



**COMPARISON BETWEEN IMPUTATION METHOD FOR
HANDLING MISSING DATA**

**AYUNIE BINTI EZADIN
NUR IZZATY BINTI CHUMIN
SITI NUR IZZATULNISA BINTI SALIT**

**BACHELOR OF SCIENCE (HONS.) STATISTICS
FACULTY OF COMPUTER AND MATHEMATICAL SCIENCES**

JANUARY 2021

ABSTRACT

This paper presents imputation method for the National Institute of Diabetes and Digestive and Kidney Diseases data from Arizona, United States. Missing data occurs in this data for five variables which are plasma glucose concentration, diastolic blood pressure, triceps skin fold thickness, serum insulin intake and body mass index (BMI). Missing data leads to problem that can cause bias and invalid conclusions to be made. This research objectives are to improve the data by filling the missing value and to compare which imputation method is better to handle missing value in a data set. In this research, imputation method and evaluation of the performance are applied for this data using Rstudio software. Five imputation methods used in this paper are Mean imputation method, K-Nearest Neighbour (KNN) imputation method, Multiple imputation method, Hot-Deck imputation method and Regression imputation method. The performance of these methods are evaluated using statistical analysis, coefficient of determination (R^2), mean-squared error (MSE), root of mean square error (RMSE), mean absolute error (MAE), index of agreement (d) and bias (B). Based on the result obtained from this research, it can be concluded that K-Nearest Neighbour imputation method is the best method among the five methods that are applied to handle the missing value. Conclusions are made as K-Nearest Neighbour (KNN) imputation method shows the best performance and has the lowest error value compared to other methods.

ACKNOWLEDGEMENT

IN THE NAME OF ALLAH, THE MOST GRACIOUS, THE MOST MERCIFUL

First and foremost, praises and thanks to the God, the Almighty, for his showers and blessings throughout our journey in completing the research work smoothly and successfully.

We would like to express our very great appreciation and sincere gratitude to our research supervisor, Madam Zaitul Anna Melisa Md Yasin, lecturer of University Teknologi MARA Negeri Sembilan, campus Seremban 3 for providing us this research topic and for giving us opportunity to do this research. We feel blessed as Madam Zaitul Anna Melisa chose our group to do this research and giving us chance to explore and learn new thing. She provides the guidance about the overall picture of this research and it really helps us in understanding and completing the research. We also would like to thank her for pointing out our mistakes, giving feedback on what we have done in this research and monitoring our progress weekly to prevent any difficulties towards the end of the due for the submission. Our supervisor has been such a great help to us in our journey and we cannot complete this research smoothly without her guide. We are also happy to work under our supervisor as she is someone that is approachable, and we are comfortable to work on this research with her.

Next, we would like to thank our lecturer who taught us code subject STA650 for last semester and teach us code subject MSP660 for this semester, Dr. Nurul Nisa' Khairol Azmi. She has been teaching us for 2 semester on how to complete our proposal and our paper work for final year project. Dr Nurul Nisa' has been guiding us on our final year project from chapter 1 to chapter 5, which includes introduction, literature review, methodology, results and discussion, and conclusion and future recommendation. We are grateful for her guidance as we can complete our paper work successfully. We also want to thank her for the template of latex that has been provided and teaching us on how to do the coding of the latex.

We are also extremely grateful to our parents for their prayers and love towards us. We are thankful for their sacrifice and effort for providing us the needs to further our study. Our family understanding to our busy schedule also have been a help for us in completing this research. They have always been supporting us physically and mentally in completing this research and always encourage us to do our best for this research. Lastly, our thank goes to all people who have help and supported us in this research directly or indirectly.

TABLE OF CONTENTS

ABSTRACT	i
ACKNOWLEDGMENT	ii
TABLE OF CONTENTS	iv
LIST OF TABLES	v
LIST OF FIGURES	vi
Chapter 1 Introduction	1
1.1 Background of Study	1
1.2 Problem Statement	4
1.3 Research Objectives	5
1.4 Significance of Study	5
1.5 Scope and Limitation	7
Chapter 2 Literature Review	8
2.1 Introduction	8
2.2 Mean	8
2.3 K-Nearest Neighbor	10
2.4 Hot-deck	12
2.5 Regression	13
2.6 Multiple Imputation	14
2.7 Maximum Likelihood (ML)	16
2.8 R Programming	18
2.9 Mean Square Error	19
2.10 Mean Absolute Error	19
2.11 Root Mean Square Error	19
2.12 Coefficient of Determination, R^2	20
2.13 Index of agreement (d)	20
2.14 Bias (B)	21
2.15 Conclusion	21
Chapter 3 Methodology	22
3.1 Introduction	22
3.2 Description of Data	23
3.3 Imputation Method	25
3.3.1 Mean Imputation Method	25
3.3.2 K-nearest Neighbour Method	25
3.3.3 Multiple Imputation Method	26
3.3.4 Hot-Deck Imputation Method	26

3.3.5	Regression Method	27
3.4	Method of Analysis	27
3.4.1	Coefficient of determination (R^2)	27
3.4.2	Mean-squared error (MSE)	28
3.4.3	Root of mean square error (RMSE)	28
3.4.4	Mean absolute error (MAE)	28
3.4.5	Index of agreement (d)	28
3.4.6	Bias (B)	29
3.5	Conclusion	29
Chapter 4	Results and Discussions	30
4.1	Introduction	30
4.2	Descriptive Analysis	30
4.3	Mean Imputation method	33
4.4	K-nearest Neighbour Method	34
4.5	Multiple Imputation Method	35
4.6	Hot-Deck Method	36
4.7	Regression Method	37
4.8	Imputation Method Performance	38
4.9	Performance of Logistic Regression Before and After Imputation Method	39
4.10	Conclusion	41
Chapter 5	Conclusion and Future Recommendation	43
5.1	Conclusion	43
5.2	Future Recommendation	44
References		46
APPENDIX A		49
APPENDIX B		49
APPENDIX C		64