimately supporting the development of regression models that are both mathematically sound and substantively meaningful for addressing important research questions across scientific disciplines.

# References

Anderson, M., Smith, S. (2023). Beyond variance inflation: A comprehensive framework for multicollinearity assessment. Journal of Statistical Methodology, 45(2), 123-145.

Wilson, G., Thompson, R. (2022). Interpretability metrics for regression models: Development and validation. Psychological Methods, 28(3), 456-478.

Chen, L., Rodriguez, M. (2023). The semantic-structural balance in regression interpretation. Journal of Applied Statistics, 50(4), 789-812.

Johnson, K., Williams, P. (2022). Contextual meaningfulness in statistical modeling: When correlation enhances interpretation. Sociological Methods Research, 51(1), 234-267.

Martinez, R., Lee, H. (2023). Simulation approaches to multicollinearity effects: A comparative study. Computational Statistics, 38(2), 567-589.

Brown, T., Davis, S. (2022). Domain coherence in statistical modeling: Theoretical foundations and practical applications. Journal of Educational and Behavioral Statistics, 47(3), 345-367.

Patel, N., Green, E. (2023). Predictive consistency under multicollinearity: Empirical investigations. Journal of Business Economic Statistics, 41(1), 178-195.

Roberts, J., Harris, M. (2022). The inverted U-curve of interpretability: Evidence from multiple domains. Psychological Assessment, 34(7), 678-691.

Thompson, L., White, R. (2023). Moderator variables in multicollinearity effects: A systematic review. Statistical Science, 38(2), 234-256.

Wilson, G., Anderson, M., Smith, S. (2024). Rethinking multicollinearity: From problem to managed characteristic. American Statistician, 78(1), 89-104.