

From Data Imputation to Data Cleaning - Automated Cleaning of Tabular Data Improves Downstream Predictive Performance

DISCLAIMER: Summarized by AI

Problem they are trying to solve / Purpose of method

The paper addresses the issue of **data quality in tabular datasets**, focusing specifically on the impact of **automated data cleaning** on downstream **predictive model performance**.

Previous problems that needed solving:

- **Real-world tabular data is often noisy**, with missing values, invalid entries, or inconsistencies that hinder predictive modeling.
- Existing approaches often focus on **data imputation** or specific cleaning rules but fail to generalize or optimize cleaning **with respect to downstream tasks** (like classification or regression).
- There's a **lack of benchmarks** and systematic evaluations to measure how cleaning impacts predictive performance.

Why the method is introduced:

- To **automatically clean tabular data** in a way that directly improves the performance of predictive models.
- To **go beyond traditional imputation** by addressing a broader range of errors and incorporating downstream task performance into the cleaning process.
- To provide a **general, extensible, and modular framework** that supports various types of cleaning and integrates with machine learning pipelines.

How does it differ from other methods?

Differences from other methods:

- Most previous methods **focus solely on imputing missing values**, whereas this approach addresses a **wider range of data errors**, including type errors, inconsistencies, and outliers.
- Traditional cleaning does not take into account **the effect of cleaning on the target machine learning task**, while this method directly **optimizes data cleaning for better predictive performance**.

Unique aspects:

- Introduces **AutoClean**, an **automated system** that selects, configures, and applies cleaning operations based on their **measured impact on downstream model accuracy**.
- Uses **search and optimization techniques** to choose the most effective cleaning operations.
- Evaluated on **64 real-world datasets**, showing that their approach significantly improves downstream performance over default cleaning strategies and state-of-the-art imputation tools.

How the method works

Simple overview:

- **AutoClean** automatically applies a pipeline of data cleaning operations, such as imputation, outlier removal, and error correction.
- It **evaluates combinations of cleaning operations** by measuring how much they improve predictive performance on a validation set.
- The system **searches for the best sequence and configuration** of operations using techniques like greedy search and ensembling.

More detailed breakdown:

1. **Modular Cleaning Primitives:**
 - The system includes multiple types of cleaning operations (e.g., imputation, value correction, outlier removal).
 - Each operation is implemented as a modular “primitive” with various parameter settings.
2. **Cleaning Policy Search:**
 - Uses a **greedy algorithm** to evaluate different sequences of cleaning operations.
 - Assesses each candidate pipeline by running a predictive model on cleaned data and measuring validation performance.
3. **Meta-Ensemble Approach:**
 - Combines multiple top-performing cleaning policies using **model ensembling** to further improve performance.
4. **Comprehensive Evaluation:**
 - Benchmarks across 64 real-world datasets.
 - Demonstrates improvements over baseline imputation and cleaning strategies.
 - Released an **open-source library and dataset suite** for further research.