

GradX – Graduation Success Predictor

Prathiksha Rumale Vishwanath

Load the Data

```
data <- read_csv("data.csv", col_names = TRUE, col_types = cols())
head(data)
```

```
## # A tibble: 6 x 1
##   Marital status;Application mode;Application order;Course;"Daytime/evening at~1
##   <chr>
## 1 1;17;5;171;1;1;122.0;1;19;12;5;9;127.3;1;0;0;1;1;0;20;0;0;0;0;0.0;0;0;0;0;0~
## 2 1;15;1;9254;1;1;160.0;1;1;3;3;3;142.5;1;0;0;0;1;0;19;0;0;6;6;6;14.0;0;0;6;6;6~
## 3 1;1;5;9070;1;1;122.0;1;37;37;9;9;124.8;1;0;0;0;1;0;19;0;0;6;0;0;0.0;0;0;6;0;0~
## 4 1;17;2;9773;1;1;122.0;1;38;37;5;3;119.6;1;0;0;1;0;0;20;0;0;6;8;6;13.428571428~
## 5 2;39;1;8014;0;1;100.0;1;37;38;9;9;141.5;0;0;0;1;0;0;45;0;0;6;9;5;12.333333333~
## 6 2;39;1;9991;0;19;133.1;1;37;37;9;7;114.8;0;0;1;1;1;0;50;0;0;5;10;5;11.8571428~
## # i abbreviated name:
## #   1: 'Marital status;Application mode;Application order;Course;"Daytime/evening attendance\t";Prev.
```

Data Cleaning

```
# Checking for missing values
missing_values <- colSums(is.na(data))
missing_values[missing_values > 0]
```

```
## named numeric(0)
```

```
# Reload with proper delimiter if needed
data <- read_csv("data.csv", sep = ";", header = TRUE)
head(data)
```

```
##   Marital.status Application.mode Application.order Course
## 1             1             17             5      171
## 2             1             15             1     9254
## 3             1              1             5     9070
## 4             1             17             2     9773
## 5             2             39             1     8014
## 6             2             39             1     9991
##   Daytime.evening.attendance. Previous.qualification
## 1                      1                      1
## 2                      1                      1
```

## 3	1	1		
## 4	1	1		
## 5	0	1		
## 6	0	19		
##	Previous.qualification..grade.	Nacionality	Mother.s.qualification	
## 1	122.0	1	19	
## 2	160.0	1	1	
## 3	122.0	1	37	
## 4	122.0	1	38	
## 5	100.0	1	37	
## 6	133.1	1	37	
##	Father.s.qualification	Mother.s.occupation	Father.s.occupation	
## 1	12	5	9	
## 2	3	3	3	
## 3	37	9	9	
## 4	37	5	3	
## 5	38	9	9	
## 6	37	9	7	
##	Admission.grade	Displaced	Educational.special.needs	Debtor
## 1	127.3	1	0	0
## 2	142.5	1	0	0
## 3	124.8	1	0	0
## 4	119.6	1	0	0
## 5	141.5	0	0	0
## 6	114.8	0	0	1
##	Tuition.fees.up.to.date	Gender	Scholarship.holder	Age.at.enrollment
## 1	1	1	0	20
## 2	0	1	0	19
## 3	0	1	0	19
## 4	1	0	0	20
## 5	1	0	0	45
## 6	1	1	0	50
##	International	Curricular.units.1st.sem..credited.		
## 1	0	0		
## 2	0	0		
## 3	0	0		
## 4	0	0		
## 5	0	0		
## 6	0	0		
##	Curricular.units.1st.sem..enrolled.	Curricular.units.1st.sem..evaluations.		
## 1	0	0		
## 2	6	6		
## 3	6	0		
## 4	6	8		
## 5	6	9		
## 6	5	10		
##	Curricular.units.1st.sem..approved.	Curricular.units.1st.sem..grade.		
## 1	0	0.00000		
## 2	6	14.00000		
## 3	0	0.00000		
## 4	6	13.42857		
## 5	5	12.33333		
## 6	5	11.85714		
##	Curricular.units.1st.sem..without.evaluations.			

```

## 1 0
## 2 0
## 3 0
## 4 0
## 5 0
## 6 0
## Curricular.units.2nd.sem..credited. Curricular.units.2nd.sem..enrolled.
## 1 0 0
## 2 0 6
## 3 0 6
## 4 0 6
## 5 0 6
## 6 0 5
## Curricular.units.2nd.sem..evaluations. Curricular.units.2nd.sem..approved.
## 1 0 0
## 2 6 6
## 3 0 0
## 4 10 5
## 5 6 6
## 6 17 5
## Curricular.units.2nd.sem..grade.
## 1 0.00000
## 2 13.66667
## 3 0.00000
## 4 12.40000
## 5 13.00000
## 6 11.50000
## Curricular.units.2nd.sem..without.evaluations. Unemployment.rate
## 1 0 10.8
## 2 0 13.9
## 3 0 10.8
## 4 0 9.4
## 5 0 13.9
## 6 5 16.2
## Inflation.rate GDP Target
## 1 1.4 1.74 Dropout
## 2 -0.3 0.79 Graduate
## 3 1.4 1.74 Dropout
## 4 -0.8 -3.12 Graduate
## 5 -0.3 0.79 Graduate
## 6 0.3 -0.92 Graduate

```

Student Enrollment by Course

```

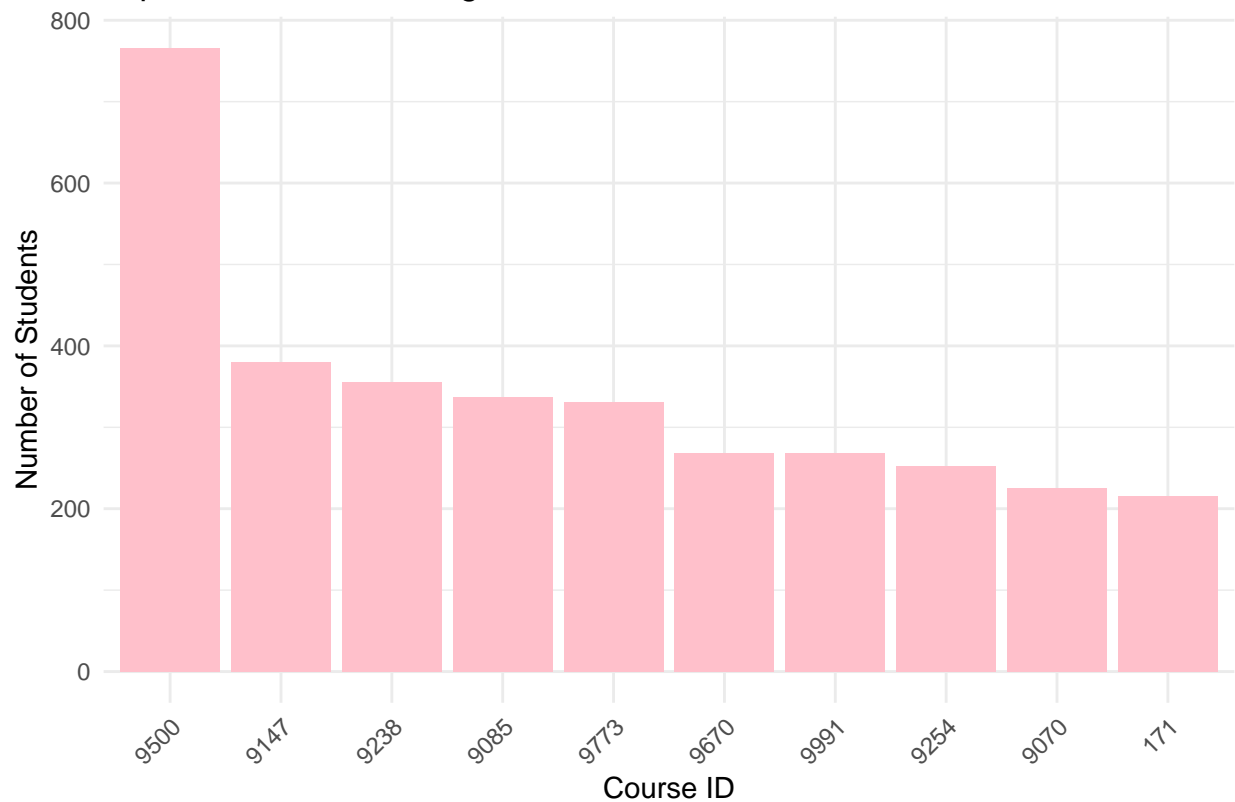
course_enrollment <- data %>%
  group_by(Course) %>%
  summarise(Number_of_Students = n()) %>%
  arrange(desc(Number_of_Students))

# Display summary
top_courses <- head(course_enrollment, 10)
print(top_courses)

```

```
## # A tibble: 10 x 2
##   Course Number_of_Students
##   <int>           <int>
## 1  9500             766
## 2  9147             380
## 3  9238             355
## 4  9085             337
## 5  9773             331
## 6  9670             268
## 7  9991             268
## 8  9254             252
## 9  9070             226
## 10 171             215
```

Top 10 Courses with Highest Student Enrollment



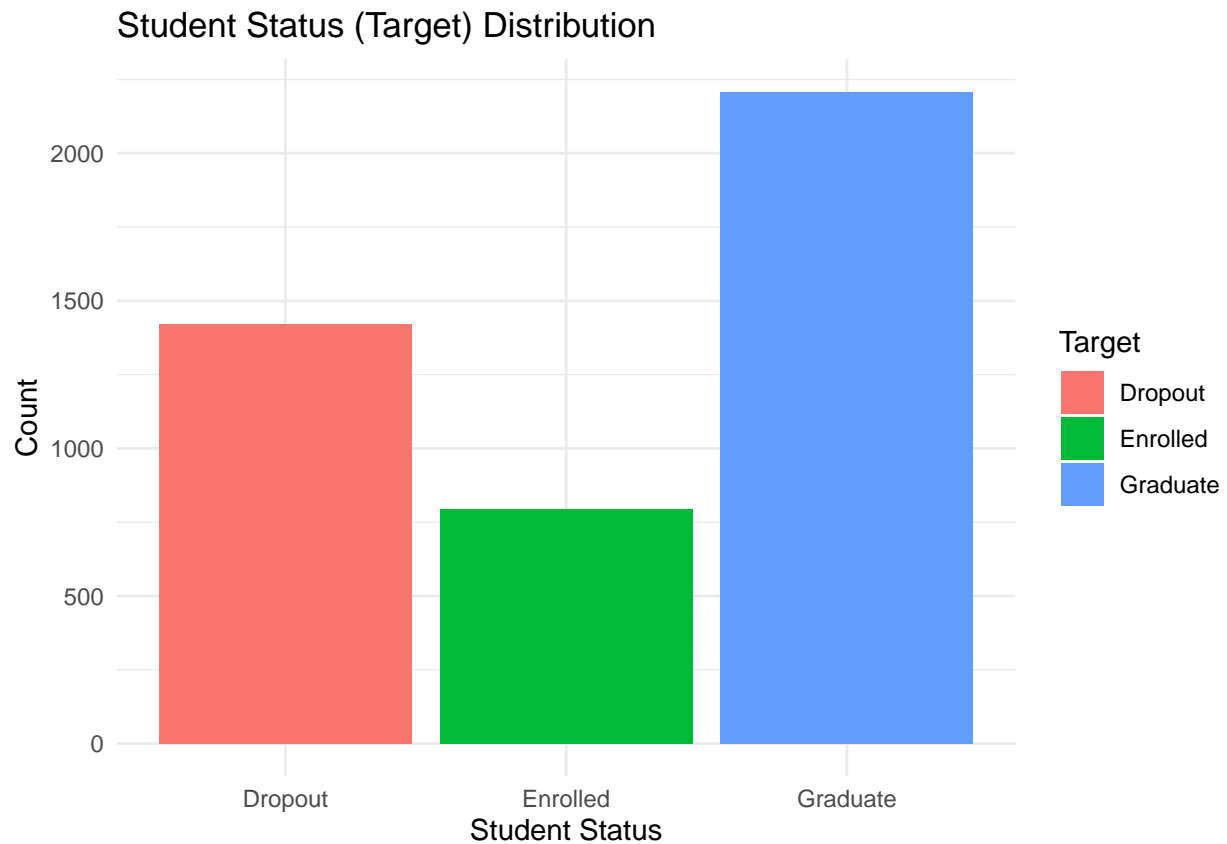
Dropout vs Graduation Rate

```
dropout_distribution <- data %>%
  group_by(Target) %>%
  summarise(Count = n())

# Display summary
dropout_distribution
```

```
## # A tibble: 3 x 2
```

```
## Target Count
## <chr> <int>
## 1 Dropout 1421
## 2 Enrolled 794
## 3 Graduate 2209
```



Statistical Analysis: Predicting Dropout

```
colnames(data)
```

```
## [1] "Marital.status"
## [2] "Application.mode"
## [3] "Application.order"
## [4] "Course"
## [5] "Daytime.evening.attendance."
## [6] "Previous.qualification"
## [7] "Previous.qualification..grade."
## [8] "Nacionality"
## [9] "Mother.s.qualification"
## [10] "Father.s.qualification"
## [11] "Mother.s.occupation"
## [12] "Father.s.occupation"
## [13] "Admission.grade"
```

```

## [14] "Displaced"
## [15] "Educational.special.needs"
## [16] "Debtor"
## [17] "Tuition.fees.up.to.date"
## [18] "Gender"
## [19] "Scholarship.holder"
## [20] "Age.at.enrollment"
## [21] "International"
## [22] "Curricular.units.1st.sem..credited."
## [23] "Curricular.units.1st.sem..enrolled."
## [24] "Curricular.units.1st.sem..evaluations."
## [25] "Curricular.units.1st.sem..approved."
## [26] "Curricular.units.1st.sem..grade."
## [27] "Curricular.units.1st.sem..without.evaluations."
## [28] "Curricular.units.2nd.sem..credited."
## [29] "Curricular.units.2nd.sem..enrolled."
## [30] "Curricular.units.2nd.sem..evaluations."
## [31] "Curricular.units.2nd.sem..approved."
## [32] "Curricular.units.2nd.sem..grade."
## [33] "Curricular.units.2nd.sem..without.evaluations."
## [34] "Unemployment.rate"
## [35] "Inflation.rate"
## [36] "GDP"
## [37] "Target"

data$Target <- factor(data$Target, levels = c("Dropout", "Graduate")) # Ensure correct ordering
model <- glm(Target ~ `Age.at.enrollment` + `Scholarship.holder` + `Gender` + `Curricular.units.1st.sem
              data = data, family = binomial)
summary(model)

##
## Call:
## glm(formula = Target ~ Age.at.enrollment + Scholarship.holder +
##      Gender + Curricular.units.1st.sem..grade. + Curricular.units.2nd.sem..grade.,
##      family = binomial, data = data)
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)      -1.469306   0.205883  -7.137 9.57e-13 ***
## Age.at.enrollment -0.045851   0.005632  -8.141 3.91e-16 ***
## Scholarship.holder  1.390197   0.126558  10.985 < 2e-16 ***
## Gender            -0.566619   0.093834  -6.039 1.56e-09 ***
## Curricular.units.1st.sem..grade.  0.003211   0.020066   0.160  0.873
## Curricular.units.2nd.sem..grade.  0.279223   0.019041  14.664 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 4859.8 on 3629 degrees of freedom
## Residual deviance: 3083.8 on 3624 degrees of freedom
## (794 observations deleted due to missingness)
## AIC: 3095.8
##

```

```
## Number of Fisher Scoring iterations: 5
```

```
vif(model)
```

```
##              Age.at.enrollment      Scholarship.holder
##              1.023303              1.029744
##              Gender Curricular.units.1st.sem..grade.
##              1.020658              2.837992
## Curricular.units.2nd.sem..grade.
##              2.848617
```

```
set.seed(42)
```

```
trainIndex <- createDataPartition(data$Target, p = 0.8, list = FALSE)
```

```
trainData <- data[trainIndex, ]
```

```
testData <- data[-trainIndex, ]
```

```
model <- glm(Target ~ Age.at.enrollment + Scholarship.holder + `Gender` + `Curricular.units.1st.sem..grade`,
             data = trainData, family = binomial)
```

```
predictions <- predict(model, newdata = testData, type = "response")
```

```
predicted_classes <- ifelse(predictions >= 0.5, "Graduate", "Dropout")
```

```
confusionMatrix(factor(predicted_classes), testData$Target)
```

```
## Confusion Matrix and Statistics
```

```
##
```

```
##           Reference
```

```
## Prediction Dropout Graduate
```

```
## Dropout      177      32
```

```
## Graduate     107     409
```

```
##
```

```
##           Accuracy : 0.8083
```

```
##           95% CI : (0.7777, 0.8363)
```

```
## No Information Rate : 0.6083
```

```
## P-Value [Acc > NIR] : < 2.2e-16
```

```
##
```

```
##           Kappa : 0.5778
```

```
##
```

```
## McNemar's Test P-Value : 3.461e-10
```

```
##
```

```
##           Sensitivity : 0.6232
```

```
##           Specificity : 0.9274
```

```
## Pos Pred Value : 0.8469
```

```
## Neg Pred Value : 0.7926
```

```
## Prevalence : 0.3917
```

```
## Detection Rate : 0.2441
```

```
## Detection Prevalence : 0.2883
```

```
## Balanced Accuracy : 0.7753
```

```
##
```

```
## 'Positive' Class : Dropout
```

```
##
```

```

# Define new student data for prediction
new_data <- data.frame(
  Age.at.enrollment = c(21, 25, 30),
  Scholarship.holder = c(1, 0, 1), # 1 = Has Scholarship, 0 = No Scholarship
  Gender = c(1, 0, 1), # 1 = Male, 0 = Female
  Curricular.units.1st.sem..grade. = c(14, 10, 12),
  Curricular.units.2nd.sem..grade. = c(15, 9, 11)
)

# Predict dropout probability
predictions <- predict(model, newdata = new_data, type = "response")

# Convert probabilities into class labels
predicted_classes <- ifelse(predictions >= 0.75, "Graduate", "Dropout")

# Display results
prediction_results <- data.frame(new_data, Predicted_Status = predicted_classes, Probability = predictions)
print(prediction_results)

```

```

##   Age.at.enrollment Scholarship.holder Gender Curricular.units.1st.sem..grade.
## 1                21                1      1                                14
## 2                25                0      0                                10
## 3                30                1      1                                12
##   Curricular.units.2nd.sem..grade. Predicted_Status Probability
## 1                      15      Graduate    0.9404663
## 2                      9      Dropout    0.4694791
## 3                     11      Graduate    0.7689142

```

Conclusion

This analysis provides insights into **student enrollment trends and dropout rates**, helping institutions understand key factors influencing student retention.