# Analytical Techniques in Data Science: Simulation and Regression Analysis for Model Evaluation

02405885

Compiled: April 20, 2024

**Github Repo:** https://github.com/taipingxian/MATH70076-Data-Science-Assessment-2

# 1 Project Description

The project begins by simulating a dataset of 100 points using a simple linear model $y=a+bx+$error, where $a=5$, $b=7$, $x$ values are randomly sampled from a uniform distribution over the range [0, 50], and the error term follows a normal distribution with mean 0 and standard deviation 3. The generated data is then used to fit a linear regression model, and the coefficients of this model are determined. The analysis includes plotting the data along with the fitted regression line, and also annotating the plot with the regression line's equation.

Subsequently, the project introduces a more complex scenario by simulating data from a quadratic model $y=a+bx+cx2+$error, which introduces nonlinearity into the dataset. Here, the coefficients $a$, $b$, and $c$ are chosen to ensure the nonlinearity is apparent in a scatterplot. Despite the underlying model being quadratic, a linear model is fit to this data to explore the effects of model mis-specification. The output and implications of fitting a linear model to nonlinear data are analyzed and displayed graphically, highlighting the concept of "best-fit" in terms of minimizing the residual sum of squares.

Through these exercises, students learn to handle and analyze data in R, understand the importance of model selection, and evaluate the fit of different statistical models. This educational experience provides practical skills in data simulation, regression analysis, and critical evaluation of model fits, crucial for aspiring data scientists.

# 2 Assessment Criteria

*In 2-3 bullet points each and at most 1 page in total, describe how your submission addresses each of the assessment criteria below. You may delete this italicised text when filling in the template.*

**Technical Competence:** Proficiency in data collection, processing, analysis, and coding.

- Data Handling Proficiency: Demonstrated expertise in generating synthetic datasets through random sampling and manipulation in R, ensuring a robust foundation for subsequent analyses.

- Advanced Coding Skills: Utilized advanced R functionalities, including packages like ggplot2 for visualization and rstanarm for Bayesian regression, showcasing a high level of coding proficiency and technical capability.

**User Interface:** Design, functionality, and usability of the final data product.

- Interactive Design: Developed a user-friendly graphical interface in R using Shiny, allowing users to interactively adjust model parameters and instantly visualize the effects on regression outcomes.

- Functional and Aesthetic Elements: Ensured the interface is not only functional but also aesthetically pleasing, improving user engagement through well-organized layouts and clear, intuitive controls.

**Analysis and Interpretation:** Depth of analysis, appropriate use of statistical methods, and meaningful interpretation.

- Depth of Statistical Analysis: Applied both simple and complex regression models to the data, providing a comprehensive comparison of model fits and discussing the implications of each in the context of data characteristics.

- Insightful Interpretation of Results: Offered detailed explanations of statistical outputs, interpreting coefficients and fit statistics to draw meaningful conclusions about the underlying data relationships.

**Presentation and Communication:** Clarity, organisation and effectiveness of written and visual communication.

- Effective Documentation: Prepared a clear and well-organized written report, using structured sections and sub-sections, bullet points for clarity, and appropriate graphs to visually support textual content.

- Communication of Complex Concepts: Simplified complex statistical concepts through effective communication strategies, ensuring that the analysis is accessible to audiences with varying levels of statistical knowledge.

**Reproducibility:** Clarity and completeness of documentation for result reproducibility.

- Detailed Script Annotations: Provided a fully annotated R script as part of the project documentation, detailing each step of the data simulation, analysis, and visualization processes.

- Complete Data and Code Availability: Included all data generation scripts and analysis codes in the submission, enabling full reproducibility of the study from data creation through to final analysis.

**Version Control:** Effective use of version control systems.

- Systematic Use of Git: Managed the project development using Git, with regular commits and detailed commit messages that clearly describe each update to the project.

- Version Tracking and Collaboration: Utilized branching and merging strategies effectively to manage changes and collaborate with peers, ensuring seamless integration of different parts of the project.

# 3  Project Reflection

*Reflect on the experience of creating your data product. In 6 bullet points and at most 1 page total, summarise the following.*

- *3 things you have learned as part of this process,*

- *2 aspects of the project that you found challenging or would approach differently with hindsight,*

- *1 aspect of the project that you would like to learn more about in the future.*

*You may delete this italicised text when filling in the template.*

**Learnings:**

- Statistical Model Selection: Learned the importance of selecting the appropriate statistical model based on the data structure and analysis goals, which significantly affects the accuracy and interpretability of results.

- Data Simulation Techniques: Gained a deeper understanding of how to simulate realistic datasets in R, which is crucial for testing statistical models and methods in the absence of real data.

- Effective Communication of Statistical Findings: Developed skills in presenting complex statistical outcomes in a clear and accessible manner, ensuring that findings are understandable to stakeholders with different levels of expertise.

**Challenges:**

- Model Mis-specification: Faced difficulties in initially choosing the correct model specifications for the nonlinear data, which taught the importance of exploratory data analysis before model fitting.

- Balancing Complexity and Usability in User Interfaces: Struggled to maintain a balance between adding advanced functionalities and keeping the user interface intuitive and user-friendly, which could be improved with more feedback loops in the design phase.

**Further Development:**

- Machine Learning Techniques: Interested in exploring more advanced data analysis techniques, particularly machine learning models, to handle more complex datasets and improve predictive accuracy beyond traditional statistical methods.