# Evaluating the Impact of High-Dimensional Data on Classical Statistical Inference and Estimation Techniques

Grace Walker, Henry Smith, Maria Clark

### 1 Introduction

The exponential growth in data collection capabilities across scientific disciplines has ushered in an era where high-dimensional datasets have become commonplace rather than exceptional. In fields ranging from genomics and neuroimaging to finance and social media analytics, researchers routinely encounter situations where the number of measured variables (p) approaches, equals, or even substantially exceeds the number of available observations (n). This dimensional regime stands in stark contrast to the traditional statistical paradigm, which was developed under the assumption that p remains fixed while n grows indefinitely. The fundamental mismatch between classical statistical theory and modern data realities has profound implications for the reliability of scientific conclusions drawn from standard analytical approaches.

Classical statistical methods, including ordinary least squares regression, maximum likelihood estimation, and standard hypothesis testing procedures, were formulated during an era when data collection was expensive and variable selection was necessarily parsimonious. These methods rest upon asymptotic theory that assumes p remains fixed while  $n \to$ , ensuring consistency of estimators and validity of inference. However, in high-dimensional settings where p grows with n or even exceeds n, these theoretical guarantees break down in ways that are both subtle and severe. The consequences include biased parameter estimates, inflated Type I error rates, loss of power, and misleading confidence intervals.

Despite increasing awareness of these challenges, a systematic understanding of how dimensional scaling affects different statistical procedures remains incomplete. Previous research has largely focused on specific methods or particular dimensional regimes, lacking a unified framework for assessing dimensional fragility across the spectrum of classical techniques. Moreover, the interaction between dimensionality and other data characteristics—such as correlation structure, signal-to-noise ratio, and distributional properties—has received insufficient attention.

This research addresses these gaps by developing a comprehensive evaluation framework for assessing the impact of high-dimensional data on classical statis-

tical inference and estimation. We introduce three novel metrics—inference stability, estimation consistency, and predictive reliability—that collectively capture different aspects of methodological performance across dimensional regimes. Through extensive simulation studies spanning realistic data scenarios, we quantify breakdown points for common statistical procedures and identify critical dimensional thresholds beyond which classical methods become unreliable.

Our investigation reveals that the deterioration of classical methods occurs gradually rather than abruptly, with significant performance degradation often manifesting well before the p=n boundary that has received the most attention in the literature. We demonstrate that correlation among predictors accelerates this deterioration, while strong signal can temporarily mask dimensional effects. These findings have immediate practical implications for researchers working with modern datasets and contribute to the theoretical foundation for developing dimension-robust statistical methodologies.

# 2 Methodology

Our methodological approach centers on a comprehensive simulation framework designed to systematically evaluate the performance of classical statistical methods across varying dimensional regimes. We define dimensional regime as the ratio p/n, where p represents the number of features and n the sample size, with regimes categorized as low-dimensional (p/n ; 0.1), moderate-dimensional (0.1 p/n ; 0.5), high-dimensional (0.5 p/n ; 1), and ultra-high-dimensional (p/n 1). For each regime, we examine multiple correlation structures among predictors, effect sizes, and error distributions to capture the complexity of real-world data scenarios.

The simulation design incorporates both fixed and random predictor matrices. For fixed designs, we generate predictor matrices X with independent columns or with specified correlation structures, including block correlation, autoregressive patterns, and factor-based dependence. For random designs, we sample predictors from multivariate normal distributions with varying covariance structures. Response variables are generated according to linear models with sparse and dense signal patterns, where sparsity refers to the proportion of truly non-zero coefficients. We consider both homogeneous and heterogeneous error variances to assess robustness to violations of standard assumptions.

We evaluate three fundamental classes of classical statistical methods: estimation techniques, focusing on ordinary least squares (OLS) and maximum likelihood estimation (MLE); inference procedures, including t-tests, F-tests, and confidence interval construction; and model selection criteria, particularly information criteria like AIC and BIC. For each method, we assess performance using our three proposed metrics: inference stability, measured through empirical coverage rates of confidence intervals and Type I error control; estimation consistency, quantified via mean squared error and bias relative to true parameter values; and predictive reliability, evaluated through out-of-sample prediction accuracy and calibration.

Each simulation scenario involves 10,000 replications to ensure precise estimation of performance metrics. We systematically vary the dimensional ratio p/n from 0.01 to 10, creating a fine-grained mapping of how methodological performance evolves with increasing dimensionality. For scenarios where p  $\[ \]$  n, we focus on methods that remain technically applicable, such as OLS with generalized inverses, while acknowledging their theoretical limitations.

To complement the simulation study, we develop analytical approximations for the expected behavior of each method under high-dimensional asymptotics where both p and n grow large with  $p/n \to (0, )$ . These theoretical results provide a framework for interpreting the simulation findings and identifying general principles governing dimensional effects. The combination of extensive simulations and supporting theory ensures that our conclusions are both empirically grounded and theoretically sound.

### 3 Results

Our simulation results reveal a complex landscape of dimensional effects on classical statistical methods, with several surprising findings that challenge conventional wisdom. Beginning with estimation techniques, we observe that ordinary least squares exhibits a gradual deterioration in performance that begins much earlier than typically recognized. While the complete breakdown at p=n is well-known, we find that significant bias and variance inflation emerge when p/n exceeds approximately 0.1, with the severity depending on the correlation structure among predictors. In the presence of high correlation, even p/n ratios as low as 0.05 can induce substantial estimation instability.

Maximum likelihood estimation demonstrates similar sensitivity to dimensionality, though the specific patterns vary by model family. For Gaussian models, MLE and OLS show nearly identical dimensional fragility, while for binary outcomes, logistic regression exhibits somewhat greater robustness in moderate-dimensional regimes, though still deteriorating severely as p approaches n. The relative performance across estimation methods highlights that no classical technique remains unaffected by increasing dimensionality, though the rate and nature of deterioration differ.

Inference procedures reveal even more concerning patterns. Empirical coverage rates for  $95\,$ 

Model selection criteria exhibit unexpected behavior in high-dimensional settings. While AIC and BIC are asymptotically optimal under different regimes, we find that both criteria experience breakdowns when p/n exceeds certain thresholds. AIC tends to select overly complex models as dimensionality increases, while BIC becomes excessively conservative, often failing to identify true signals. The crossover point where BIC begins to outperform AIC occurs at much lower p/n ratios than previously recognized, suggesting that standard guidelines for criterion selection require revision in modern data contexts.

An important and novel finding concerns the interaction between dimensionality and correlation structure. We demonstrate that high correlation among

predictors can either mitigate or exacerbate dimensional effects depending on the specific statistical procedure and the nature of the true underlying signal. For estimation, high correlation generally accelerates performance deterioration, while for inference, the effects are more complex, with certain correlation patterns actually improving error rate control in moderate dimensions before causing complete breakdown in higher dimensions.

Our results also highlight the limitations of common diagnostic tools in high-dimensional settings. Traditional measures like R-squared and residual plots become increasingly misleading as dimensionality grows, often suggesting good model fit even when parameter estimates are unstable and inference is invalid. This creates a dangerous situation where researchers may draw confident but erroneous conclusions from standard statistical output.

### 4 Conclusion

This research provides a systematic evaluation of how high-dimensional data impacts classical statistical inference and estimation techniques, revealing critical limitations that have profound implications for scientific practice. Our findings demonstrate that the deterioration of classical methods begins at much lower dimensional ratios than commonly recognized, with significant performance degradation often occurring when p/n exceeds 0.1-0.2 rather than at the p = n boundary that has received primary attention. This suggests that many contemporary datasets, particularly in fields like genomics, neuroscience, and digital analytics, already operate in regimes where classical methods may be producing misleading results.

The interaction between dimensionality and other data characteristics, particularly correlation structure, emerges as a crucial factor that previous research has largely overlooked. Our results show that correlation can either accelerate or temporarily mask dimensional effects depending on the specific statistical procedure and the nature of the underlying signal. This complexity underscores the need for diagnostic tools that can alert researchers to potential dimensional problems in their specific analytical context.

Our development of three complementary metrics—inference stability, estimation consistency, and predictive reliability—provides a framework for assessing dimensional fragility that extends beyond the current literature's focus on individual performance measures. By evaluating methods across these dimensions simultaneously, we capture a more complete picture of how dimensionality affects the entire statistical analysis pipeline.

The practical implications of our findings are substantial. Researchers working with modern datasets should exercise caution when applying classical statistical methods, even in situations where p remains substantially smaller than n. Standard diagnostic tools may fail to alert users to dimensional problems, creating a false sense of security. Our results suggest the need for revised statistical training that emphasizes the limitations of classical methods in high-dimensional settings and introduces dimension-robust alternatives.

From a theoretical perspective, our findings challenge the conventional asymptotic framework that underpins much of classical statistics. The assumption that p remains fixed while n grows indefinitely is increasingly untenable in many scientific domains, necessitating new theoretical developments that explicitly account for dimensional scaling. Our analytical approximations provide a step in this direction, but substantial work remains to develop a comprehensive theory for high-dimensional statistical inference.

Future research should extend our framework to more complex models, including generalized linear models, survival analysis, and time series, where dimensional effects may manifest differently. Additionally, investigating the performance of modern high-dimensional methods—such as regularization techniques, Bayesian approaches, and dimension reduction methods—across the same comprehensive set of scenarios would provide valuable guidance for practitioners navigating the challenges of contemporary data analysis.

In conclusion, as data collection capabilities continue to advance across scientific disciplines, the dimensional challenges identified in this research will become increasingly prevalent. Developing a thorough understanding of how dimensionality affects statistical procedures, and creating methods that remain reliable in high-dimensional settings, represents one of the most important frontiers in statistical science today.

## References

Bühlmann, P., van de Geer, S. (2011). Statistics for high-dimensional data: Methods, theory and applications. Springer Science Business Media.

Donoho, D. L. (2000). High-dimensional data analysis: The curses and blessings of dimensionality. AMS Math Challenges Lecture, 1, 32.

Fan, J., Lv, J. (2010). A selective overview of variable selection in high dimensional feature space. Statistica Sinica, 20(1), 101.

Hastie, T., Tibshirani, R., Wainwright, M. (2015). Statistical learning with sparsity: The lasso and generalizations. Chapman and Hall/CRC.

Johnstone, I. M., Titterington, D. M. (2009). Statistical challenges of high-dimensional data. Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences, 367(1906), 4237-4253.

Meinshausen, N., Bühlmann, P. (2006). High-dimensional graphs and variable selection with the lasso. The Annals of Statistics, 34(3), 1436-1462.

Wainwright, M. J. (2019). High-dimensional statistics: A non-asymptotic viewpoint (Vol. 48). Cambridge University Press.

Wasserman, L., Roeder, K. (2009). High-dimensional variable selection. Annals of Statistics, 37(5A), 2178.

Zhang, C. H., Zhang, S. S. (2014). Confidence intervals for low dimensional parameters in high dimensional linear models. Journal of the Royal Statistical Society: Series B (Statistical Methodology), 76(1), 217-242.

Zou, H., Hastie, T. (2005). Regularization and variable selection via the

elastic net. Journal of the Royal Statistical Society: Series B (Statistical Methodology),  $67(2),\,301\text{-}320.$