

*Assignment 2: Written Assessment –
Classification Modelling Case Study*

ANALYSIS WRITE UP

Machine Learning – DAT-5303 – SFMBANDDD1
Professor: Chase Kusterer

March 18th, 2019

Anant Kumar

Table of Contents

<u>INTRODUCTION</u>	<u>3</u>
<u>INSIGHTS AND IMPLEMENTATION RECOMMENDATIONS.....</u>	<u>3</u>
TITLE	ERROR! BOOKMARK NOT DEFINED.
CULTURE	3
OTHER VARIABLES	3
<u>CONCLUSIONS.....</u>	<u>4</u>
<u>APPENDIX.....</u>	<u>5</u>

Introduction

The dataset for Game of Thrones (GoT) available for our assessment is 26 columns carrying qualitative, numerical and factor variables. Many of these variables are having missing values to the tune of 99% and also some mis-spelled words have been found. The assignment is to predict data for column 'isAlive' which is a binary data and hence a problem of logistic regression.

Key Insights

The dataset of GoT was scrubbed and found some of the followings are highly unsuitable for analysis. Before start of any kind of imputation, flags for missing values and outliers were created.

Title: 52% missing value with 263 unique values in the column which is 28% of the available record for the column. If we make it average, we can say that no value has more than four repetition in the column, Though, in actual, some of the value is highly repetitive but most of these are so sparse. We need to fill in the missing values and also lower the variable count in the column. To achieve these, I treated the missing value as "Unknown" and all those sparse entries (anything appearing less than 15 times) as "Other" to make the column suitable for analysis.

Culture: 65% missing values and just 9% unique values. Treated the same way as above.

dateOfBirth: Since a column is available for age of the characters, I preferred to discard the data.

mother: 99% missing values, 85% unique records from the available set of data. With such a high number of missing value, I found the data unsuitable for any kind of analysis.

father: 99% missing values and 80% unique records. Discarded the data for inclusion in analysis for the reasons.

heir: Discarded same as above for the same reasons. Missing values – 98%

spouse: It is name with 92% unique values. Preferred not to use in the analysis.

isAliveMother, isAliveFather, isAliveHeir, isAliveSpouse: all are binary data and can not be imputed. After creating missing flag.

age: Age has 77% missing values. It might prove to be an important variable and needs to imputed. To impute, going with mean or median age of the available dataset was one way. But, in my analysis, I observed, the median age of characters appearing first in different volumes of the book are different. Median age for character appearing first in first book, 2nd book, 3rd book, 4th book, 5th book and probably 6th (all those characters who did not appear in any of the first 5 books) book is 34, 21, 21, 27, 55 and 45 in that order. So while imputing data for age, I considered which book the character appeared first.

After completion of data cleaning, data were split (0.9: 0.1 train-test split) for training and testing of models. Models were initiated and fitted for the Logistic regression, KNN, Decision Tree, Random Forest, Support vector classifier and Gradient Boosting classifier. The model average score after CV = 3 are 0.777, 0.766, 0.722, 0.788, 0.766 and 0.739 respectively.

Conclusions

The result obtained from the above model building are consistent with accuracy around 0.78. The will likely to improve with more future engineering.

Appendix









