

Survival Analysis and Machine Learning Methodologies for Bone Marrow Cancer

Yitian Cai, Rudra Guin, Jorge Portugal

12/17/2021

Contents

Introduction	1
Data Description	1
Data Preprocessing	1
Analysis	2
Survival Analysis	2
Machine Learning Modeling and Analysis	2

```
## Working directory (to be set manually) - yt ultimately remove
## setwd("C:/Users/Rudra Guin/Documents/UMich/Biostatistics_625/biostat625_final_project")
```

Introduction

For this project, our group chose to perform statistical analysis using the Clinical Data Set from the Clinical Package of bone marrow cancer data. It was collected from the National Cancer Institute (NCI) genomic data portal. The version of the dataset we used was released by October 29, 2021. The bone marrow cancer can damage the bones, immune system, kidneys, and red blood cell count. The dataset contains 16,029 different cases of bone marrow cancer. We planned on using it for survival analysis and machine learning methodologies.

Data Description

The dataset contains 47 clinical variables. At the beginning, we tried to include the biospecimen data in order to increase the dimension of our dataset, but we found that most of those genomics data were missing, which would not provide us too much help. So, we decided to stick with the clinical data of bone marrow. Vital status means whether the patients were alive or dead at the time the data were collected. Our dataset contains 9,036 alive cases and 3,551 dead cases. Others were either missing or recorded as “unknown” or “not reported”.

Data Preprocessing

The initial dataset has several missing values among the different features. So, before going into analysis and modeling, we first clean up and pre-process the data. First, we remove all the observations that have missing values for vital status, because vital status is the one we wish to forecast or predict. Then, we decided to remove variables with more than 8,000 missing values. Because it means that more than half of

the observations in these variables were missing, which we believe will cause too much inaccuracy after data imputation. After that, we used `mice()` to impute those missing values. The `mice` package implements a method to deal with missing data and it generates Multivariate Imputations by Chained Equations. And finally, we randomly split the dataset into two where the training set contains 70% of the data, and the test set contains the rest 30% of the data.

```
##
## Attaching package: 'mice'

## The following object is masked from 'package:stats':
##
##   filter

## The following objects are masked from 'package:base':
##
##   cbind, rbind

## Loading required package: ggplot2
## Loading required package: lattice
## randomForest 4.6-14

## Type rfNews() to see new features/changes/bug fixes.

##
## Attaching package: 'randomForest'

## The following object is masked from 'package:ggplot2':
##
##   margin
```

Analysis

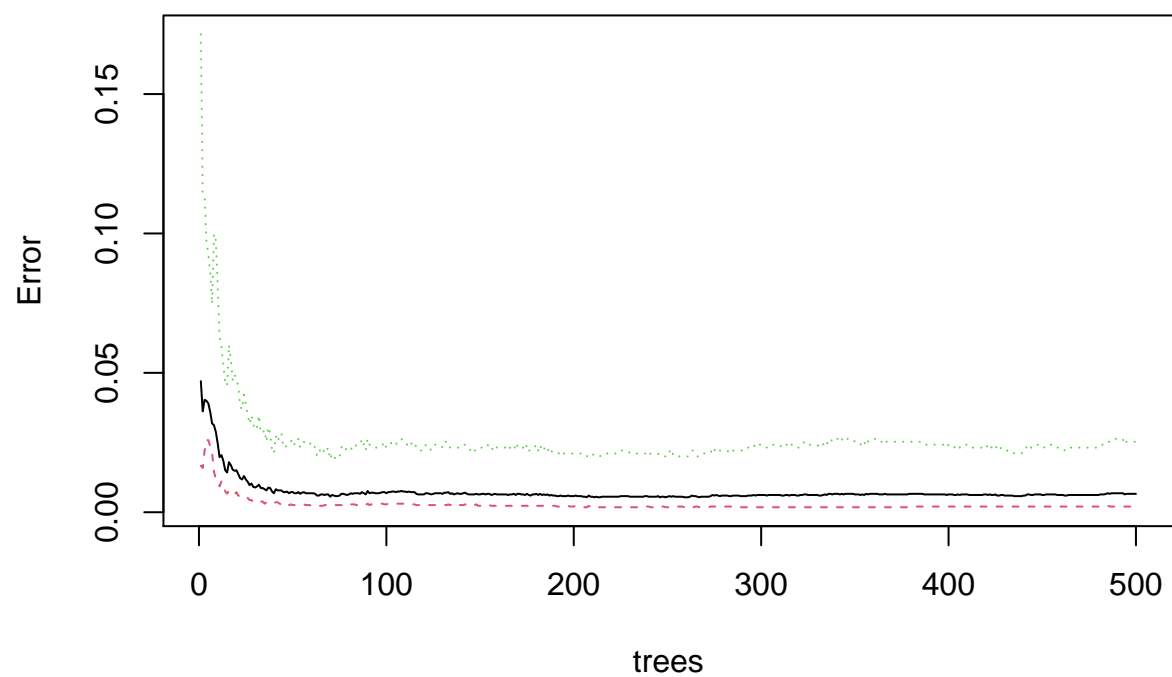
Survival Analysis

We first conducted some survival analysis

Machine Learning Modeling and Analysis

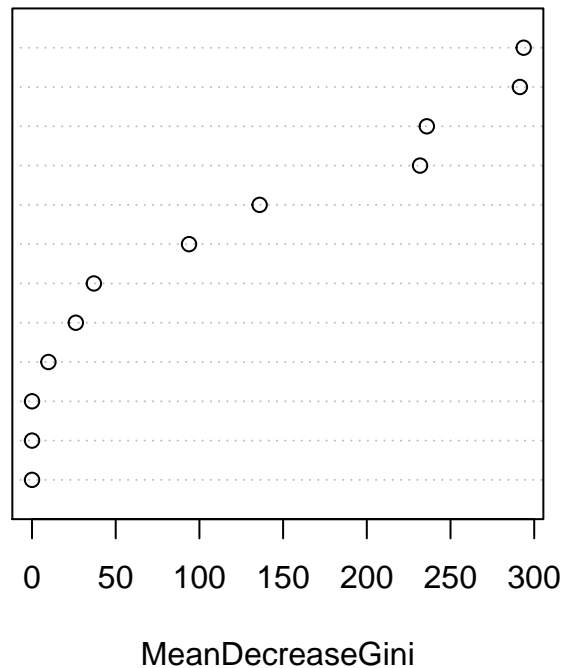
We also performed some machine learning modeling and evaluated the effectiveness of the models' prediction power of the vital status of certain patients.

Random Forest



Importance of the Variables

days_to_last_known_disease_status
 days_to_last_follow_up
 days_to_birth
 age_at_diagnosis
 age_at_index
 iss_stage
 race
 gender
 ethnicity
 project_id
 site_of_resection_or_biopsy
 tissue_or_organ_of_origin



Results:

```

p1 <- predict(rf, train_rf)
confusionMatrix(p1, train_rf$vital_status)

```

Confusion Matrix and Statistics

```

##
##           Reference
## Prediction Alive Dead
##      Alive 3896    1
##      Dead    0  948
##
##           Accuracy : 0.9998
##           95% CI : (0.9989, 1)
##      No Information Rate : 0.8041
##      P-Value [Acc > NIR] : <2e-16
##
##           Kappa : 0.9993
##
##  Mcnemar's Test P-Value : 1
##
##           Sensitivity : 1.0000
##           Specificity : 0.9989
##      Pos Pred Value : 0.9997
##      Neg Pred Value : 1.0000
##           Prevalence : 0.8041
##      Detection Rate : 0.8041

```

```

##      Detection Prevalence : 0.8043
##      Balanced Accuracy : 0.9995
##
##      'Positive' Class : Alive
##

p2 <- predict(rf, test_rf)
confusionMatrix(p2, test_rf$vital_status)

## Confusion Matrix and Statistics
##
##              Reference
## Prediction Alive Dead
##      Alive 1657    2
##      Dead    3  409
##
##              Accuracy : 0.9976
##              95% CI : (0.9944, 0.9992)
##      No Information Rate : 0.8015
##      P-Value [Acc > NIR] : <2e-16
##
##              Kappa : 0.9924
##
##  Mcnemar's Test P-Value : 1
##
##              Sensitivity : 0.9982
##              Specificity : 0.9951
##              Pos Pred Value : 0.9988
##              Neg Pred Value : 0.9927
##              Prevalence : 0.8015
##              Detection Rate : 0.8001
##      Detection Prevalence : 0.8011
##              Balanced Accuracy : 0.9967
##
##      'Positive' Class : Alive
##

```