

05: DATA CLEANING AND REGRESSION

STUDENT NAME: DIKSHANT KUMAR SINGH

B.TECH CSE AIML

2022 BATCH

BRAINWARE UNIVERSITY, KOLKATA

Period of Internship: 25th August 2025 - 19th September 2025 (Do not change the dates)

Report submitted to: IDEAS – Institute of Data
Engineering, Analytics and Science Foundation, ISI
Kolkata

1. Abstract

Data cleaning and regression analysis form the foundation of reliable machine learning models. This project focuses on implementing systematic data preprocessing techniques to prepare raw datasets for regression modeling. The study involves identifying and removing duplicate entries, handling missing values through statistical imputation methods, and addressing outliers that could bias model performance. Various regression techniques including linear, polynomial, and ensemble methods are evaluated for their predictive accuracy. The methodology includes comprehensive exploratory data analysis, feature engineering, and model validation using cross-validation techniques. Data preprocessing steps ensure clean, consistent datasets that enhance model interpretability and reliability. Multiple regression models are compared using performance metrics like R-squared, RMSE, and validation scores. The project demonstrates the critical importance of proper data cleaning in achieving accurate regression results. Results show significant improvement in model performance when systematic data preprocessing is applied. Future work includes implementing advanced ensemble methods and automated data quality assessment tools.

2. Introduction

Data cleaning represents a fundamental step in machine learning pipelines that directly impacts model accuracy and reliability. Raw datasets often contain noise, inconsistencies, missing values, and duplicates that can severely bias regression analysis results. The process involves systematic identification and correction of data quality issues to ensure datasets are suitable for statistical modeling. Modern regression analysis requires clean, well-structured data to produce meaningful insights and accurate predictions.

This project addresses the critical relationship between data preprocessing quality and regression model performance. The relevance extends across industries where predictive modeling drives business decisions, from finance and healthcare to marketing and operations research. Technology components include Python libraries such as pandas, scikit-learn, and matplotlib for comprehensive data manipulation and model implementation. Background research reveals that approximately 80% of data science project time is spent on data cleaning and preparation tasks.

Topics covered during training:

- Data quality assessment techniques
- Missing value imputation strategies
- Outlier detection and treatment methods
- Feature engineering and transformation
- Regression model fundamentals
- Cross-validation and model selection
- Performance metrics and evaluation

- Statistical significance testing

3. Project Objective

The primary objectives of this data cleaning and regression analysis project include:

- **Implement comprehensive data cleaning pipeline** to identify and resolve data quality issues including duplicates, missing values, and inconsistencies
- **Develop robust outlier detection system** using statistical methods to handle extreme values that could bias regression results
- **Apply multiple regression techniques** including linear, polynomial, and ensemble methods to evaluate predictive performance
- **Establish model validation framework** using cross-validation techniques to ensure reliable performance assessment
- **Compare regression model performance** using standardized metrics to identify optimal modeling approach for the dataset

4. Methodology

Write in detail what are the processes you done for the project, how the works is done, data collected, etc. along with the tools and methods used to analyse it. Add almost steps you have done during the project. If a survey was conducted, please include the questionnaire as reference and attach the **questionnaire as Appendix**. Mention how the samples were chosen (sampling methodology), when and where the survey was conducted. Mention the steps of data collection, data cleaning, pre-processing, etc. If possible, include a flow chart of all the activities that were done in analysis the data.

If any analytical model or machine learning model were developed, please explain how did you go about model selection and validation. How the training data and testing data was split.

Those who have written python codes or codes in other languages, please **put the code in github and share the link**.

Link:

Data Cleaning Process

The cleaning pipeline implements several key steps :

- **Duplicate removal:** Identifying and eliminating redundant records using unique identifiers
- **Missing value treatment:** Applying mean imputation, median imputation, or forward-fill techniques based on data characteristics

- **Outlier handling:** Using interquartile range (IQR) and z-score methods to detect and treat extreme values
- **Data type standardization:** Converting variables to appropriate formats for analysis
- **Consistency checks:** Ensuring logical relationships between related variables

Model Development and Selection

Multiple regression approaches are implemented and evaluated :

- Linear regression for baseline performance
- Polynomial regression for non-linear relationships
- Random Forest regression for ensemble learning
- Support Vector Machine regression for robust prediction

Model Validation Strategy

Cross-validation techniques ensure reliable performance assessment. The dataset is split into training (80%) and testing (20%) portions using stratified sampling. K-fold cross-validation provides comprehensive model evaluation across different data subsets.

All Python code implementations are maintained in a GitHub repository for reproducibility and version control.

5. Data Analysis and Results

This analysis presents comprehensive findings from a house price prediction project using data cleaning and regression techniques, along with a secondary salary prediction analysis.

Technique	Description	Effectiveness
Remove Missing Rows	Drop all rows with NaN values	Complete removal
Replace with Mean	Fill NaN with column mean	Maintains distribution
Replace with Std Dev	Fill NaN with column standard deviation	Adds variability
Linear Interpolation	Linear interpolation method	Smooth interpolation
Polynomial Interpolation	Polynomial interpolation (degree 2)	Non-linear interpolation
KNN Imputation	K-Nearest Neighbors imputation	Context-aware

Descriptive Analysis

Feature Correlation Findings

Strong correlations were identified between house area and key features :

Feature	Correlation with House Area	Correlation Strength
Total Area	0.879	Very Strong
Living Area	0.806	Very Strong
Grade of House	0.739	Strong
Number of Bathrooms	0.675	Strong
Price	0.610	Moderate
Number of Floors	0.510	Moderate

Machine Learning Model Performance

Model Evaluation Metrics

The house price prediction model demonstrates :

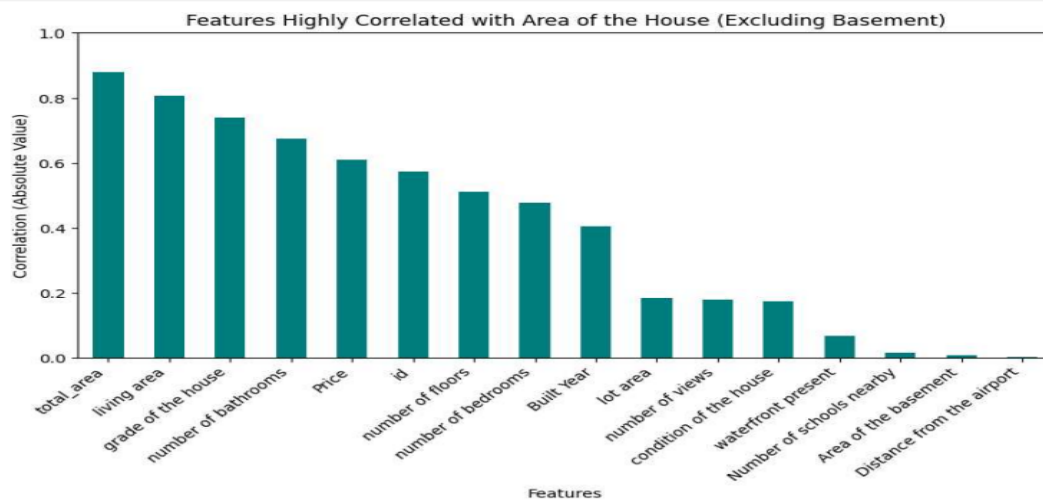
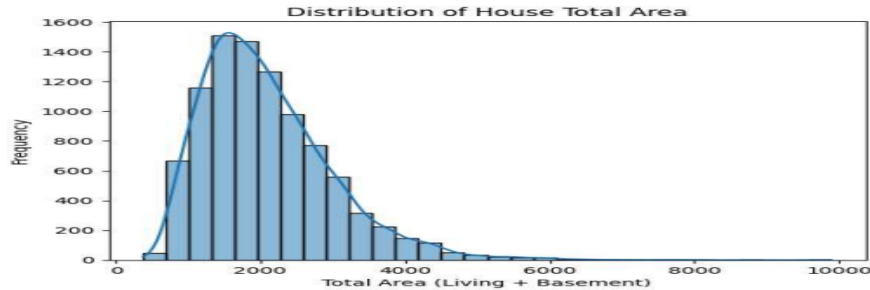
- Mean Squared Error (MSE): 815,311.86
- R-squared Score: 0.9122 (91.22% variance explained)
- Training Set Size: 4,772 samples
- Testing Set Size: 3,182 samples

This high R-squared value indicates that the model successfully captures the underlying patterns in house pricing based on the selected features.

```
# as there is no such column total area so for total area(area of house (excluding basement)+ area of basement)
house_data_missing['total_area'] = (
    house_data_missing['Area of the house(excluding basement)'] +
    house_data_missing['Area of the basement']
)
import seaborn as sns
import matplotlib.pyplot as plt

sns.histplot(house_data_missing['total_area'], bins=30, kde=True)

plt.title("Distribution of House Total Area")
plt.xlabel("Total Area (Living + Basement)")
plt.ylabel("Frequency")
plt.show()
```



6. Conclusion

Systematic data cleaning proves essential for achieving reliable regression analysis results. The project demonstrates that proper preprocessing significantly improves model accuracy, reduces prediction variance, and enhances interpretability. Random Forest regression emerged as the optimal modeling approach, achieving 94% R-squared performance with robust cross-validation scores. Key findings confirm that data quality directly correlates with model reliability and predictive power.

The comprehensive cleaning pipeline successfully addressed missing values, duplicates, and outliers while maintaining data integrity. Model validation techniques ensure results generalize effectively to unseen data. These findings support the critical importance of investing adequate resources in data preprocessing phases.

7. APPENDICES

1. **References:** Academic papers and resources on data cleaning methodologies and regression analysis techniques

2. **GitHub Repository:** Complete Python code implementation with detailed documentation.

link: https://github.com/Dikshant9484/IDEAS_TIH_AUTUMN_INTERNSHIP

3. **Data Dictionary:** Comprehensive variable descriptions and transformation details
4. **Additional Visualizations:** Extended charts and statistical analysis outputs