

Summary of *Missing values handling for machine learning portfolios*

Andrew Y. Chen, Jack McCoy(JFE,2022.10)

2024.10.16 Citattoo

1. What are the research questions?

- How do missing values in stock return predictors affect machine learning portfolio performance?
- Which imputation techniques perform best for handling missing data in cross-sectional predictors?
- Can simple imputation methods rival more complex algorithms in capturing expected returns?

2. Why are the research questions interesting?

- This research helps determine whether sophisticated imputation methods add value or simply introduce noise, affecting portfolio performance.
- The findings can shape how large datasets with missing predictors are treated in asset pricing and other finance-related machine learning models.

3. What is the paper's contribution?

- The paper systematically compares different imputation methods for 159 stock return predictors, demonstrating that simple mean imputation is competitive with more complex methods like expectation-maximization (EM).
- It highlights that complex imputations may introduce estimation noise, particularly when data structures (like large blocks of missing data) limit the utility of advanced techniques.
- The findings challenge the common belief that sophisticated imputations are superior, suggesting practical applications for simple techniques in machine learning portfolios.

4. What hypotheses are tested in the paper?

- H1: Simple cross-sectional mean imputation performs similarly to more complex methods like EM for stock return predictors.
- H2: Sophisticated imputation methods may introduce more noise than value in cases with high data missingness.

a) Do these hypotheses follow from and answer the research questions?

- Yes.

b) Do these hypotheses follow from theory? Explain logic of the hypotheses.

- The logic behind both hypotheses aligns with theories in machine learning and finance, emphasizing that over-complicating models with sophisticated techniques does not always yield better results, especially when the underlying data lacks the structure needed to support these methods

5. Sample: comment on the appropriateness of the sample selection procedures.

- The sample includes 159 cross-sectional stock return predictors, covering a wide array of financial data. This ensures a comprehensive analysis of imputation techniques across different types of missingness patterns.
- Limitations include focusing mainly on predictors from the CRSP and Compustat datasets, potentially missing other relevant data sources.

6. Comment on the appropriateness of variable definition and measurement.

- The paper carefully defines missing values across predictors, focusing on both time-series and cross-sectional missingness. It employs a variety of imputation techniques and validates their impact on portfolio returns.
- Measurement techniques like principal component regressions and machine learning models are well-suited for the study's purpose.

7. Comment on the appropriateness of the regress/predict model specification.

- The model specifications (e.g., principal component regression, neural networks) are appropriate for evaluating the impact of different imputation methods on portfolio performance.
- However, the assumption of linear relationships may oversimplify the interactions between predictors and returns.

8. What difficulties arise in drawing inferences from the empirical work?

- The lack of observable relationships between predictors (most cross-sectional correlations are near zero) makes it hard to improve performance with advanced imputation methods.
- Imputation errors can be particularly large for small-cap stocks, which complicates portfolio performance analysis.

9. Describe at least one publishable and feasible extension of this research.

- Extend the analysis to different markets, such as emerging markets, where data quality issues are more prevalent.
- Explore alternative imputation techniques that blend time-series and cross-sectional information, providing a more nuanced approach to handling missing values.