



HANDLING MISSING DATA TO IMPROVE GENERALIZATION PERFORMANCE OF MACHINE LEARNING CLASSIFIER

R.Savithiri.¹

C.Kalaiarasi² K. Varalakshmi³ A.Vijayanarayanan⁴

1, 2, 3, 4, Assistant Professor COMPUTER SCIENCE AND ENGINEERING
PERI INSTITUTE OF TECHNOLOGY

Abstract

In supervised learning, missing values usually appear in the training set. The missing values in a dataset may generate bias, affecting the quality of the supervised learning process or the performance of classification algorithms. These imply that a reliable method for dealing with missing values is necessary. In this project, we analyze the difference between iterative imputation of missing values and single imputation in real-world applications. We propose an iterative imputation method, in which each missing attribute-value is iteratively filled using a predictor constructed from the known values and predicted values of the missing attribute-values from the previous iterations. Meanwhile, we demonstrate that it is reasonable to consider the imputation ordering for patching up multiple missing attribute values, and therefore introduce a method for imputation ordering. We experimentally show that our approach significantly outperforms some standard machine learning methods for handling missing values in classification tasks.

Introduction

The missing data problem is arguably the most common issue encountered by machine learning practitioners when analyzing real-world data. In many applications ranging from gene expression in computational biology to survey responses in social sciences, missing data is present to various degrees. As many statistical models and machine learning algorithms rely on complete data sets, it is key to handle the missing data appropriately. Missing data problem is a common issue in most real-world studies. Since most statistical models and data-dependent machine learning (ML) algorithms could only handle complete datasets, the issue of how to approach missing values plays an important role in statistical inferences. Let Y be an $(N \times K)$ data matrix with i -th row $y_i = (y_{i1}, y_{i2}, \dots, y_{iK})$ where y_{ij} is the value of j -th feature for the i -th sample. Define the subset of observed values as Y_{obs} and missing values as Y_{mis} . Also, let $M = [m_{ij}]$ be the missing indicator matrix, where m_{ij} indicates whether y_{ij} is missing or not. Rubin (1976) defines three different missing mechanisms according to the conditional probability of the missing $\{m_{ij} = 1\}$, given the data. The mechanism of missing data is completely at random (MCAR) if the probability of missingness is

independent of all data values, missing or observed, $P(m_{ij} = 0 | Y) = g(\phi)$, $i = 1, \dots, N$, $j = 1, \dots, K$, where $g(\cdot)$ is a known link function and ϕ is the vector of unknown mechanism parameters. The missing mechanism is called missing at random (MAR) if the probability of missingness depends only on the observed data values, $P(m_{ij} = 0 | Y) = g(Y_{\text{obs}}; \phi)$, $i = 1, \dots, N$, $j = 1, \dots, K$. Finally, the mechanism is called missing not at random (MNAR) when the probability of missingness may also depend on the unobserved data even after conditioning on the observed ones. The missing mechanism for the likelihood inferences is ignorable when the MCAR or MAR assumptions hold with the additional condition of disjoint parameter spaces of the missing mechanism and the data model (see Little and Rubin, 2014; Tsiatis, 2007, for more details). One simple approach to analyze incomplete data is complete case (CC) analysis which discards all incomplete cases. This approach is logical only if the missing rate is considerably small or the missing data mechanism is MCAR (Little and Rubin, 2014). However, if the missing mechanism is MAR or MNAR or the missing rate is considerably high, the CC approach could highly influence statistical results. This is due to the fact that CC analysis makes no use of observed features of an incomplete case.

SYSTEM ANALYSIS EXISTING SYSTEM

In data mining process the biggest task of data preprocessing is missing value imputation. Imputation is a statistical process of replacing missing data with substituted values. Many clinical diagnostic datasets are usually incomplete. Excluding incomplete dataset from the original dataset can bring more problem than simplification. In this paper the machine learning techniques for missing value imputation have been explored using Ionosphere data from UCI repository. The data imputation problem has been approached using well-known machine learning techniques.

DEMERITS

Some algorithms cannot handle missing values properly, while some techniques give efficient results to estimate the missing values. It is very important to handle missing data because many machine learning algorithm performances reduce due to missing values.

PROPOSED SYSTEM

In the dataset, there are few missing values (yet found to be hyper parameter), and pre-processing with such missing values is a common yet challenging problem. Re-substitution will give biased results from the data to be observed for HD diagnosis and will certainly affect the value of the learning process in Machine Learning. In this approach, we impute each missing data attribute value by predicting its data value from non-missing data attributes. The experiments are conducted on benchmark medical datasets missing values ranging from 1.98% to 50.65% and compared with iterative imputation.

MERITS

In this approach, we impute each missing data attribute value by predicting its data value from non-missing data attributes. The experiments are conducted on benchmark medical datasets missing values ranging from 1.98% to 50.65% and compared with iterative imputation.

Data Sets

A. Dataset Obesity With Attributes and Class

Estimation of obesity levels based on eating habits and physical condition Data Set. Dataset include data for the estimation of obesity levels in individuals from the countries of Mexico, Peru and Colombia, based on their eating habits and physical condition. This data contains 17 attributes and 2111 records.

B. Imputation & Mice Implementation

After the missing value has been successfully created on purpose, the next step is to substitute the missing value with the imputation technique. This project will perform imputation with 2 techniques, which are Statistical and MICE .Done by using the fancy impute library , two ways to impute :KNN or K-Nearest Neighbor MICE or Multiple Imputation by Chained Equation for fancy impute : <https://pypi.org/project/fancyimpute/> . In fancy impute, the MICE algorithm is named as Iterative Imputer. The steps taken are: Copy data_mv to obesity_mice_mpute Initialize Iterative Imputer with name mice_imputer Impute using fit_transform on data

C. Correlation Coefficient

The correlation coefficient is a statistical measure of the strength of the relationship between the relative movements of two variables. The values range between -1.0 and 1.0. A correlation of -1.0 shows a perfect negative correlation, while a correlation of 1.0 shows a perfect positive correlation. A correlation of 0.0 shows no linear relationship between the movements of the two variables. variable value (y) based on a given independent variable (x). Hence, the name is Linear Regression.

$$r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2}}$$

It has two types:

- Simple Linear Regression

If a single independent variable is used to predict the value of a numerical dependent variable, then such a Linear Regression algorithm is called Simple Linear Regression.

- Multiple Linear regressions

If more than one independent variable is used to predict the value of a numerical dependent variable, then such a Linear Regression algorithm is called Multiple Linear Regression.

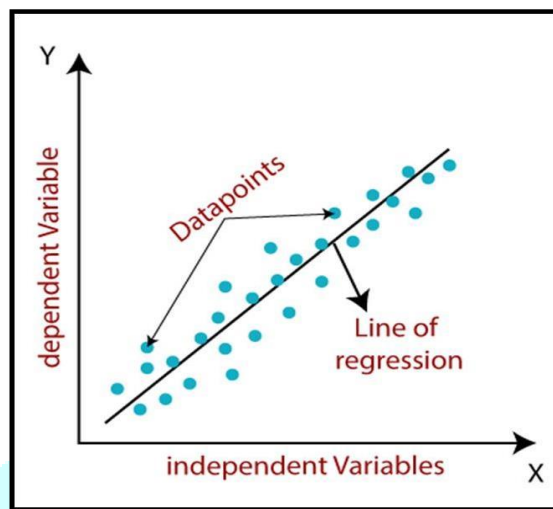


Fig 1: Regression

B. Logistic regression Algorithm

Logistic regression is basically a supervised classification algorithm. In a classification problem, the target variable (or output), y , can take only discrete values for a given set of features (or inputs), X . The model builds a regression model to predict the probability that a given data entry belongs to the category numbered as "1".

MICE ALGORITHM

A. LINEAR REGRESSION ALGORITHM

Linear Regression is a machine learning algorithm based on supervised learning. It performs a regression task. Linear regression performs the task to predict a dependent

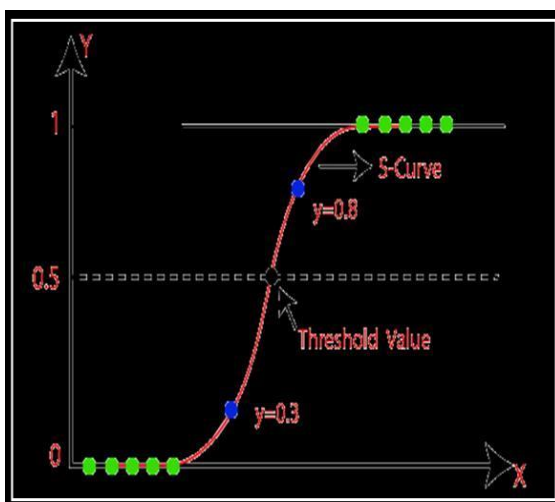


Fig 2: Linear Regression

C. RIDGE CLASSIFIER ALGORITHM

A Ridge regressor is basically a regularized version Linear Regressor. i.e. to the original cost function of linear regressor we add a regularized term that forces the learning algorithm to fit the data and helps to keep the weights lower a possible. The regularized term has the parameter, alpha "which controls the regularization of the model i.e helps in reducing the variance of the estimates. Cost Function for Ridge Regressor.

$$J(\theta) = \frac{1}{m}(X\theta - Y)^2 + \alpha \frac{1}{2}(\theta)^2$$

CONCLUSION

The choice of missing handling methodology has a significant impact on the clinical interpretation of the accompanying statistic analyses. With missing data, the choice of whether to impute or not, and choice of imputation method, can influence the clinical conclusion drawn from a regression model. We recommend researchers to perform a sensitivity analysis including at least MICE method. Often the common methods like mean, median, mode, frequent data and constant doesn't provide the correct data for the missing values. The model is only as good as the data, so having a complete dataset with proper data is a must; consider using MICE algorithm when you need to impute missing data.

FUTURE ENHANCEMENT

Missing data imputation using MICE was successfully performed on Missing at Random data using the fancy impute library. The results show that MICE performance is better than Statistical Imputation in terms of MSE and RMSE for numeric data, and Accuracy for categorical data. The MICE idea of using information from other columns proved helpful in replacing missing values. MICE is able to perform well on numeric and categorical data types. It is believed that MICE will produce the more excellent imputation when the correlation of each column is greater.

REFERENCES

- Little
- [1] Little & Rubin; "Missing data and imputation methods" Paper on Mean Imputation 1987.
 - [2] Zoubin Ghahramani and Michael I. Jordan ; "Supervised learning from incomplete data via an EM approach" 1994.
 - [3] Van Buuren, Stef Groothuis-Oudshoorn, Karin "Mice: Multivariate Imputation by Chained Equations in R" 2011.
 - [4] A. P. Dempster, N. M. Laird and D. B. Rubin "Maximum Likelihood from Incomplete Data via

theEMAAlgorithm” 1977.

- [5] Little, R.J.A. and Rubin, D.B. (1987) Statistical Analysis with Missing Data. John Wiley & Sons, New York.
- [6] Y. Burgette and Reiter “Multiple Imputation: A Review of Practical and Theoretical Findings” 2010.
- [7] Daniel J. Stekhoven and Peter Bühlmann “Miss Forest—non-parametric missing value imputation for mixed-type data” 2010
- [8] Dimitris Bertsimas & Bart Van Parys; “Sparse High-Dimensional Regression: Exact Scalable Algorithms and Phase Transitions” 2014
- [9] Dimitris Bertsimas & Jack Dunn “Optimal classification trees” 2017
- [10] Zhang Et Al & Cai Et Al; “From Predictive Methods to Missing Data Imputation: An Optimization Approach” 2018.

