

# Data Mining for Student Success and Perseverance

*Sameer Bhatnagar*

*Jonathan Guillemette*

*Micheal Dugdale*

*Sahir Bhatnagar*

*Nathaniel Lasry*

*2017-05-22*



# Contents

<b>1</b>	<b>Preface</b>	<b>5</b>
<b>2</b>	<b>Introduction</b>	<b>7</b>
<b>3</b>	<b>Literature</b>	<b>9</b>
<b>4</b>	<b>Descriptive Statistics</b>	<b>11</b>
<b>5</b>	<b>Methods Centered on Determining Predictive Factors</b>	<b>13</b>
<b>6</b>	<b>Methods centered on predicting at-risk students</b>	<b>15</b>
6.1	Example one . . . . .	15
6.2	Example two . . . . .	15
<b>7</b>	<b>Comparisons</b>	<b>17</b>
<b>8</b>	<b>Final Words</b>	<b>19</b>



# Chapter 1

## Preface

This report summarizes the work done by our team on using college registration records at three different anglophone CEGEPS in Montreal in order to find predictors of attrition.



## Chapter 2

# Introduction





## Chapter 3

# Literature

The most important relevant work for this project is (Jorgensen et al., 2009).



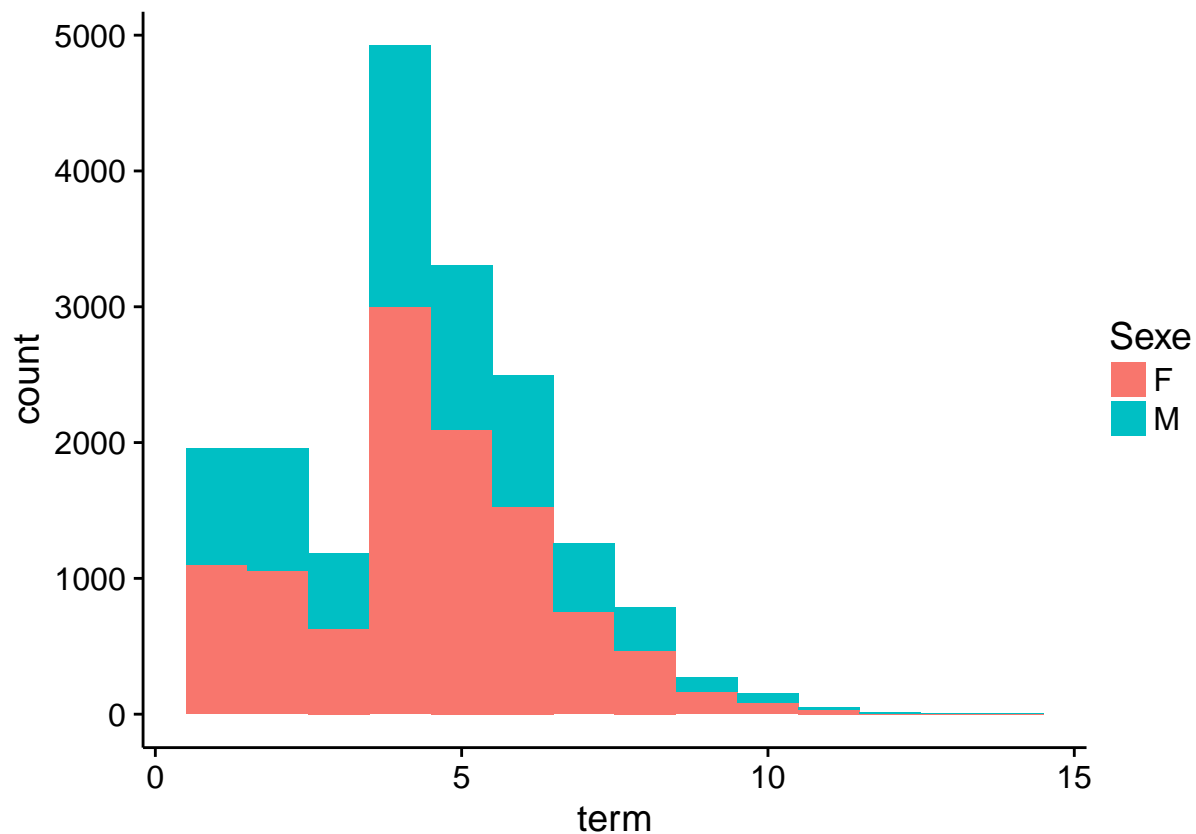
## Chapter 4

# Descriptive Statistics

Here in we will describe - the data set - the methods by which we label students at risk - the distributions of at-risk students by - demographic indicators - registration record indicators

For example, how many what is the distribution of terms spent by a student at Dawson College, disaggregated by gender?

```
library(ggplot2)
library(cowplot)
library(data.table)
load('bin/data/labelled_students.Rdata')
past_students=students_last_session[status!='current']
ggplot(data = past_students,aes(term,fill=Sexe))+geom_histogram(binwidth = 1)
```





## Chapter 5

# Methods Centered on Determining Predictive Factors

This chapter will be focused on methods which have a sound probabilistic framework, and allow for inference into the statistical importance of predictive factors.

- Logistic Regression
- Mixed Effects Models



## Chapter 6

# Methods centered on predicting at-risk students

This chapter will explore the use of machine learning algorithms whose primary focus is prediction. This comes at the expense of model interpretability, and this tradeoff will be discussed herein as well.

- Decision Trees and Random Forests
- Neural Networks

### 6.1 Example one

### 6.2 Example two





## Chapter 7

# Comparisons

This chapter will compare the effectiveness of the methods developed in the previous two chapters, and compare them to more basic approaches to identifying students at-risk.



## Chapter 8

# Final Words

We have finished a nice book.



# Bibliography

Jorgensen, S., Fichten, C., and Havel, A. (2009). *Predicting the At Risk Status of College Students: Males and Students with Disabilities. Final Report Presented to PAREA, Spring 2009.*