



**NAVAL
POSTGRADUATE
SCHOOL**

MONTEREY, CALIFORNIA

THESIS

**PREDICTIVE STATISTICAL MODELING OF NAVAL
RESERVE OFFICERS TRAINING CORPS ATTRITION**

by

Zachary H. Swenson

December 2020

Thesis Advisor:

Ruriko Yoshida

Second Reader:

Jefferson Huang

Approved for public release. Distribution is unlimited.

THIS PAGE INTENTIONALLY LEFT BLANK

REPORT DOCUMENTATION PAGE			<i>Form Approved OMB No. 0704-0188</i>
Public reporting burden for this collection of information is estimated to average 1 hour per response, including the time for reviewing instruction, searching existing data sources, gathering and maintaining the data needed, and completing and reviewing the collection of information. Send comments regarding this burden estimate or any other aspect of this collection of information, including suggestions for reducing this burden, to Washington headquarters Services, Directorate for Information Operations and Reports, 1215 Jefferson Davis Highway, Suite 1204, Arlington, VA 22202-4302, and to the Office of Management and Budget, Paperwork Reduction Project (0704-0188) Washington, DC 20503.			
1. AGENCY USE ONLY (Leave blank)	2. REPORT DATE December 2020	3. REPORT TYPE AND DATES COVERED Master's thesis	
4. TITLE AND SUBTITLE PREDICTIVE STATISTICAL MODELING OF NAVAL RESERVE OFFICERS TRAINING CORPS ATTRITION			5. FUNDING NUMBERS
6. AUTHOR(S) Zachary H. Swenson			
7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES) Naval Postgraduate School Monterey, CA 93943-5000			8. PERFORMING ORGANIZATION REPORT NUMBER
9. SPONSORING / MONITORING AGENCY NAME(S) AND ADDRESS(ES) N/A			10. SPONSORING / MONITORING AGENCY REPORT NUMBER
11. SUPPLEMENTARY NOTES The views expressed in this thesis are those of the author and do not reflect the official policy or position of the Department of Defense or the U.S. Government.			
12a. DISTRIBUTION / AVAILABILITY STATEMENT Approved for public release. Distribution is unlimited.			12b. DISTRIBUTION CODE A
13. ABSTRACT (maximum 200 words) Attrition and the retention of talent is a concern for all organizations, but especially in the military. It is important to understand what factors influence attrition so that organizational leadership can optimize its purpose. This research examines attrition within Naval Reserve Officer Training Corps (NROTC) between 2013 and 2020. Using data provided by Naval Education Training Command (NETC), predictive statistical models aim to demonstrate what types of demographic, academic, and performance-based factors are important to predicting attrition among NROTC midshipmen. Modeling NROTC attrition behavior could enhance the NROTC Scholarship selection process, improve and organize attrition tracking in NROTC, and inform NETC on where resources should be allocated in order to improve program function. Furthermore, methods for tracking and predicting attrition, and for defining variable importance toward prediction, serve to inform other programs in decision-making processes surrounding attrition and talent retention.			
14. SUBJECT TERMS Naval Education Training Command, NETC, Naval Reserve Officer Training Corps, NROTC, attrition, predictive statistical modeling, classification			15. NUMBER OF PAGES 57
			16. PRICE CODE
17. SECURITY CLASSIFICATION OF REPORT Unclassified	18. SECURITY CLASSIFICATION OF THIS PAGE Unclassified	19. SECURITY CLASSIFICATION OF ABSTRACT Unclassified	20. LIMITATION OF ABSTRACT UU

THIS PAGE INTENTIONALLY LEFT BLANK

Approved for public release. Distribution is unlimited.

**PREDICTIVE STATISTICAL MODELING OF NAVAL RESERVE OFFICERS
TRAINING CORPS ATTRITION**

Zachary H. Swenson
Ensign, United States Navy
BSCE, Villanova University, 2019

Submitted in partial fulfillment of the
requirements for the degree of

MASTER OF SCIENCE IN OPERATIONS RESEARCH

from the

**NAVAL POSTGRADUATE SCHOOL
December 2020**

Approved by: Ruriko Yoshida
Advisor

Jefferson Huang
Second Reader

W. Matthew Carlyle
Chair, Department of Operations Research

THIS PAGE INTENTIONALLY LEFT BLANK

ABSTRACT

Attrition and the retention of talent is a concern for all organizations, but especially in the military. It is important to understand what factors influence attrition so that organizational leadership can optimize its purpose. This research examines attrition within Naval Reserve Officer Training Corps (NROTC) between 2013 and 2020. Using data provided by Naval Education Training Command (NETC), predictive statistical models aim to demonstrate what types of demographic, academic, and performance-based factors are important to predicting attrition among NROTC midshipmen. Modeling NROTC attrition behavior could enhance the NROTC Scholarship selection process, improve and organize attrition tracking in NROTC, and inform NETC on where resources should be allocated in order to improve program function. Furthermore, methods for tracking and predicting attrition, and for defining variable importance toward prediction, serve to inform other programs in decision-making processes surrounding attrition and talent retention.

THIS PAGE INTENTIONALLY LEFT BLANK

Table of Contents

1 Introduction	1
1.1 Background: Naval Reserve Officer Training Corps	1
1.2 Research Objectives	3
1.3 Structure.	3
2 Background	5
2.1 Literature Review	5
2.2 Advanced Data Analysis Methods and Definitions	6
3 Methodology and Modeling	11
3.1 Methodology Overview.	11
3.2 Data Description	12
3.3 Cleaned Data	14
3.4 Modeling	16
4 Results and Analysis	21
4.1 Model Performance	21
4.2 Model Inference.	27
4.3 Research Questions	33
5 Conclusion	35
5.1 Summary	35
5.2 Recommendations for Future Research.	35
List of References	37
Initial Distribution List	39

THIS PAGE INTENTIONALLY LEFT BLANK

List of Figures

Figure 2.1	Diagram on Statistical Learning: Unsupervised versus Supervised	7
Figure 2.2	Diagram on Statistical Learning: Classification versus Regression	8
Figure 2.3	Diagram on Statistical Learning: Parametric versus Non-parametric	9
Figure 2.4	Summary Diagram on Statistical Learning Algorithms	10
Figure 3.1	Methodology Overview	11
Figure 3.2	Visual of Commission Categorical Variable	15
Figure 3.3	Visual of DOR Categorical Variable	15
Figure 3.4	Visual of Nuke Categorical Variable	16
Figure 3.5	Model Classifiers and Response Variables	17
Figure 3.6	Training and Testing Data as Percentage of Total Available Data .	17
Figure 3.7	Training and Testing Data as Percentage of Total Available Data Source: James et al. (2013).	18
Figure 4.1	Confusion Matrix Example. Source: James et al. (2013).	21
Figure 4.2	Confusion Matrix and Model Performance Metrics: Model 1 . . .	23
Figure 4.3	Confusion Matrix and Model Performance Metrics: Model 2 . . .	23
Figure 4.4	Confusion Matrix and Model Performance Metrics: Model 3 . . .	24
Figure 4.5	Model Performance Metrics: Commission Model Comparison . .	24
Figure 4.6	Confusion Matrix and Model Performance Metrics: Model 4 . . .	24
Figure 4.7	Confusion Matrix and Model Performance Metrics: Model 5 . . .	25
Figure 4.8	Confusion Matrix and Model Performance Metrics: Model 6 . . .	25

Figure 4.9	Model Performance Metrics: Drop on Request (DOR) Model Comparison	25
Figure 4.10	Example Receiver Operating Characteristic (ROC) Curve Source: James et al. (2013).	27
Figure 4.11	ROC Curves: Commission Models	28
Figure 4.12	Area Under the Curve (AUC): Commission Models	28
Figure 4.13	ROC Curves: DOR Models	29
Figure 4.14	AUC: DOR Models	29
Figure 4.15	Decision Tree: Commission (Model 1)	30
Figure 4.16	Decision Tree: DOR (Model 4)	31
Figure 4.17	Random Forests (RF) Variable Importance (Mean Decrease of Accuracy): Commission (Model 3)	32
Figure 4.18	RF Variable Importance (Mean Decrease of Accuracy): DOR (Model 6)	33

List of Acronyms and Abbreviations

ADABOOST	Adaptive Boosting
AUC	Area Under the Curve
CART	Classification and Regression Tree
CP	College Programmer
DOR	Drop on Request
FITREP	Fitness Report
FN	False Negative
FP	False Positive
GPA	Grade Point Average
N	Negative
N*	Negative (Predicted)
NETC	Naval Education and Training Command
NPS	Naval Postgraduate School
NROTC	Naval Reserve Officer Training Corps
P	Positive
P*	Positive (Predicted)
PFA	Physical Fitness Assessment
PFT	Physical Fitness Test
RF	Random Forests

ROC	Receiver Operating Characteristic
SAT	Scholastic Aptitude Test
TN	True Negative
TP	True Positive
USNA	U.S. Naval Academy
URL	Unrestricted Line
USMC	U.S. Marine Corps
USW	Undersea Warfare
USN	U.S. Navy

Executive Summary

Attrition plays a pivotal role in any organization, especially in the military. Desired attrition levels vary among military commands—it is often based on the organization’s mission. For example, with a desire to be a small, highly qualified, and supremely motivated force, the Special Warfare community generally desires high attrition rates in their selection and training process. Other Navy communities, such as Naval Reserve Officer Training Corps (NROTC), target lower attrition rates. Regardless of the target level, organizational leadership is interested in understanding the factors that influence attrition—so that they can optimize its purpose.

This research examines attrition in NROTC from 2013 to 2020. The research objective is to understand what types of demographic, academic, and performance-based data influence attrition of NROTC midshipmen. To address this objective, predictive statistical models were trained on midshipmen data to classify commissions and drop on requests in NROTC.

The data used for this research was collected and provided by Naval Education and Training Command (NETC). The data is comprised of demographic (race, sex, ethnicity, scholarship status, Navy-option/Marine-option), academic (Grade Point Average, Scholastic Aptitude Test, Tier/major), and performance-based (aptitude score, Physical Fitness Assessment score, Nuclear Submarine Officer eligibility) measures.

Model performance and model inference results indicate RF models to be the most effective classifier on this data set for both NROTC commissions and DORs. The most important predictors for the RF classification models are scholarship status, aptitude score, Grade Point Average (GPA), and Physical Fitness Assessment (PFA) score. Further, academic and performance-based data prove far more effective as predictors—for both commissions and DORs—than demographic data. Lastly, eligibility to serve as a Nuclear Submarine Officer is not an important predictor for the trained models, despite it being highly correlated with commissioning.

This research demonstrates what types of factors influence attrition in NROTC, establishes analytic strategies for examining attrition in any organization, and aims to motivate decision-makers to collect meaningful data that can inform leadership on what factors lead to

attrition.

Acknowledgments

First, I want to thank my thesis advising team. Professor Yoshida and Professor Huang: thank you for your time and guidance. Your continued support was instrumental to my research.

Thank you to the entire Naval Postgraduate School team. My professors, program officer, academic associate, fellow students, and the entire NPS staff: thank you for the last 18 months. I am so grateful to have enjoyed this opportunity to earn my master's in such a special place.

Finally, to my family and friends. Dad, Mom, Elle, and Maddie: thank you for your unwavering love and support. To my roommates and friends: thank you for making my time at NPS and in Monterey so memorable.

THIS PAGE INTENTIONALLY LEFT BLANK

CHAPTER 1:

Introduction

This chapter introduces the premise and purpose of this research. The following sections provide necessary background information on Naval Reserve Officer Training Corps (NROTC) and outline research objectives, motivations, and structure.

1.1 Background: Naval Reserve Officer Training Corps

As the Navy's most prolific officer commissioning program, NROTC Scholarship Program exists to train and commission officers to serve in the Navy's Unrestricted Line (URL), Nurse Corps, and (Naval Education and Training Command 2020). NROTC is a college-based commissioning program with 63 NROTC units representing 160 different private and public universities nationwide.

NROTC scholarship awardees are selected, generally out of high school, via a highly competitive national application process. Students are responsible for gaining acceptance to a university on their own merit—then, assuming the university offers NROTC training, the scholarship is applied to that school on behalf of the awardee. Scholarship benefits include the full price of tuition, stipends for books and some living expenses, uniforms, and invaluable training and education. Unlike the U.S. Naval Academy (USNA), NROTC does not cover room and board costs. Naval Education and Training Command (NETC) is the overseeing command of NROTC.

NROTC is a multi-year program that runs concurrently with a student's normal college or university educational course of study. In addition to a normal academic workload leading to a Baccalaureate degree, Navy ROTC students attend classes in Naval Science, participate in the Navy ROTC unit for drill, physical training, and other activities, and are generally taught the leadership principles and high ideals of a military officer. During the summer break between school years, Navy ROTC students participate in a variety of training activities. These sessions help students understand various career options as well as familiarize them with a military life. (Naval Education and Training Command 2020)

Not all NROTC midshipmen are on scholarship. While it is a commissioning requirement to be on some type of Navy-funded scholarship, some midshipmen begin their NROTC journey as a College Programmer (CP). CP midshipmen do not receive any scholarship benefit outlined above—instead, they participate in NROTC without financial compensation. With high academic, physical, and leadership performance, CP midshipmen can earn the NROTC scholarship during their tenure (Department of Defense 2019).

1.1.1 Attrition Policy

On the contrary, not all midshipmen who start in NROTC graduate and commission. The motivation for this research is to examine if there are demographic, academic, performance-based data that can be used to predict whether a midshipman will commission or not. For the purposes of this research, the various reasons for not commissioning are classified into two broad categories.

- Unit-motivated attrition
- Midshipman-motivated attrition

Midshipmen dismissed from NROTC due to academic, medical, physical, disciplinary, or aptitude-based issues are classified as unit-motivated attrition. At the individual level, attrition of this type is unfortunate. At the institutional level, some amount of unit-motivated attrition is desired. Midshipman-motivated attrition describes disenrollment prompted by the student, also known as Drop on Request (DOR). At the individual level, this type of attrition is unfortunate, but possibly the best decision for the student. However, at the institutional level, attrition of this type is usually not desirable. A DOR represents an instance where a midshipman, that the Navy is both invested in and interested in retaining, leaves NROTC by their own volition.

It is important to note that there is certainly overlap between the two defined attrition categories. In the instance where a midshipman DORs due to their poor individual performance leading to discontent in NROTC, a valid argument could suggest that if this midshipman did not DOR own their own, they may have been dismissed soon thereafter. Nevertheless, examining the factors associated with DORs should be of great interest to the Navy, since at the time of attrition, the Navy was not forcing them out.

NROTC attrition, of any type or quantity, costs the Navy money—the amount of money is based on when, where, and why the midshipman dropped. The average cost per commission through NROTC is just under \$141,000 in present dollars (Office 1990). NROTC scholarship midshipmen who drop out any time before the start of their sophomore year are not forced to recoup any of the scholarship benefits. Beyond that time-frame, both unit-motivated and midshipman-motivated dropouts are contractually bound to pay back the Navy for all academic-associated costs afforded, unless it is waived for some extenuating circumstance. Benefits can be recouped financially or through enlisted military service (Department of Defense 2019).

1.2 Research Objectives

This thesis builds a predictive model that informs NETC about what demographic, academic, performance, and environmental factors influence NROTC attrition. As justified above, we focus on two critical attrition types: Commissions and DORs. Our purpose is to examine what characteristics and metrics are useful in predicting NROTC commissions—and if those factors maintain importance in predicting DORs. Furthermore, we examine if eligibility to serve as a Nuclear Submarine Officer (which requires higher academic than general NROTC mandate) serves as a strong predictor for a midshipman commissioning or dropping on request. In summary, the research is guided by the three main research questions.

1. What demographic, academic, and performance-based metrics are important to predicting whether a NROTC midshipman will commission?
2. What demographic, academic, and performance-based metrics are important to predicting whether a NROTC midshipman will DOR?
3. Is eligibility to serve as a Nuclear Submarine Officer an important factor in predicting whether a NROTC midshipman will commission or DOR?

1.3 Structure

Chapter 2 reviews a previous Naval Postgraduate School (NPS) thesis on NROTC attrition and summarizes the uses, advantages, and disadvantages of various types of statistical modeling. Chapter 3 outlines the research methodology, including data description and modeling strategies. Chapter 4 presents analytic results. Chapter 5 summarizes the thesis

and makes conclusions and recommendations for further research.

CHAPTER 2: Background

The first section of this chapter reviews a NPS thesis on NROTC Attrition. Its objectives, methods, and relevant findings motivate this study and demonstrate the need to assess NROTC attrition with more advanced data analysis techniques. The chapter proceeds to summarize general categories of data analysis and machine learning, providing necessary background information on the methods utilized in the following chapter.

2.1 Literature Review

Cahill (1993) examines attrition behavior among Midshipmen with four-year national NROTC scholarships between 1983 and 1987. Using data from NETC with a population size of 21,496 NROTC Midshipmen, Cahill (1993) provides summary statistics on attrition by various factors. Data set variables included: race/ethnicity (white, black, Hispanic, other), gender (male, female), academic major (technical, non-technical), and type of attrition (motivational, disciplinary, ineptitude, personal).

2.1.1 Research Objectives and Methods

Cahill (1993) was motivated by three research questions.

1. Who is dropping out of NROTC and why?
2. Are there similarities between attrition behavior of NROTC Midshipmen and civilian college dropouts?
3. Is NROTC attrition influenced by the service obligation policy?

To address the first research question, bivariate data analysis of NROTC attrition behavior was conducted using the Statistical Analysis System (SAS Institute 2017). Summary statistics were presented in order to compare NROTC attrition statistics of various demographics. With respect to the second research question, comparative college dropout behavior was vague, descriptive, and sourced through outside research. Summary statistics for NROTC attrition by each demographic was compared to parallel statistics for civilian college students. Finally, the third research question was investigated by evaluating what year (freshman,

sophomore, junior, or senior) midshipmen dropped out of NROTC and comparing it to the year when service obligation begins. In NROTC, this deadline was shifted forward from after sophomore year to after freshman year in 1983.

2.1.2 Relevant Findings

Descriptive statistical analysis of NROTC attrition from 1983 to 1987 found that men disenrolled at a higher rate than women. While men have a higher propensity for academic related dropouts than women, female midshipmen proportionally received more involuntary disenrollments than their male counterparts. Furthermore, overall attrition rates for racial and ethnic minorities were higher than those for white students. However, minority dropouts were less common at units with higher racial and ethnic diversity. Academic dropouts were more common amongst black midshipmen than Hispanic or white midshipmen. All demographic attrition patterns mirrored that of civilian college dropouts. Academic major was shown to have strong influence over NROTC attrition. Over half of all non-technical majors dropped from NROTC, compared to the 36 percent attrition rate for technical majors. Finally, over 64 percent of were classified as motivational or by choice of the student and not the unit. Today, that decision is referred to as a DOR. The major share of motivational dropouts occur right before the obligation deadline, leading to the conclusion that the service policy greatly influences when midshipmen drop.

2.2 Advanced Data Analysis Methods and Definitions

Cahill's research summarized general trends and speculated at possible reasons for attrition. By applying more advanced analytic tools to a more comprehensive, descriptive, and recent data set, this research aims to expand on Cahill's work. Using the machine learning algorithms described in the following sections, this research is dedicated to creating a predictive model that informs NETC about what demographic, academic, performance, and environmental factors influence NROTC attrition.

2.2.1 Supervised versus Unsupervised Learning

Statistical learning broadly describes a vast set of tools to understand data (James et al. 2013). The composition and availability of input and output data, as well as research objectives, dictate what types of statistical learning strategies are appropriate. Under the broad umbrella

of statistical learning, there are two general categories dictated by the availability of output data: supervised (predictive) and unsupervised (descriptive) learning. Supervised statistical learning involves building a statistical model for predicting or inference with an output variable based on one or more inputs (James et al. 2013). Unsupervised statistical learning, when there are inputs but no supervising outputs, learns relationships from the structure and patterns present in the input data (James et al. 2013). Figure 2.1 illustrates the distinction between unsupervised and supervised learning.

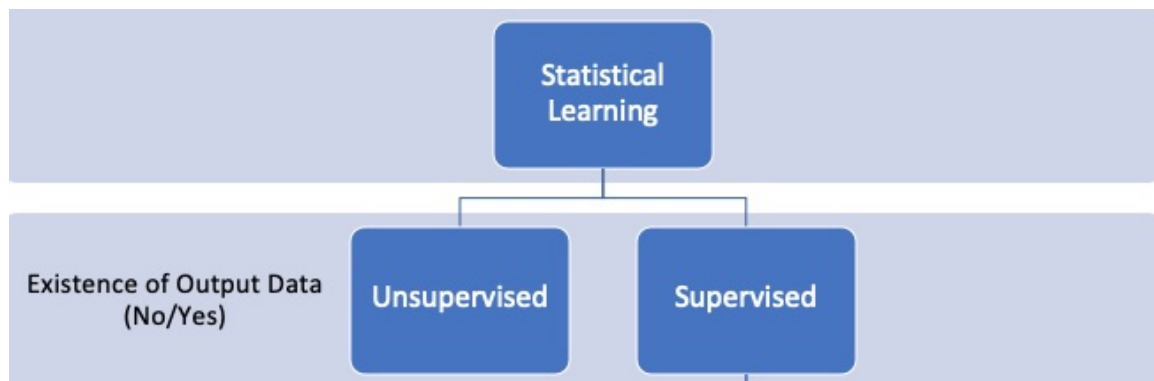


Figure 2.1. Diagram on Statistical Learning: Unsupervised versus Supervised

2.2.2 Classification versus Regression

Under both supervised learning are two general types of analysis that are dictated by the type of response variable: classification (categorical outputs) and regression (continuous outputs) shown in Figure 2.2. Classification algorithms aim to segregate observations into finitely many and more than one discrete categories (James et al. 2013). Common types of classifications are the belonging to a certain category (Yes/No) or descriptive categories, such as origin of cars (USA, Germany, Japan, etc.). Regression is used for continuous, quantitative outputs, such as the number of miles driven on a car before mechanical failure (James et al. 2013). Figure 2.2 illustrates the distinction between classification and regression.

2.2.3 Parametric versus Non-parametric Learning

There are two types of statistical learning that are dictated by the level of understanding of the relationship between inputs and outputs prior to analysis: parametric and non-parametric. Parametric learning assumes a form of the relationship between input and output data, such

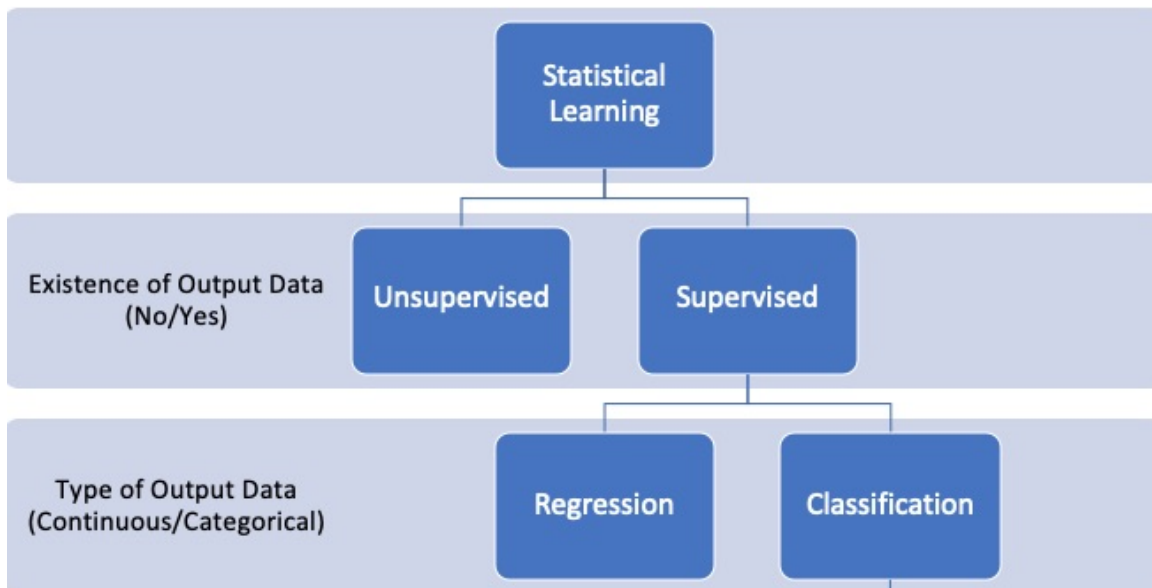


Figure 2.2. Diagram on Statistical Learning: Classification versus Regression

as a linear combination of input variables multiplied by unique coefficients (James et al. 2013). The parametric algorithms then aim to determine the value of these coefficients. Examples of parametric algorithms include linear regression, logistic regression, linear discriminant analysis, and simple neural networks. The strengths to parametric learning is in its simplicity, speed, and ability to perform on smaller data sets; the analysis is easier to understand and interpret. However, parametric learning is constrained by the assumption of the function having a certain form, which limits its ability to map complex relationships between inputs and outputs. Non-parametric learning does not assume anything about the form of the relationship between inputs and outputs, enabling algorithms to illustrate more complex relationships (James et al. 2013). While they require more data and computational investment, non-parametric learning generally produces more flexible and reliable results. However, this can sometimes adversely result in over-fitting and/or difficulty interpreting results. Examples of non-parametric learning algorithms are decision trees, k-nearest neighbors, and support vector machines. Figure 2.3 illustrates the distinction between parametric and non-parametric modeling.

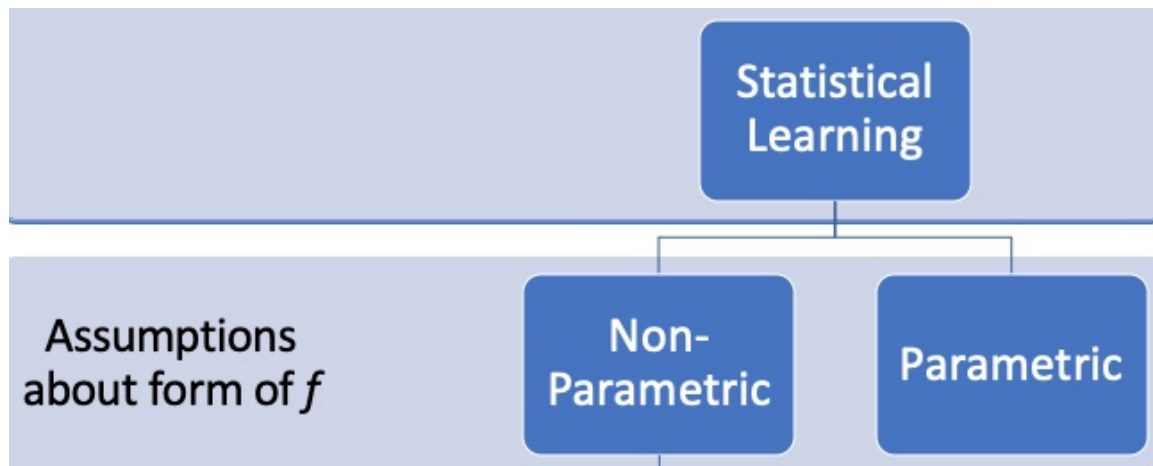


Figure 2.3. Diagram on Statistical Learning: Parametric versus Non-parametric

2.2.4 Decision Trees: CART, Random Forests and Adaptive Boosting

One of the most popular supervised, non-parametric machine learning algorithms is decision trees. Decision trees sort data into subsets. At the node of each split, or branch, is a rule that describes the split; the tree grows with branches subject to various rules until the data is categorized correctly. Decision trees can be used for classification or regression (James et al. 2013). Decision trees are intuitive, visually pleasing, and easy to interpret.

In order to reduce variance and enhance performance, decision trees can be aggregated using methods such as bagging, boosting, and random forests (James et al. 2013). Random Forests are an ensemble of many decision trees generated on random subsets of the data with random subsets of predictors. Random forests reduce variance and improve accuracy, robustness, and efficiency. Further, their composition can be used for self-evaluation such as internal error estimation and variable importance.

The efficiency and accuracy of decision trees can also be improved through adaptive boosting (James et al. 2013). Within classification, adaptive boosting produces a weighted linear combination of weaker classifiers, which when combined, achieves higher accuracy. Starting with the first classifier, adaptive boosting weighs incorrectly classified observations heavily, so that the next classifier learns from the previous one. The final model integrates the results from all of the classifiers to produce the highest performing model. Figure 2.4 displays a summary diagram of supervised and unsupervised learning algorithms as well as parametric

and non-parametric modeling.

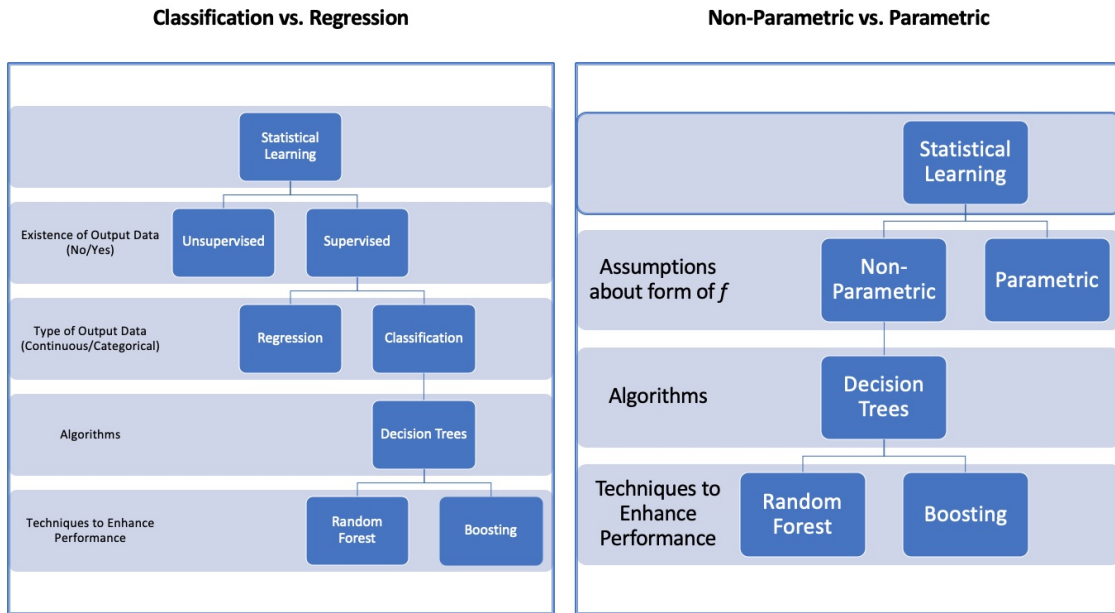


Figure 2.4. Summary Diagram on Statistical Learning Algorithms

CHAPTER 3: Methodology and Modeling

This chapter presents the methodology and modeling strategy for analyzing NROTC Attrition data in this thesis. The following sections describe the data, summarize general data cleaning techniques, introduce specifically constructed response and proxy variables, and outline modeling approaches.

3.1 Methodology Overview

Figure 3.1 illustrates the methodology used to create analytical models capable of predicting attrition behavior of NROTC midshipmen based on available demographic and performance data. There were six basic steps to analysis.

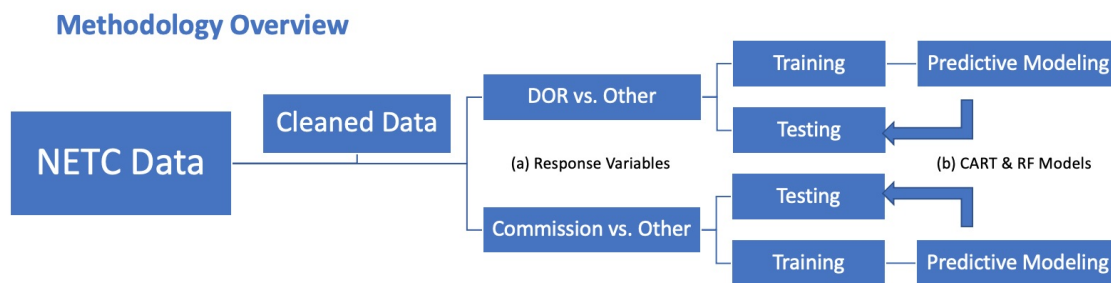


Figure 3.1. Methodology Overview

1. Obtain NROTC midshipmen demographic, performance, and attrition data.
2. Clean the data to make it analyzable.
3. Create categorical response variables of interest.
4. Split clean data into training and testing sets.
5. Fit various models on training data.
6. Assess model performance on testing data.

The remaining sections of the chapter describes each of these steps.

3.2 Data Description

The data used for this research was obtained through NETC. It contains demographic, performance, and attrition data for each and every midshipmen enrolled in NROTC at any time between 2013 and 2020. Each midshipmen is populated in the data every semester enrolled in NROTC—for example, a midshipmen who participates in 4 years of NROTC, which is standard for those who graduate and commission, appears in the data eight times. For this standard midshipman, demographic data across the eight observations remains constant (unless they transfer universities or something else out of the ordinary), while performance data changes over time. Of the eight observations, attrition data will only populate once – when the midshipmen commissions or withdraws from NROTC. Midshipmen who participate in less than (or more than) 4 years in NROTC only have data observations from their time enrolled – the example of the "standard midshipmen" is simply provided to clarify the structure of the original data provided by NETC. Data was compiled by NETC in Microsoft Excel (Microsoft Organization 2019) and then imported into RStudio (RStudio Team 2020). This resulted in 117,787 original data entries.

3.2.1 Data Variables

The following describes each variable in the data provided by NETC. Not every variable listed was analyzed, as many are repetitive and/or not related to attrition—those ultimately used in modeling NROTC Attrition are in bold letters.

- *YRMO*: Year and month that data was entered
- *SCHOOL.CODE*: Numeric code associated with a given NROTC unit
- *STUDENT.TYPE*: Midshipman
- *SCHOOL.NAME*: Name of University hosting NROTC unit
- *STUDENT.ID*: Numeric ID associated with specific NROTC midshipman
- *RECORD.STATUS*: Active (currently enrolled in NROTC/Inactive
- ***SEX*: Male/Female**
- ***OPTION*: Navy/Marine**
- ***PROGRAM.TYPE*: Scholarship/College Program**
- *PROGRAM.CODE*: Two digit numeric/alphabetical code associated with a specific type of scholarship, college program status, or advanced standing
- *PROGRAM.CODE.DESCRPTION*: Description of *PROGRAM.CODE*, Navy Nurse

- Option distinguished
- **RACE:** Caucasian, African American, Multiple Races, Native Hawaiian and other Pacific Islander, American Indian, Asian, decline to answer
 - **HISP:** Hispanic/Non-Hispanic
 - **ESTIMATED.COMMISSIONING.DATE:** Year and month of estimated commissioning
 - **NAVAL.SCIENCE.YEAR:** 4/C, 3/C, 2/C, 1/C
 - **GPA.CUM.LATEST:** Latest recorded cumulative Grade Point Average (GPA)
 - **APT.FINAL.LATEST:** Latest recorded aptitude score, calculated via Fitness Report (FITREP)
 - **DATE.ATTRITE:** Year and month that the midshipmen left NROTC to commission or disenroll
 - **ATTRITE.CODE:** Four digit numeric/alphabetical code associated with the type of attrition
 - **ATTRITE.CATEGORY:** Commission, Academic, Physical, Inaptitude, DOR, Disciplinary, Other
 - **ATTRITION.DESCRPTION:** Description associated with *ATTRITE.CATEGORY*
 - **TIER:** Category of academic major: Tier 1 (engineering programs of Navy interest), Tier 2 (other engineering, math, and science programs) or Tier 3 (foreign language and remaining programs) (Department of Defense 2019)
 - **DATE.REPORTED:** Year and month midshipmen enrolled in NROTC
 - **SAT.MATH:** Scholastic Aptitude Test (SAT) Math score
 - **SAT.VERB:** SAT Verbal score
 - **SAT.COMP:** SAT Composite score (Math + Verbal)
 - **CURRENT.PFA.DATE:** Physical Fitness Assessment (PFA) Cycle of latest completed PFA
 - **CURRENT.PFA.SCORE:** Latest PFA/Physical Fitness Test (PFT) score (0-100 for Navy, 0-300 for Marines)
 - **CURRENT.PFA.STATUS:** Pass or Fail

3.3 Cleaned Data

Because we are interested in understanding the relationship between demographic and performance features and types of attrition, we cleaned the data (originally 117,787 entries) to only include observations that had populated data for *ATTRITE.CATEGORY*. The result was 18,313 data entries matching 18,313 distinct midshipmen who had enrolled in NROTC at any time between 2012-2020 and exited the program (via commission or attrition). We further condensed the data to only include observations where all categories of predictive interest (GPA, Tier, Aptitude score, PFA score, SAT score) contained data. Unfortunately, many observations did not include information from these categories, reducing the size of the final data set to 5,924 entries. Finally, Navy PFA and Marine PFT scores were standardized so that Navy-Option and Marine-Option midshipmen had consistent and analytically comparable performance data.

3.3.1 Proxy and Response Variables

Commission versus Other

Based on the research questions, we are interested in understanding the influence of various demographic and performance data on attrition behavior in NROTC. For modeling this relationship, instead of using *ATTRITION.CATEGORY* as the response variable, a categorical variable called *COMMISSION* was created. *ATTRITION.CATEGORY* has seven possible values (see Data Description). With the intention of optimizing model accuracy, efficiency, and specificity, *COMMISSION* only have two responses (Yes/No) and is formulated via the data in *ATTRITION.CATEGORY*. Figure 3.2 illustrates how the *COMMISSION* variable was derived. The creation of the *COMMISSION* variable concentrates our analysis on what demographic and performance measures are important to predicting if a midshipman will commission, rather than distinguishing between seven different types of attrition.

Drop on Request versus Other

We are also interested in understanding the influence of various demographic and performance data on whether or not a midshipman will DOR. Categorical variable *DOR* was created to explore this relationship. Similar to the *COMMISSION* variable, it is a Yes/No category derived from the *ATTRITION.CATEGORY* data, illustrated by Figure 3.3. Using

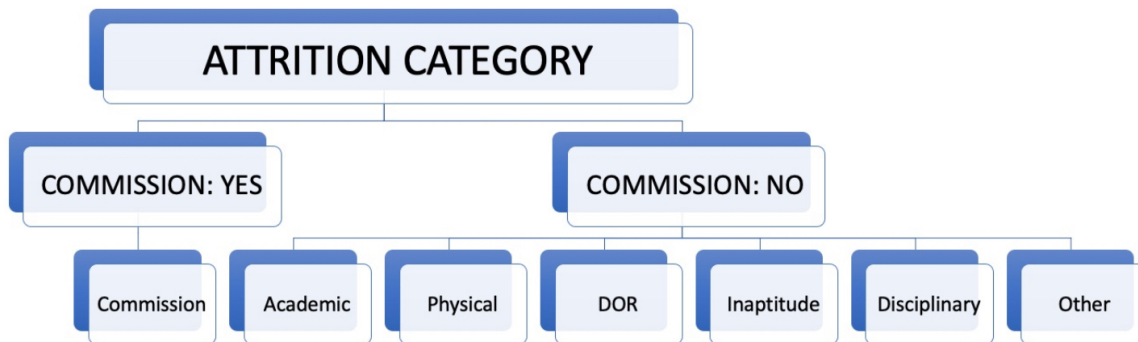


Figure 3.2. Visual of Commission Categorical Variable

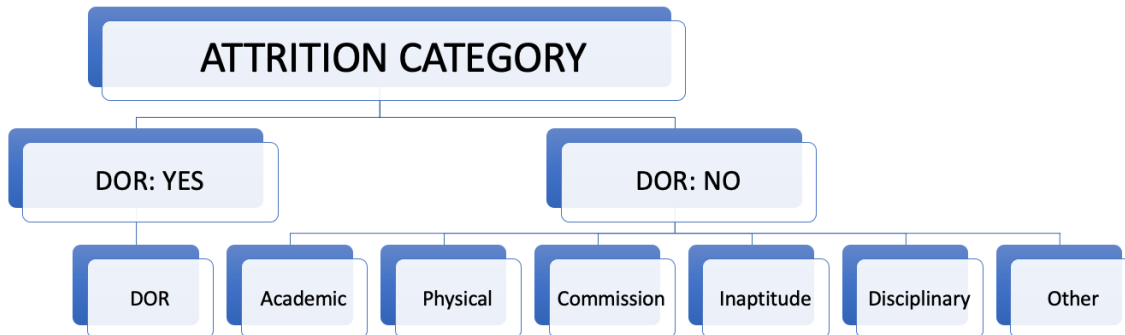


Figure 3.3. Visual of DOR Categorical Variable

DOR as a response variable substantiates understanding of what demographic and performance features are associated with a midshipmen deciding to DOR. Having two models (one that predicts commissions and one that predicts DORs) will illustrate if the critical variables to each model are the same or different.

Nuke Eligible

To address research question 3, we created proxy variable *NUKE*. Derived from data variables *GPA.CUM.LATEST*, *OPTION*, *PROGRAM.CODE*, and *TIER*, categorical variable *NUKE* (Yes/No) indicates if a midshipman was eligible to be drafted into the Submarine Officer or Surface Warfare Officer (Nuclear) communities. We used it as a categorical predicting variable in the Commission and DOR models. A Navy-Option, non-Nurse midship-

man that met the following Tier-specific GPA requirements were deemed "Nuke Eligible" (Department of Defense 2019).

- Tier 1: GPA > 2.9
- Tier 2: GPA > 3.0
- Tier 3: GPA > 3.2

Figure 3.4 displays the formulation of the *NUKE* categorical variable.

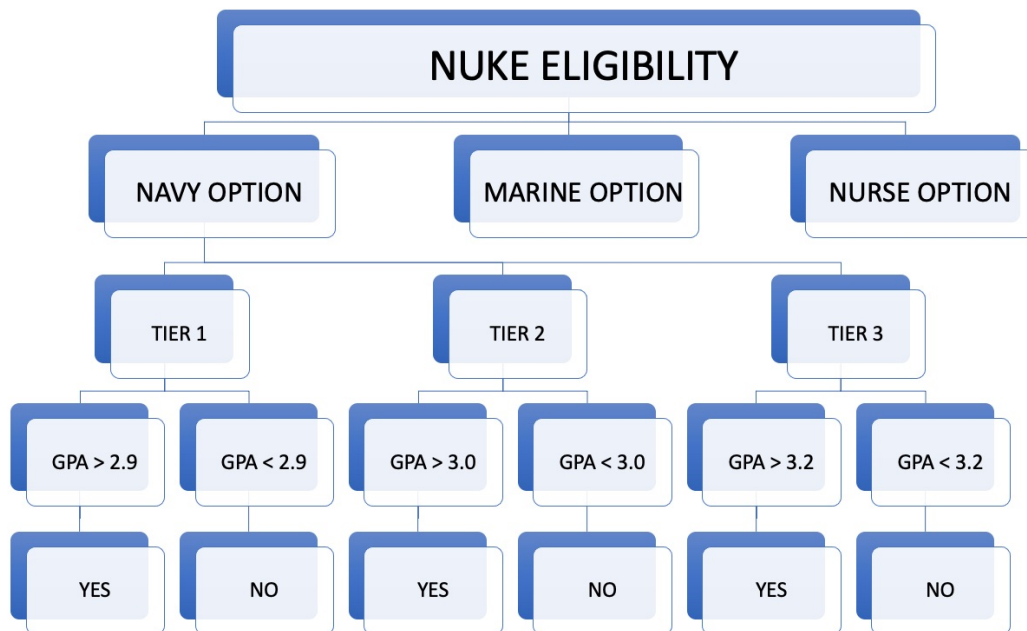


Figure 3.4. Visual of Nuke Categorical Variable

3.4 Modeling

We fit distinct models on training data to predict two different responses: Commission (Yes/No) and DOR (Yes/No). Figure 3.5 summarizes the classifiers and response for each trained model. We assessed model performance on a smaller testing subset of the data.

3.4.1 Training Data

As illustrated in 3.6, clean data was divided into random and representative training and testing sets. Training data is a set of observations used to train, or teach, a model to predict a

Model	Classifier Type	Classifier Function	Response
1	Decision Tree	rpart()	COMMISSION
2	Adaptive Boosting	adaboost()	COMMISSION
3	Random Forest	randomForest()	COMMISSION
4	Decision Tree	rpart()	DOR
5	Adaptive Boosting	adaboost()	DOR
6	Random Forest	randomForest()	DOR

Figure 3.5. Model Classifiers and Response Variables

Cleaned NROTC Attrition Data

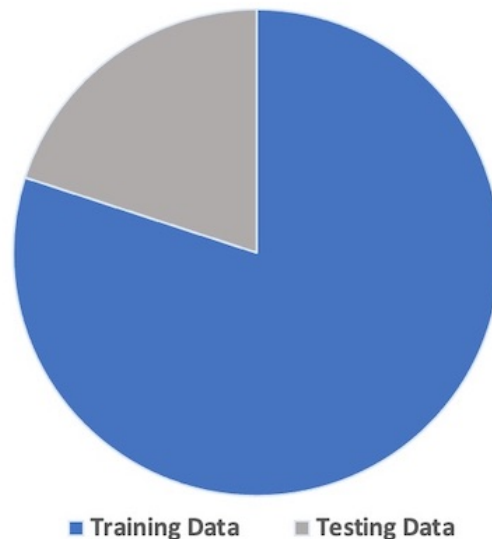


Figure 3.6. Training and Testing Data as Percentage of Total Available Data

response (James et al. 2013). For this research, 80 percent of the clean data was designated for training the models, with the remaining 20 percent reserved for testing.

3.4.2 Decision Tree Classifier

Models 1 and 4 predict Commission and DOR, respectively, using the *rpart()* classifier from the *rpart* Rstudio package (Therneau and Atkinson 2019). Decision trees partition data observations with rules aimed towards partitioning the entire data set. While not always

extremely accurate, decision trees are intuitive, visually pleasing, and easy to interpret (James et al. 2013). See Figure 3.7 for an example of a decision tree.

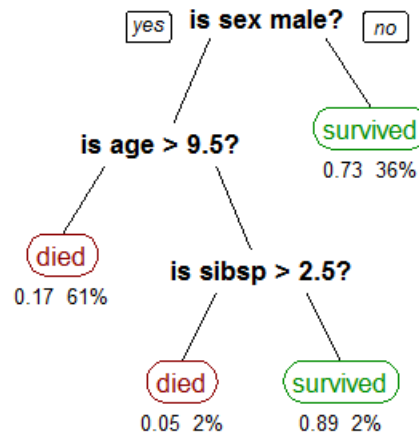


Figure 3.7. Training and Testing Data as Percentage of Total Available Data
Source: James et al. (2013).

3.4.3 Adaptive Boosting Classifier

Models 2 and 5 predict Commission and DOR, respectively, using the *adaboost()* classifier from the *fastAdaboost* Rstudio package (Chatterjee 2016). To improve classification accuracy of decision trees, adaptive boosting produces a weighted linear combination of many, weaker decision trees. Adaptive boosting weighs incorrectly classified observations heavily, so that the next classifier learns from the previous one. The final model integrates the results from all of the classifiers to produce the highest performing model – however, it’s product becomes much less intuitive and interpretable (James et al. 2013).

3.4.4 Random Forests Classifier

Models 3 and 6 predict Commission and DOR, respectively, using the *randomForest()* classifier from the *randomForest* R package (Liaw and Wiener 2002). Random Forests are an ensemble of many decision trees generated on random subsets of the data with random subsets of predictors. Random forests reduce variance and improve accuracy, robustness,

and efficiency. Further, their composition can be used for self-evaluation such as internal error estimation and variable importance—however, it’s output is also less intuitive and interpretable (James et al. 2013).

3.4.5 Testing Data

Finally, we assessed model performance. Testing data was not used in model training—it was specifically reserved to evaluate model accuracy (James et al. 2013). The trained model used training data inputs to predict the response—model performance was evaluated based on how accurately the predictions are (as compared to the known outputs of the testing data).

THIS PAGE INTENTIONALLY LEFT BLANK

CHAPTER 4: Results and Analysis

This chapter presents modeling results and analysis of NROTC Attrition. The following sections summarize model performance, review practical inferences gleaned from modeling, and address each of the research questions.

4.1 Model Performance

Each subsection describes a particular model performance metric: its definition, use, and significance. Then, performance metric results are displayed and analyzed.

4.1.1 Confusion Matrix

Under supervised learning classification, where a response variable is both categorical and known for all of the data, a confusion matrix is used to determine how well a model predicts the response variable of the testing data subset (James et al. 2013). A confusion matrix displays both the true and predicted response variable values in a 2x2 matrix. Figure 4.1 displays an example confusion matrix, with the following important terms associated within each quadrant of the matrix (James et al. 2013).

- True Negative (TN) denotes testing data observations that were predicted *NO* and are truly *NO*.
- False Negative (FN) represents testing data observations that were predicted *NO* but were actually *YES*.

		<i>Predicted class</i>		
		<i>- or Null</i>	<i>+ or Non-null</i>	Total
<i>True class</i>	<i>- or Null</i>	True Neg. (TN)	False Pos. (FP)	N
	<i>+ or Non-null</i>	False Neg. (FN)	True Pos. (TP)	P
Total		N*	P*	

Figure 4.1. Confusion Matrix Example. Source: James et al. (2013).

- False Positive (FP) designates testing data observations that were predicted *YES* but were actually *NO*.
- True Positive (TP) are testing data observations that were both predicted *YES* and are truly *YES*.
- Positive (P) is the total number of true positives in the testing data.
- Negative (N) is the total number of true negatives in the testing data.
- Positive (Predicted) (P*) is the total number of predicted positives by the model.
- Negative (Predicted) (N*) is the total number of predicted negatives by the model.

A confusion matrix is particularly useful in assessing model performance because it displays so much more than just classification accuracy. Using the values displayed in 4.1, analysts can calculate the following performance metrics.

- Accuracy Rate: How often are the true values congruent with predicted values?

$$\text{Accuracy} = (\text{TN} + \text{TP}) / (\text{N} + \text{P}). \quad (4.1)$$

- Error Rate: How often are true values not congruent with predicted values?

$$\text{Error} = (\text{FN} + \text{FP}) / (\text{N} + \text{P}). \quad (4.2)$$

- Sensitivity: How often are actual positives predicted positive?

$$\text{Sensitivity} = \text{TP} / \text{P}. \quad (4.3)$$

- Specificity: How often are actual negatives predicted negative?

$$\text{Specificity} = \text{TN} / \text{N}. \quad (4.4)$$

- Type I Error: How often are actual negatives predicted positive?

$$\text{Type I Error} = \text{FP} / \text{N}. \quad (4.5)$$

- Type II Error: How often are actual positives predicted negative?

$$\text{Type II Error} = \text{FN} / \text{P}. \quad (4.6)$$

Each performance metric is uniquely useful in understanding the model’s strengths and shortcomings, as well as drawing model inferences. The level of importance for each of these metrics is situational—it depends on the subject, data structure, and research objectives.

Results: Commission

As described in Chapter 3, Models 1, 2 and 3 predict commissions in NROTC using decision tree, adaptive boosting, and random forests algorithms, respectively. Figure 4.2 displays the confusion matrix and model performance metrics (derived from the confusion matrix) for Model 1; Figure 4.3 displays the confusion matrix and model performance metrics for Model 2; Figure 4.4 displays the confusion matrix and model performance metrics for Model 3. Figure 4.5 compares performance metrics for each of the models that predict commissions, with the most desirable results bold and italicized.

<u>Confusion Matrix</u>			
n = 1,185	Actual: No	Actual: Yes	Total
Predicted: NO	119	39	158
Predicted: YES	77	950	1,027
Total	196	989	1,185
<u>Performance Metrics</u>			
Accuracy	0.9021	Error	0.0979
Sensitivity	0.6071	Specificity	0.9606
Type I Error	0.3928	Type II Error	0.0394

Figure 4.2. Confusion Matrix and Model Performance Metrics: Model 1

<u>Confusion Matrix</u>			
n = 1,185	Actual: No	Actual: Yes	Total
Predicted: NO	118	42	160
Predicted: YES	78	947	1,025
Total	196	989	1,185
<u>Performance Metrics</u>			
Accuracy	0.8987	Error	0.1013
Sensitivity	0.6020	Specificity	0.9575
Type I Error	0.3980	Type II Error	0.0425

Figure 4.3. Confusion Matrix and Model Performance Metrics: Model 2

<u>Confusion Matrix</u>			
n = 1,185	Actual: No	Actual: Yes	Total
Predicted: NO	123	34	157
Predicted: YES	73	955	1,028
Total	196	989	1,185
<u>Performance Metrics</u>			
Accuracy	0.9097	Error	0.0903
Sensitivity	0.6276	Specificity	0.9656
Type I Error	0.3724	Type II Error	0.0344

Figure 4.4. Confusion Matrix and Model Performance Metrics: Model 3

Model	Classifier	Accuracy	Error	Sensitivity	Specificity	Type I Error	Type II Error
1	rpart()	0.9021	0.0979	0.6071	0.9606	0.3928	0.0394
2	adaboost()	0.8987	0.1013	0.6020	0.9575	0.3980	0.0425
3	randomForest()	0.9097	0.0903	0.6276	0.9656	0.3724	0.0344

Figure 4.5. Model Performance Metrics: Commission Model Comparison

Results: DOR

Models 4, 5 and 6 predict DORs in NROTC using decision tree, adaptive boosting, and random forests algorithms, respectively. Figure 4.6 displays the confusion matrix and model performance metrics (derived from the confusion matrix) for Model 4; Figure 4.7 displays the confusion matrix and model performance metrics for Model 5; Figure 4.8 displays the confusion matrix and model performance metrics for Model 6. Figure 4.9 compares performance metrics for each of the models that predict DORs, with the most desirable results bold and italicized.

<u>Confusion Matrix</u>			
n = 1,185	Actual: No	Actual: Yes	Total
Predicted: NO	1,030	88	1,118
Predicted: YES	26	41	67
Total	1,056	129	1,185
<u>Performance Metrics</u>			
Accuracy	0.9038	Error	0.0962
Sensitivity	0.9754	Specificity	0.3178
Type I Error	0.0246	Type II Error	0.6822

Figure 4.6. Confusion Matrix and Model Performance Metrics: Model 4

<u>Confusion Matrix</u>			
n = 1,185	Actual: No	Actual: Yes	Total
Predicted: NO	1,021	83	1,104
Predicted: YES	35	46	81
Total	1,056	129	1,185
<u>Performance Metrics</u>			
Accuracy	0.9004	Error	0.0996
Sensitivity	0.9669	Specificity	0.3566
Type I Error	0.0331	Type II Error	0.6434

Figure 4.7. Confusion Matrix and Model Performance Metrics: Model 5

<u>Confusion Matrix</u>			
n = 1,185	Actual: No	Actual: Yes	Total
Predicted: NO	1,028	79	1,107
Predicted: YES	28	50	78
Total	1,056	129	1,185
<u>Performance Metrics</u>			
Accuracy	0.9097	Error	0.0903
Sensitivity	0.9735	Specificity	0.3875
Type I Error	0.0265	Type II Error	0.6125

Figure 4.8. Confusion Matrix and Model Performance Metrics: Model 6

Model	Classifier	Accuracy	Error	Sensitivity	Specificity	Type I Error	Type II Error
4	rpart()	0.9038	0.0962	0.9754	0.3178	0.0246	0.6822
5	adaboost()	0.9004	0.0996	0.9669	0.3566	0.0331	0.6434
6	randomForest()	0.9097	0.0903	0.9735	0.3875	0.0265	0.6125

Figure 4.9. Model Performance Metrics: DOR Model Comparison

4.1.2 ROC Curve

As a Bayesian classifier, each of the trained models in this thesis use a posterior probability threshold of 0.5 to assign an observation to a given category. A 0.5 posterior probability threshold is optimal in producing the lowest overall error rate (James et al. 2013). However, it is useful to display model performance when that posterior probability threshold is

varied. Modifying the threshold can affect which type of error (type I or type II) dominates. Selecting a threshold other than 0.5 should be motivated by domain knowledge (James et al. 2013). For example, let's say type I error is associated with much higher cost than type II error—varying the posterior probability threshold, while possibly increasing the overall error rate, could decrease type I error rate, and therefore decrease overall cost.

Receiver Operating Characteristic (ROC) curves graphically display model performance, specifically type I and type II error, as the threshold of posterior probability is varied (James et al. 2013). An example ROC curve is displayed in Figure 4.10. An ideal ROC curve, one representing a perfect classifier, would have an Area Under the Curve (AUC) of 1.0; a no information classifier's ROC curve, represented by the dashed line in Figure 4.10 has an AUC equal to 0.5. In practice, AUC values closer to 1.0 signify stronger classifiers.

Results: Commission

Figure 4.11 displays the ROC curves for Models 1, 2, and 3. Figure 4.12 tabulates the AUC values, with the most desired value bold and italicized. Stronger model performance is associated with higher AUC values (maximum of 1.0). As seen in Figure 4.12, Model 3 performs best using this metric.

Results: DOR

Figure 4.13 displays the ROC curves for Models 4, 5, and 6. Figure 4.14 tabulates the AUC values, with the most desired value bold and italicized. As seen in Figure 4.12, Model 6 performs best using this metric.

4.1.3 Conclusion: Model Performance

We assessed model performance using confusion matrices and ROC curves. Both performance metrics established Model 3 as the superior model for predicting NROTC commissions; both performance metrics confirmed Model 6 as the optimal NROTC DOR-predicting model. Model 3 and Model 6 used the Random Forests (RF) classifier.

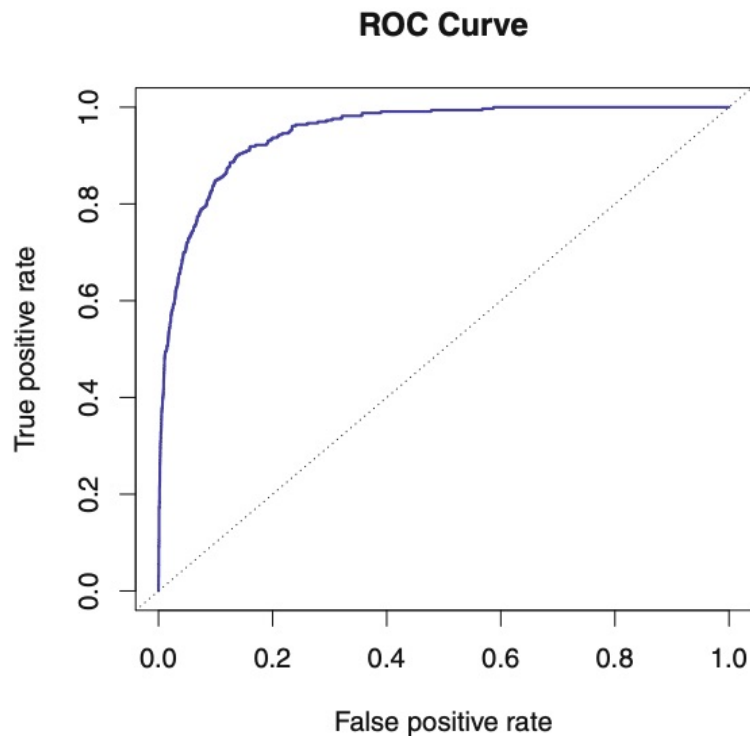


Figure 4.10. Example ROC Curve Source: James (2013). Although the posterior probability thresholds are not shown, a ROC curve demonstrates the relationship between type I and type II error for a classifier under different thresholds. A curve representing a perfect classifier appears as a vertical line hugging the y-axis connected with a horizontal line at the very top of the plot. The dashed line is a curve for a no information classifier—one that is no more helpful than randomly assigning observations into distinct categories.

4.2 Model Inference

The following section summarizes the practical inferences gleaned from modeling. The level and method of interpretation for each classifier type is different—each proceeding subsection describes these differences and then illustrates the information about NROTC attrition derived from each classifier.

4.2.1 Decision Tree

From the decision tree classifier, information about the relationship between data inputs and outputs is represented graphically in a tree-format (James et al. 2013). Decision trees visually

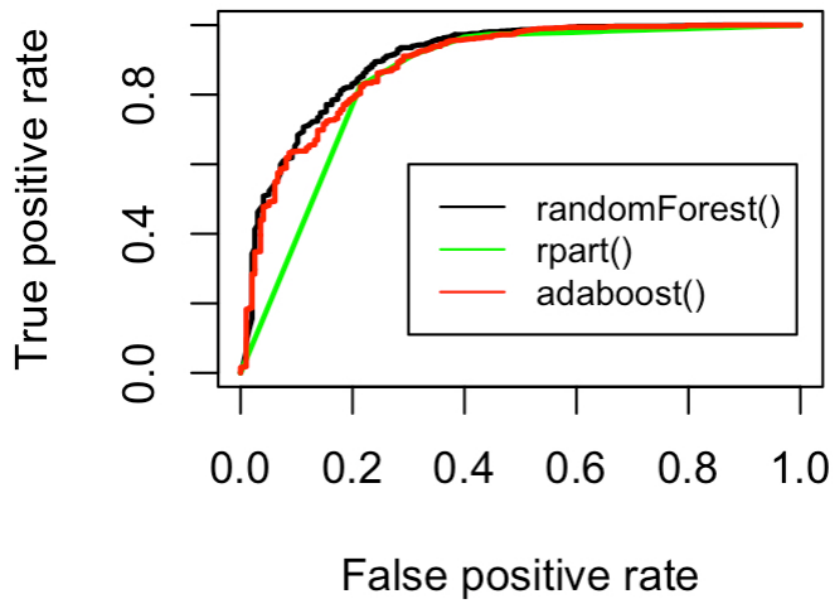


Figure 4.11. ROC Curves: Commission Models

Model	Classifier	AUC
1	rpart()	0.8474
2	adaboost()	0.8871
3	randomForest()	0.9011

Figure 4.12. AUC: Commission Models

indicate which data inputs are most important to predicting the response variable. The first split, or node, at the very top of the tree represents the most significant input variable for classification; each subsequent node signifies the most important input for classification for the specific subgroup that the algorithm is evaluating. Beyond variable importance, decision trees visually illustrate the classification process—using the input values for a given data observation, one could follow a path down the tree, ultimately ending with classification. This can be done with every data observation—greatly enhancing how interested parties can interpret classification results and draw practical conclusions.

Results: Commission

Figure 4.15 portrays the classification process of Model 1.

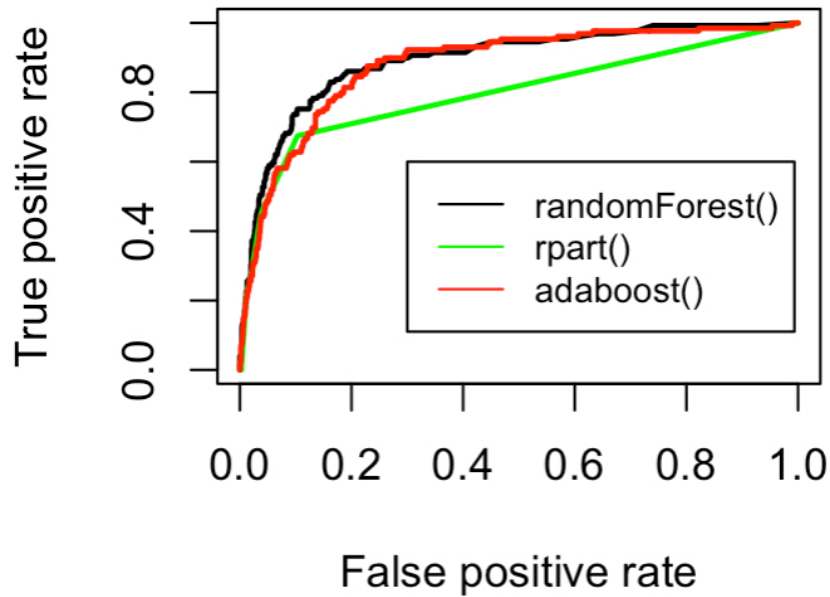


Figure 4.13. ROC Curves: DOR Models

Model	Classifier	AUC
4	rpart()	0.7970
5	adaboost()	0.8796
6	randomForest()	0.8925

Figure 4.14. AUC: DOR Models

The tree clearly illustrates each branch of the tree: the variable, its cutoff value, and the proportion of observations following each branch. Figure 4.15 shows the most critical variable to predicting NROTC commissions to be *PROGRAM.TYPE*—a categorical variable with midshipmen deemed as either scholarship or college program midshipmen. Below this categorization, the decision tree branched off other important factors: aptitude, GPA, PFA score, and option (Navy/Marine).

Results: DOR

Figure 4.16 portrays the classification process of Model 4. It shows the most critical variable to predicting DORs in NROTC to be *PROGRAM.TYPE*—a categorical variable with midshipmen deemed as either scholarship or college program midshipmen. Below this cat-

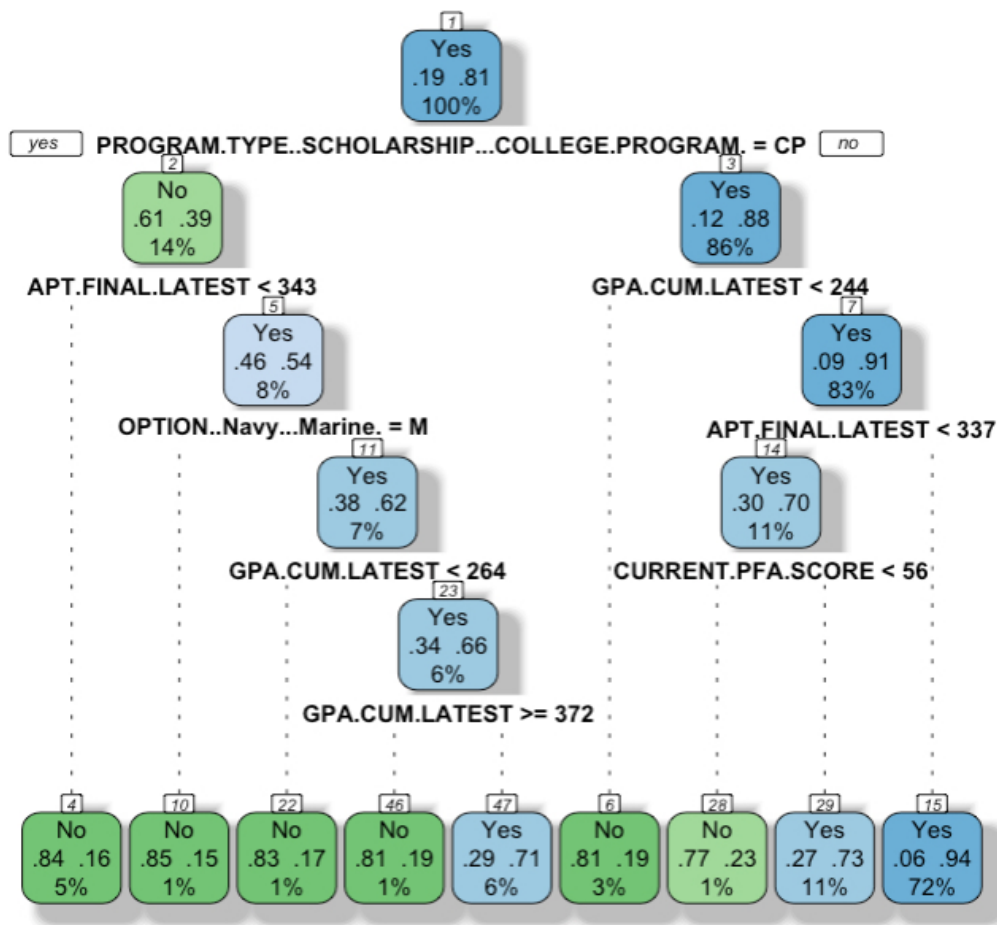


Figure 4.15. Decision Tree: Commission (Model 1)

egorization, the decision tree branched off GPA, indicating this metric as the second most important factor to predicting whether a midshipman will DOR or not. Further down the tree, classification decisions branch off aptitude and option (Navy/Marine).

4.2.2 RF: Variable Importance

Since adaptive boosting and random forests algorithms are aggregations of many decision trees, their classification process is not as easy to interpret or graphically display (James et al. 2013). However, the structure of the random forests classifier enables self-evaluation and determination of variable importance. The *importance()* function calculates two measures of variable importance: mean decrease of accuracy and mean decrease gini index. "The first

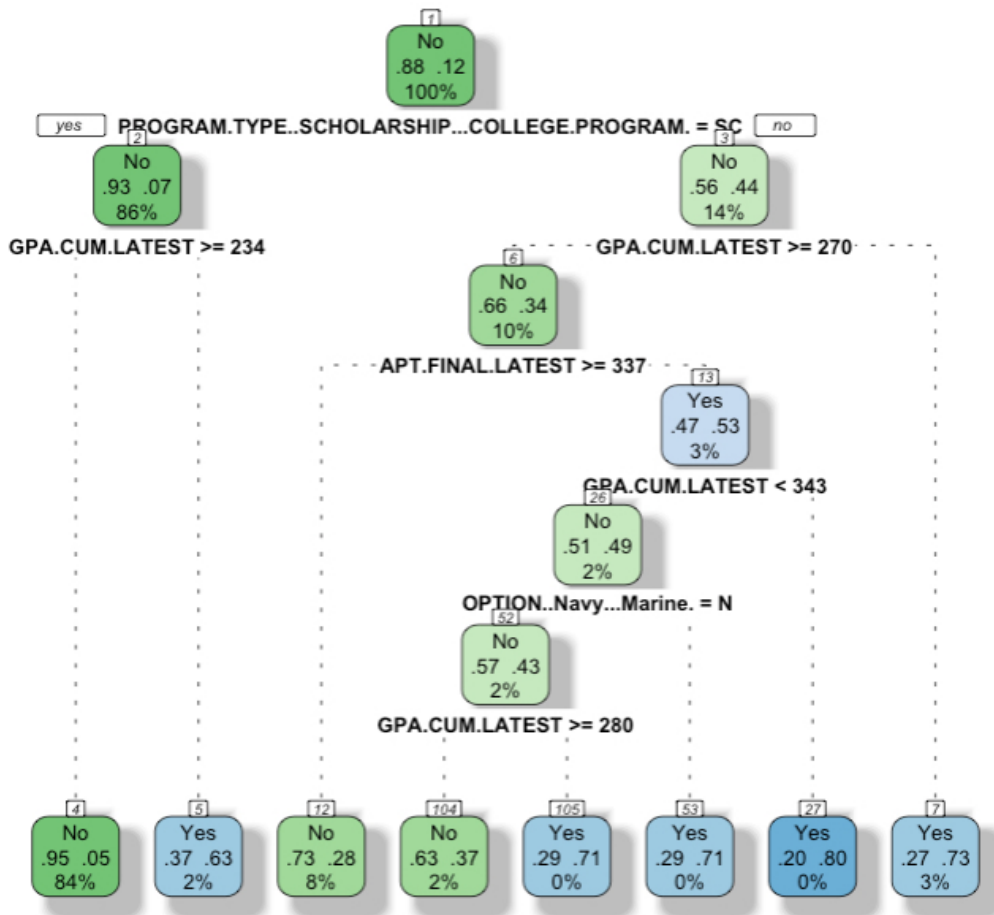


Figure 4.16. Decision Tree: DOR (Model 4)

measure is computed from permuting OOB data: For each tree, the prediction error on the out-of-bag portion of the data is recorded (error rate for classification, MSE for regression). Then the same is done after permuting each predictor variable. The difference between the two are then averaged over all trees, and normalized by the standard deviation of the differences. The second measure is the total decrease in node impurities from splitting on the variable, averaged over all trees. For classification, the node impurity is measured by the Gini index" (RStudio Team 2020). Variable importance plots in this thesis will display mean decrease of accuracy.

Results: Commission

Figure 4.17 shows variable importance for Model 3.

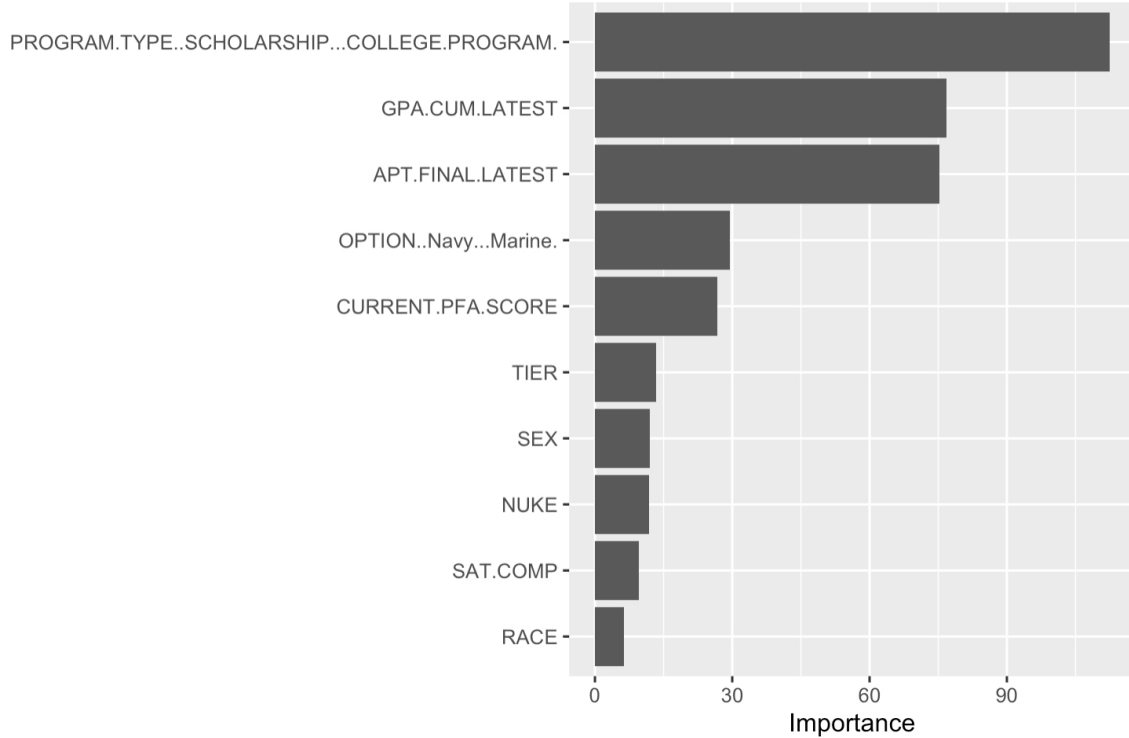


Figure 4.17. RF Variable Importance (Mean Decrease of Accuracy): Commission (Model 3)

Figure 4.17 plots program type (scholarship/college program) as the most important factor to predicting NROTC commissions, followed by GPA, aptitude, option (Navy/Marine), and PFA score.

Results: DOR

Figure 4.18 shows variable importance for Model 3. It plots program type (scholarship/college program) as the most important factor to predicting NROTC commissions, followed by GPA, aptitude, option (Navy/Marine), and PFA score.

4.2.3 Conclusion: Model Inference

Both methods of model inference, decision tree and variable importance, indicate that the most influential factors in predicting NROTC commissions and DORs are program type

(scholarship/college program), GPA, aptitude score, PFA score, and option (Navy/Marine). The last section will answer the research questions based on these results.

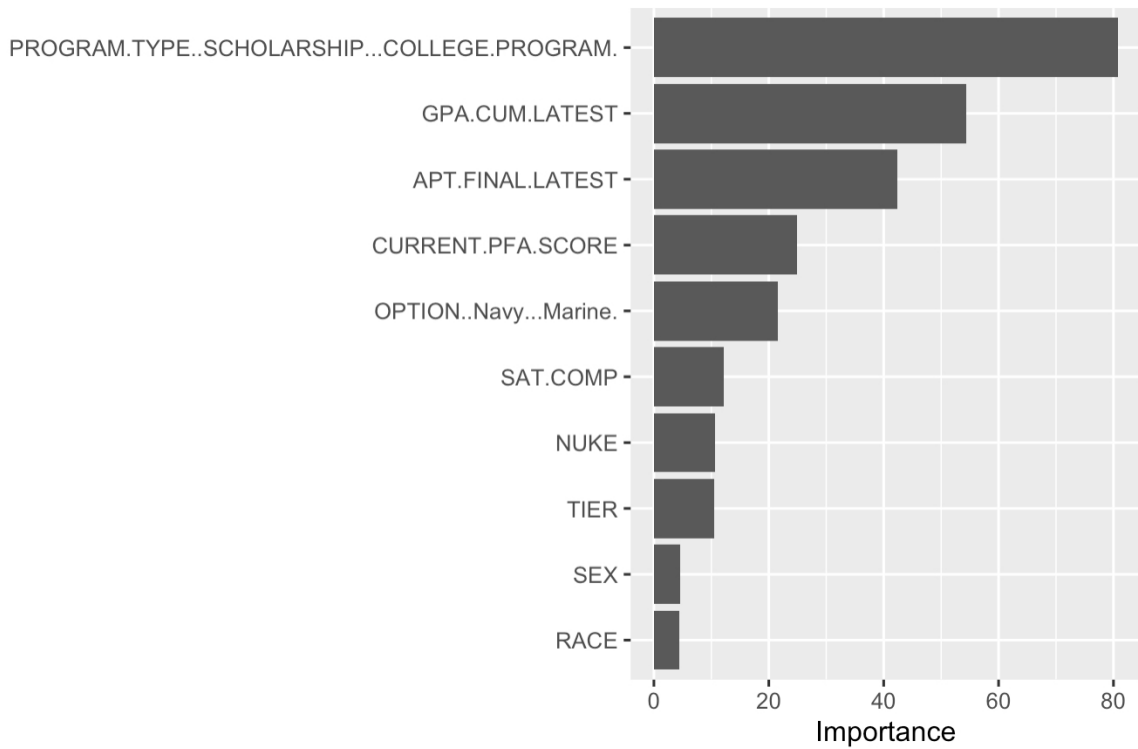


Figure 4.18. RF Variable Importance (Mean Decrease of Accuracy): DOR (Model 6)

4.3 Research Questions

This section directly answers each research question posed by this thesis based on the trained models, their performance, and inference. Considering it is a commissioning requirement to be on some type of Navy-funded scholarship, program type (scholarship/college program) is not relevant for research question 1, and therefore not addressed.

4.3.1 Research Question 1

What demographic, academic, and performance-based metrics are important to predicting whether a NROTC midshipmen will commission? Predictive modeling and analysis of NROTC attrition concludes that academic and performance-based metrics, namely

GPA and aptitude score, are the strongest predictors of whether NROTC midshipmen will commission. PFA score is another important performance-based factor, but is not as critical to the model as GPA or aptitude. Statistical analysis demonstrates that academic and performance-based metrics are far more important to predicting commissions than demographic data such as sex, race, or ethnicity.

4.3.2 Research Question 2

What demographic, academic, and performance-based metrics are important to predicting whether a NROTC midshipmen will DOR? Modeling and analysis indicates that factors important to predicting commissions are also useful to predicting DORs. First, program type (scholarship/college program) is the most important factor for predicting DORs, indicating that midshipmen are more motivated to DOR if they are not on scholarship. GPA and aptitude score are the strongest predictors of DORs. PFA score is also important, and seems to be more important for the DOR models than commission models. Academic and performance-based metrics are far more important to predicting DORs than demographic data such as sex, race, or ethnicity.

4.3.3 Research Question 3

Is eligibility to serve as a Nuclear Submarine Officer an important factor in predicting whether a NROTC midshipmen will commission or DOR? While Nuke eligibility is highly correlated with commissioning and overall high performance in NROTC, it is not an important factor in predicting commissions or DORs (as compared to other data input categories discussed).

CHAPTER 5: Conclusion

This chapter summarizes the thesis and provides recommendations for future research around attrition in the military.

5.1 Summary

This thesis models NROTC attrition from 2013 to 2020 using demographic, academic, and performance-based data. Trained on midshipmen data provided by NETC, classification models predict NROTC commissions and DORs. Model performance metrics, such as classification accuracy and ROC curves, demonstrate which models are most effective in predicting their respective response variable. Model inference strategies such as decision trees and variable importance values indicate which demographic, academic, and performance-based data are critical to model performance.

Research results and analysis indicate RF models to be the most effective classifier for both NROTC commissions and DORs. The most important predictors for the RF classification models are scholarship status, aptitude score, GPA, and PFA score. Academic and performance-based data prove far more effective as predictors—for both commissions and DORs—than demographic data. Lastly, eligibility to serve as a Nuclear Submarine Officer is not an important predictor for the trained models, despite it being highly correlated with commissioning.

5.2 Recommendations for Future Research

We focused on relationships between demographic, academic, and performance-based data and attrition solely in NROTC—there is vast potential for attrition-based research across all sects of the military. Attrition is a key part of all organizations, especially the military. It is important to understand which factors influence attrition so organizational leadership can optimize its purpose. This research should motivate organizations to collect meaningful data that can inform leadership on what factors lead to attrition. Furthermore, the types of classifiers used in this research were limited by the majority of predictors being categorical

variables—with more expansive data sets and varied descriptors, models could be trained with other classifiers such as k-nearest neighbors and support vector machines.

List of References

- Chatterjee S (2016) *fastAdaboost: a Fast Implementation of Adaboost*. URL <https://CRAN.R-project.org/package=fastAdaboost>, r package version 1.0.0.
- Department of Defense (2019) Regulations for officer development. Technical Report NSTC M-1533.2D, Great Lakes, IL, <https://www.netc.navy.mil/Portals/46/NSTC/cmd-docs/manuals/NSTC/20M-1533.2D/20Regulations/20for/20Officer/20Development/20v18/20Final.pdf>.
- James G, Witten D, Hastie T, Tibshirani R (2013) *An Introduction to Statistical Learning with Applications in R* (Springer, New York).
- Liaw A, Wiener M (2002) Classification and regression by randomforest. URL <https://CRAN.R-project.org/doc/Rnews/>.
- Microsoft Organization (2019) Microsoft Excel. <https://office.microsoft.com/excel>.
- Naval Education and Training Command (2020) Naval reserve officer training corps: About. Accessed November 5, 2020, <https://www.netc.navy.mil/Commands/Naval-Service-Training-Command/NROTC/About/>.
- Office CB (1990) Officer commissioning programs: Costs and officer performance. Technical report, Congressional Budget Office, Washington, DC, <https://www.cbo.gov/sites/default/files/101st-congress-1989-1990/reports/90-cbo-028.pdf>.
- RStudio Team (2020) RStudio: Integrated Development Environment for R. <https://rstudio.com/>.
- SAS Institute (2017) Statistical Analysis System. http://www.sas.com/en_us/home.html.
- Therneau T, Atkinson B (2019) *rpart: Recursive Partitioning and Regression Trees*. URL <https://CRAN.R-project.org/package=rpart>, r package version 4.1-15.

THIS PAGE INTENTIONALLY LEFT BLANK

Initial Distribution List

1. Defense Technical Information Center
Ft. Belvoir, Virginia
2. Dudley Knox Library
Naval Postgraduate School
Monterey, California