

Final Project

Ling Jin

2024-12-01

Introduction

This project analyzes the “Higher Education Students Performance Evaluation” dataset to identify factors affecting academic performance. Various regression and classification methods will be applied. The main questions of interest are:

1. What are the strongest predictors of cumulative GPA?
2. How do socio-economic factors influence student performance?
3. How do educational habits correlate with academic outcomes?

```
# Load Libraries
library(tidyverse)
```

```
-- Attaching core tidyverse packages ----- tidyverse 2.0.0 --
v dplyr      1.1.4      v readr      2.1.5
v forcats    1.0.0      v stringr    1.5.1
v ggplot2    3.5.1      v tibble     3.2.1
v lubridate  1.9.3      v tidyr      1.3.1
v purrr      1.0.2
-- Conflicts ----- tidyverse_conflicts() --
x dplyr::filter() masks stats::filter()
x dplyr::lag()     masks stats::lag()
i Use the conflicted package (<http://conflicted.r-lib.org/>) to force all conflicts to become errors
library(caret)
```

Loading required package: lattice

Attaching package: 'caret'

The following object is masked from 'package:purrr':

lift

```
library(randomForest)
```

randomForest 4.7-1.2

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

Attaching package: 'randomForest'

The following object is masked from 'package:dplyr':

```
combine
```

The following object is masked from 'package:ggplot2':

```
margin
```

```
library(ggplot2)
library(corrplot)
```

```
corrplot 0.95 loaded
```

```
library(knitr)
library(pROC)
```

Type 'citation("pROC")' for a citation.

Attaching package: 'pROC'

The following objects are masked from 'package:stats':

```
cov, smooth, var
```

Load Data

```
# Load the dataset
data <- read.csv("Higher Education Students Performance Evaluation.csv")

# Remove unique identifier STUDENT.ID
data <- data %>% select(-STUDENT.ID)

# View the structure
str(data)
```

```
'data.frame':  145 obs. of  32 variables:
 $ X1      : int  2 2 2 1 2 2 1 1 2 2 ...
 $ X2      : int  2 2 2 1 2 2 2 1 1 1 ...
 $ X3      : int  3 3 2 1 1 2 2 2 3 2 ...
 $ X4      : int  3 3 3 3 3 3 4 3 3 3 ...
 $ X5      : int  1 1 2 1 2 2 2 1 2 2 ...
 $ X6      : int  2 2 2 2 2 2 2 1 1 2 ...
 $ X7      : int  2 2 2 1 1 2 2 1 1 1 ...
 $ X8      : int  1 1 2 2 3 2 1 2 1 3 ...
 $ X9      : int  1 1 4 1 1 1 1 2 1 4 ...
 $ X10     : int  1 1 2 2 4 1 3 3 3 2 ...
 $ X11     : int  1 2 2 1 3 3 1 4 2 1 ...
 $ X12     : int  2 3 2 2 3 3 3 3 4 2 ...
 $ X13     : int  3 2 2 5 2 2 1 1 2 3 ...
 $ X14     : int  1 1 1 1 1 1 1 1 1 1 ...
 $ X15     : int  2 2 2 2 2 2 2 4 2 2 ...
 $ X16     : int  5 1 1 1 4 3 4 3 4 3 ...
 $ X17     : int  3 2 2 3 2 1 2 1 1 2 ...
 $ X18     : int  2 2 1 1 1 1 2 2 2 2 ...
 $ X19     : int  2 2 2 2 1 2 2 2 2 2 ...
 $ X20     : int  1 1 1 1 1 1 2 1 1 1 ...
```

```

$ X21      : int  1 1 1 1 1 1 1 1 1 1 ...
$ X22      : int  1 1 1 1 1 1 2 1 1 2 ...
$ X23      : int  1 1 1 1 2 1 1 3 1 1 ...
$ X24      : int  1 1 1 2 1 1 1 1 1 1 ...
$ X25      : int  3 3 2 3 2 1 3 3 3 2 ...
$ X26      : int  2 2 2 2 2 2 3 2 2 2 ...
$ X27      : int  1 3 1 2 2 1 3 2 2 2 ...
$ X28      : int  2 2 1 1 1 2 3 1 2 2 ...
$ X29      : int  1 2 2 3 2 4 4 1 4 1 ...
$ X30      : int  1 3 2 2 2 4 4 1 3 2 ...
$ COURSE.ID: int  1 1 1 1 1 1 1 1 1 1 ...
$ GRADE     : int  1 1 1 1 1 2 5 2 5 0 ...

```

```
summary(data)
```

X1		X2		X3		X4		X5	
Min.	:1.000	Min.	:1.0	Min.	:1.000	Min.	:1.000	Min.	:1.000
1st Qu.:	:1.000	1st Qu.:	:1.0	1st Qu.:	:2.000	1st Qu.:	:3.000	1st Qu.:	:1.000
Median :	:2.000	Median :	:2.0	Median :	:2.000	Median :	:3.000	Median :	:2.000
Mean :	:1.621	Mean :	:1.6	Mean :	:1.945	Mean :	:3.572	Mean :	:1.662
3rd Qu.:	:2.000	3rd Qu.:	:2.0	3rd Qu.:	:2.000	3rd Qu.:	:4.000	3rd Qu.:	:2.000
Max.	:3.000	Max.	:2.0	Max.	:3.000	Max.	:5.000	Max.	:2.000
X6		X7		X8		X9		X10	
Min.	:1.0	Min.	:1.000	Min.	:1.000	Min.	:1.000	Min.	:1.000
1st Qu.:	:1.0	1st Qu.:	:1.000	1st Qu.:	:1.000	1st Qu.:	:1.000	1st Qu.:	:1.000
Median :	:2.0	Median :	:2.000	Median :	:1.000	Median :	:1.000	Median :	:2.000
Mean :	:1.6	Mean :	:1.579	Mean :	:1.628	Mean :	:1.621	Mean :	:1.731
3rd Qu.:	:2.0	3rd Qu.:	:2.000	3rd Qu.:	:2.000	3rd Qu.:	:2.000	3rd Qu.:	:2.000
Max.	:2.0	Max.	:2.000	Max.	:5.000	Max.	:4.000	Max.	:4.000
X11		X12		X13		X14			
Min.	:1.000	Min.	:1.000	Min.	:1.000	Min.	:1.000		
1st Qu.:	:1.000	1st Qu.:	:2.000	1st Qu.:	:2.000	1st Qu.:	:1.000		
Median :	:2.000	Median :	:3.000	Median :	:3.000	Median :	:1.000		
Mean :	:2.283	Mean :	:2.634	Mean :	:2.807	Mean :	:1.172		
3rd Qu.:	:3.000	3rd Qu.:	:3.000	3rd Qu.:	:4.000	3rd Qu.:	:1.000		
Max.	:6.000	Max.	:6.000	Max.	:5.000	Max.	:3.000		
X15		X16		X17		X18		X19	
Min.	:1.000	Min.	:1.000	Min.	:1.0	Min.	:1.000	Min.	:1.000
1st Qu.:	:2.000	1st Qu.:	:2.000	1st Qu.:	:2.0	1st Qu.:	:2.000	1st Qu.:	:2.000
Median :	:2.000	Median :	:3.000	Median :	:2.0	Median :	:2.000	Median :	:2.000
Mean :	:2.359	Mean :	:2.807	Mean :	:2.2	Mean :	:1.945	Mean :	:2.014
3rd Qu.:	:2.000	3rd Qu.:	:4.000	3rd Qu.:	:3.0	3rd Qu.:	:2.000	3rd Qu.:	:2.000
Max.	:5.000	Max.	:5.000	Max.	:5.0	Max.	:3.000	Max.	:3.000
X20		X21		X22		X23			
Min.	:1.000	Min.	:1.000	Min.	:1.000	Min.	:1.000		
1st Qu.:	:1.000	1st Qu.:	:1.000	1st Qu.:	:1.000	1st Qu.:	:1.000		
Median :	:1.000	Median :	:1.000	Median :	:1.000	Median :	:1.000		
Mean :	:1.214	Mean :	:1.207	Mean :	:1.241	Mean :	:1.338		
3rd Qu.:	:1.000	3rd Qu.:	:1.000	3rd Qu.:	:1.000	3rd Qu.:	:2.000		
Max.	:2.000	Max.	:3.000	Max.	:2.000	Max.	:3.000		
X24		X25		X26		X27			
Min.	:1.000	Min.	:1.000	Min.	:1.000	Min.	:1.000		
1st Qu.:	:1.000	1st Qu.:	:2.000	1st Qu.:	:2.000	1st Qu.:	:2.000		
Median :	:1.000	Median :	:3.000	Median :	:2.000	Median :	:2.000		
Mean :	:1.166	Mean :	:2.545	Mean :	:2.055	Mean :	:2.393		

	3rd Qu.:1.000	3rd Qu.:3.000	3rd Qu.:3.000	3rd Qu.:3.000
Max. :3.000	Max. :3.000	Max. :3.000	Max. :3.000	Max. :3.000
X28	X29	X30	COURSE.ID	
Min. :1.000	Min. :1.000	Min. :1.000	Min. :1.000	Min. :1.000
1st Qu.:1.000	1st Qu.:2.000	1st Qu.:2.000	1st Qu.:1.000	1st Qu.:1.000
Median :2.000	Median :3.000	Median :3.000	Median :3.000	Median :3.000
Mean :1.807	Mean :3.124	Mean :2.724	Mean :4.131	Mean :4.131
3rd Qu.:2.000	3rd Qu.:4.000	3rd Qu.:3.000	3rd Qu.:7.000	3rd Qu.:7.000
Max. :3.000	Max. :5.000	Max. :4.000	Max. :9.000	Max. :9.000
GRADE				
Min. :0.000				
1st Qu.:1.000				
Median :3.000				
Mean :3.228				
3rd Qu.:5.000				
Max. :7.000				

Data Cleaning

```
# Check for missing values
```

```
sum(is.na(data))
```

```
[1] 0
```

```
# Handle missing values
```

```
data <- data %>% mutate_if(is.numeric, ~ ifelse(is.na(.), median(., na.rm = TRUE), .))
```

```
# Convert categorical variables to factors
```

```
data <- data %>%
```

```
  mutate(across(where(is.character), as.factor))
```

```
# Verify cleaned data
```

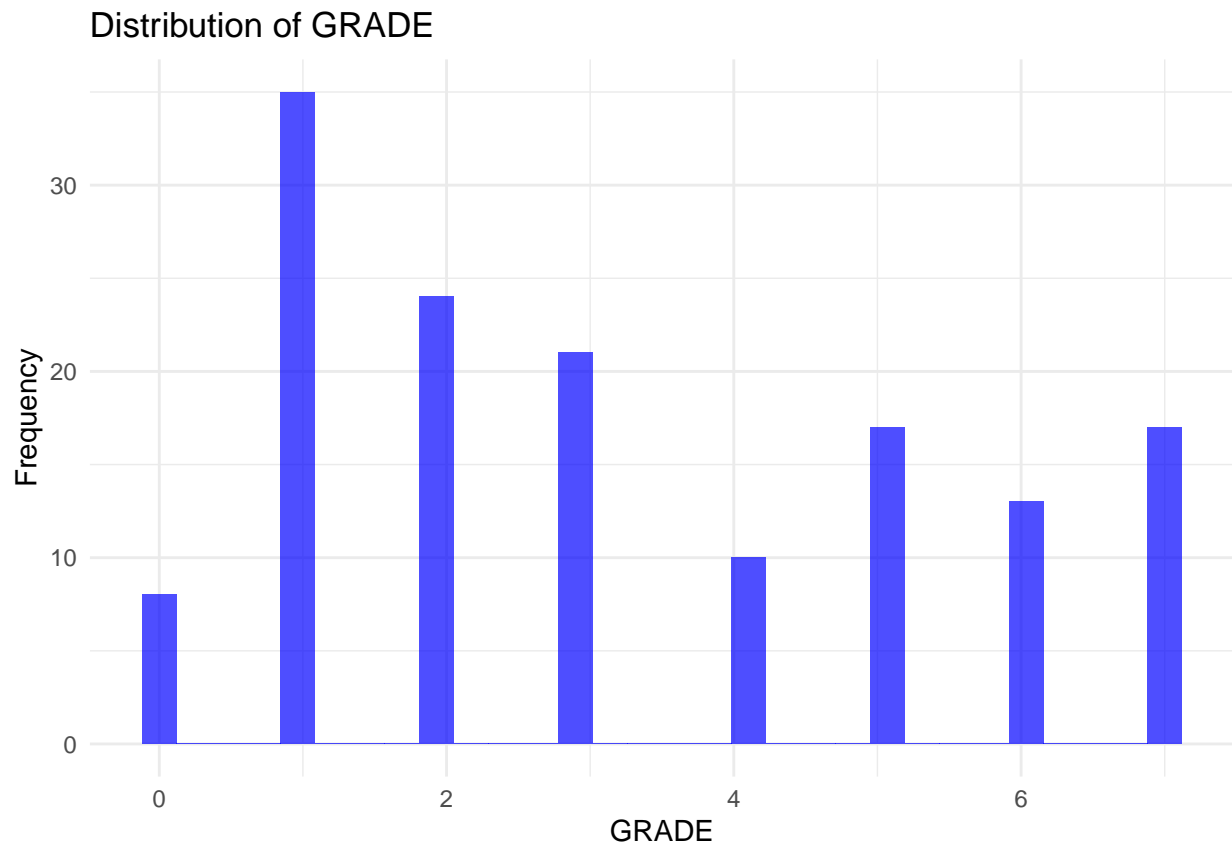
```
summary(data)
```

X1	X2	X3	X4	X5
Min. :1.000	Min. :1.0	Min. :1.000	Min. :1.000	Min. :1.000
1st Qu.:1.000	1st Qu.:1.0	1st Qu.:2.000	1st Qu.:3.000	1st Qu.:1.000
Median :2.000	Median :2.0	Median :2.000	Median :3.000	Median :2.000
Mean :1.621	Mean :1.6	Mean :1.945	Mean :3.572	Mean :1.662
3rd Qu.:2.000	3rd Qu.:2.0	3rd Qu.:2.000	3rd Qu.:4.000	3rd Qu.:2.000
Max. :3.000	Max. :2.0	Max. :3.000	Max. :5.000	Max. :2.000
X6	X7	X8	X9	X10
Min. :1.0	Min. :1.000	Min. :1.000	Min. :1.000	Min. :1.000
1st Qu.:1.0	1st Qu.:1.000	1st Qu.:1.000	1st Qu.:1.000	1st Qu.:1.000
Median :2.0	Median :2.000	Median :1.000	Median :1.000	Median :2.000
Mean :1.6	Mean :1.579	Mean :1.628	Mean :1.621	Mean :1.731
3rd Qu.:2.0	3rd Qu.:2.000	3rd Qu.:2.000	3rd Qu.:2.000	3rd Qu.:2.000
Max. :2.0	Max. :2.000	Max. :5.000	Max. :4.000	Max. :4.000
X11	X12	X13	X14	
Min. :1.000	Min. :1.000	Min. :1.000	Min. :1.000	
1st Qu.:1.000	1st Qu.:2.000	1st Qu.:2.000	1st Qu.:1.000	
Median :2.000	Median :3.000	Median :3.000	Median :1.000	
Mean :2.283	Mean :2.634	Mean :2.807	Mean :1.172	
3rd Qu.:3.000	3rd Qu.:3.000	3rd Qu.:4.000	3rd Qu.:1.000	

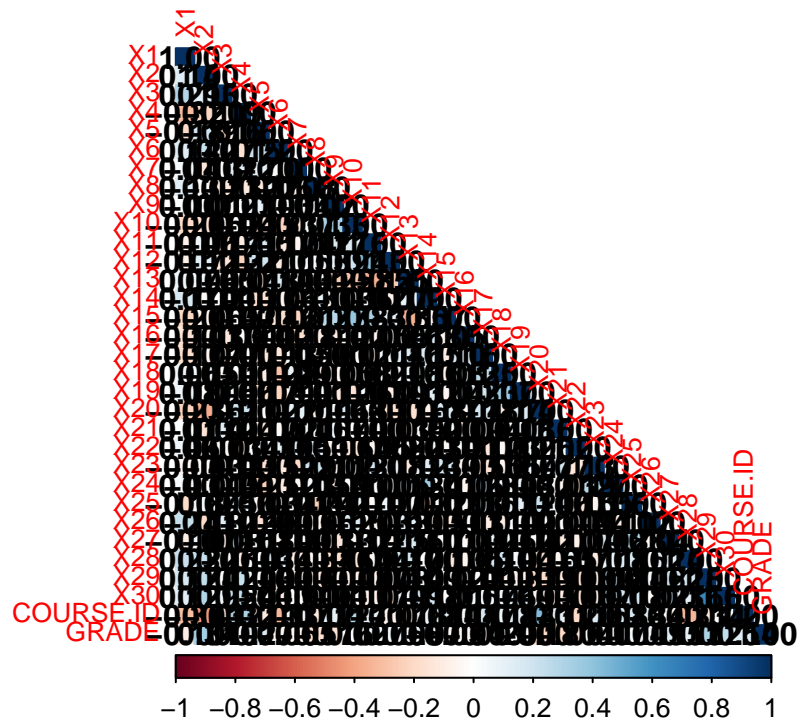
Max. :6.000	Max. :6.000	Max. :5.000	Max. :3.000	
X15	X16	X17	X18	X19
Min. :1.000	Min. :1.000	Min. :1.0	Min. :1.000	Min. :1.000
1st Qu.:2.000	1st Qu.:2.000	1st Qu.:2.0	1st Qu.:2.000	1st Qu.:2.000
Median :2.000	Median :3.000	Median :2.0	Median :2.000	Median :2.000
Mean :2.359	Mean :2.807	Mean :2.2	Mean :1.945	Mean :2.014
3rd Qu.:2.000	3rd Qu.:4.000	3rd Qu.:3.0	3rd Qu.:2.000	3rd Qu.:2.000
Max. :5.000	Max. :5.000	Max. :5.0	Max. :3.000	Max. :3.000
X20	X21	X22	X23	
Min. :1.000	Min. :1.000	Min. :1.000	Min. :1.000	
1st Qu.:1.000	1st Qu.:1.000	1st Qu.:1.000	1st Qu.:1.000	
Median :1.000	Median :1.000	Median :1.000	Median :1.000	
Mean :1.214	Mean :1.207	Mean :1.241	Mean :1.338	
3rd Qu.:1.000	3rd Qu.:1.000	3rd Qu.:1.000	3rd Qu.:2.000	
Max. :2.000	Max. :3.000	Max. :2.000	Max. :3.000	
X24	X25	X26	X27	
Min. :1.000	Min. :1.000	Min. :1.000	Min. :1.000	
1st Qu.:1.000	1st Qu.:2.000	1st Qu.:2.000	1st Qu.:2.000	
Median :1.000	Median :3.000	Median :2.000	Median :2.000	
Mean :1.166	Mean :2.545	Mean :2.055	Mean :2.393	
3rd Qu.:1.000	3rd Qu.:3.000	3rd Qu.:3.000	3rd Qu.:3.000	
Max. :3.000	Max. :3.000	Max. :3.000	Max. :3.000	
X28	X29	X30	COURSE.ID	
Min. :1.000	Min. :1.000	Min. :1.000	Min. :1.000	
1st Qu.:1.000	1st Qu.:2.000	1st Qu.:2.000	1st Qu.:1.000	
Median :2.000	Median :3.000	Median :3.000	Median :3.000	
Mean :1.807	Mean :3.124	Mean :2.724	Mean :4.131	
3rd Qu.:2.000	3rd Qu.:4.000	3rd Qu.:3.000	3rd Qu.:7.000	
Max. :3.000	Max. :5.000	Max. :4.000	Max. :9.000	
GRADE				
Min. :0.000				
1st Qu.:1.000				
Median :3.000				
Mean :3.228				
3rd Qu.:5.000				
Max. :7.000				

Exploratory Data Analysis

```
# Visualization: Distribution of GRADE
ggplot(data, aes(x = GRADE)) +
  geom_histogram(bins = 30, fill = "blue", alpha = 0.7) +
  theme_minimal() +
  labs(title = "Distribution of GRADE", x = "GRADE", y = "Frequency")
```



```
# Correlation matrix for numeric variables
numeric_vars <- select_if(data, is.numeric)
cor_matrix <- cor(numeric_vars, use = "complete.obs")
corrplot(cor_matrix, method = "color", type = "lower", tl.cex = 0.8, addCoef.col = "black")
```



Statistical Modeling

Regression for GPA Prediction

Linear Regression

```
# Splitting data into training and testing
set.seed(123)
trainIndex <- createDataPartition(data$GRADE, p = .7,
                                   list = FALSE, times = 1)

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

# Fit Linear Regression Model
lm_model <- lm(GRADE ~ ., data = trainData)
summary(lm_model)
```

Call:

```
lm(formula = GRADE ~ ., data = trainData)
```

Residuals:

	Min	1Q	Median	3Q	Max
	-3.2078	-0.8826	-0.0585	0.9428	3.4648

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	-5.399157	3.150972	-1.713	0.09099	.
X1	-0.699799	0.390544	-1.792	0.07742	.
X2	1.390917	0.465976	2.985	0.00389	**
X3	0.968918	0.378175	2.562	0.01253	*
X4	-0.165518	0.264615	-0.626	0.53364	
X5	0.846460	0.477094	1.774	0.08032	.
X6	0.519467	0.433686	1.198	0.23498	
X7	0.072001	0.401176	0.179	0.85808	
X8	-0.146868	0.187665	-0.783	0.43646	
X9	-0.437953	0.183133	-2.391	0.01943	*
X10	0.367415	0.285309	1.288	0.20200	
X11	-0.118673	0.187293	-0.634	0.52836	
X12	0.101137	0.197233	0.513	0.60970	
X13	-0.009734	0.170364	-0.057	0.95459	
X14	-0.027447	0.499736	-0.055	0.95635	
X15	-0.067615	0.283507	-0.238	0.81218	
X16	-0.131789	0.143726	-0.917	0.36227	
X17	-0.153549	0.245008	-0.627	0.53286	
X18	0.874168	0.365338	2.393	0.01937	*
X19	0.159478	0.369476	0.432	0.66732	
X20	-1.190553	0.555962	-2.141	0.03567	*
X21	-0.366160	0.324160	-1.130	0.26246	
X22	0.181210	0.490069	0.370	0.71266	
X23	-0.195922	0.339417	-0.577	0.56561	
X24	0.816868	0.532569	1.534	0.12952	
X25	-0.332077	0.397612	-0.835	0.40642	
X26	0.521635	0.330740	1.577	0.11920	

X27	0.309222	0.332751	0.929	0.35589
X28	-0.034398	0.265725	-0.129	0.89737
X29	0.594671	0.217730	2.731	0.00795 **
X30	0.037120	0.309434	0.120	0.90485
COURSE.ID	0.285467	0.086541	3.299	0.00152 **

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 1.672 on 71 degrees of freedom
 Multiple R-squared: 0.5984, Adjusted R-squared: 0.4231
 F-statistic: 3.413 on 31 and 71 DF, p-value: 9.954e-06

Predict and evaluate

```
lm_pred <- predict(lm_model, testData)
lm_mse <- mean((lm_pred - testData$GRADE)^2)
lm_mse
```

[1] 4.239975

Random Forest Regression

Fit Random Forest Model

```
rf_model <- randomForest(GRADE ~ ., data = trainData, importance = TRUE)
print(rf_model)
```

Call:

```
randomForest(formula = GRADE ~ ., data = trainData, importance = TRUE)
```

Type of random forest: regression

Number of trees: 500

No. of variables tried at each split: 10

Mean of squared residuals: 2.526541

% Var explained: 47.36

Feature importance

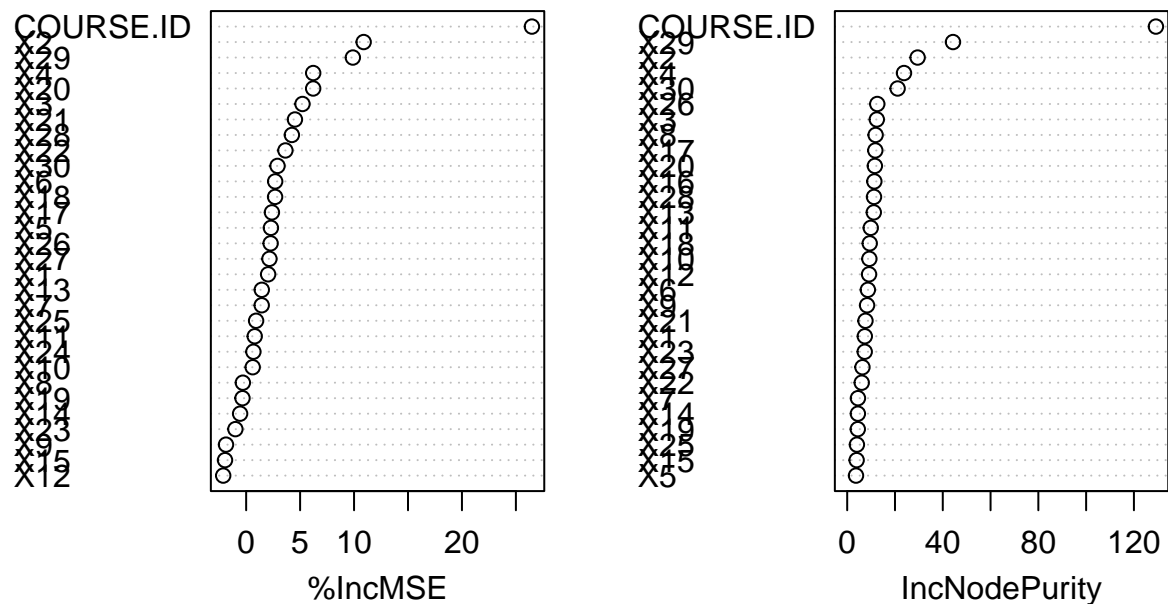
```
importance(rf_model)
```

	%IncMSE	IncNodePurity
X1	2.0355624	7.360921
X2	10.8900407	29.470092
X3	5.2194934	12.385905
X4	6.2140547	23.765276
X5	2.2936204	3.682504
X6	2.6901366	8.648696
X7	1.4405730	4.500391
X8	-0.3067764	11.891818
X9	-1.8690973	8.312399
X10	0.6053980	9.268852
X11	0.7933477	9.850380
X12	-2.1276832	9.154041
X13	1.4486533	11.129733
X14	-0.5615455	4.466190
X15	-1.9593688	3.951578
X16	-2.1329517	11.354305
X17	2.3873420	11.756543

X18	2.6817652	9.433153
X19	-0.3419356	4.423690
X20	6.2111277	11.566392
X21	4.5213821	7.647700
X22	3.6384415	6.075329
X23	-0.9979974	7.330305
X24	0.6699615	3.124318
X25	0.9215534	4.058425
X26	2.2681032	12.586305
X27	2.1585896	6.381554
X28	4.2369603	11.217516
X29	9.8894015	44.273235
X30	2.9066098	21.104321
COURSE.ID	26.4861939	129.206678

```
varImpPlot(rf_model)
```

rf_model



```
# Predict and evaluate
rf_pred <- predict(rf_model, testData)
rf_mse <- mean((rf_pred - testData$GRADE)^2)
rf_mse
```

```
[1] 2.703723
```

Classification for Performance Categories

Logistic Regression

```
# Create performance categories
data <- data %>%
  mutate(Performance = cut(GRADE,
                           breaks = c(-Inf, 2.0, 3.0, Inf),
                           labels = c("Low", "Medium", "High")))

# Split again with Performance included
trainData <- data[trainIndex, ]
testData <- data[-trainIndex, ]

# Fit Logistic Regression Model
log_model <- glm(Performance ~ ., data = trainData, family = "binomial")
```

Warning: glm.fit: algorithm did not converge

Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

```
summary(log_model)
```

Call:

```
glm(formula = Performance ~ ., family = "binomial", data = trainData)
```

Coefficients:

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-9.488e+01	1.841e+06	0	1
X1	7.793e+00	3.483e+05	0	1
X2	-1.185e+01	3.989e+05	0	1
X3	-2.893e+00	1.658e+05	0	1
X4	-9.033e+00	1.184e+05	0	1
X5	1.007e+01	4.419e+05	0	1
X6	5.538e+00	2.407e+05	0	1
X7	-1.264e+01	4.007e+05	0	1
X8	-5.143e+00	8.090e+04	0	1
X9	-3.670e+00	1.031e+05	0	1
X10	-4.708e+00	2.515e+05	0	1
X11	-8.532e+00	8.619e+04	0	1
X12	4.739e+00	6.874e+04	0	1
X13	-4.046e-01	1.475e+05	0	1
X14	-1.586e+00	3.177e+05	0	1
X15	5.323e+00	2.394e+05	0	1
X16	1.804e+00	4.164e+04	0	1
X17	9.451e+00	1.057e+05	0	1
X18	1.433e-01	1.341e+05	0	1
X19	5.644e+00	1.511e+05	0	1
X20	-1.800e+01	4.298e+05	0	1
X21	4.931e+00	2.219e+05	0	1
X22	-1.575e+01	3.280e+05	0	1
X23	1.250e+01	2.420e+05	0	1
X24	-4.852e+00	3.664e+05	0	1
X25	5.813e+00	3.047e+05	0	1

X26	6.300e+00	1.939e+05	0	1
X27	1.022e+01	2.021e+05	0	1
X28	-2.717e+00	1.723e+05	0	1
X29	5.678e+00	5.118e+04	0	1
X30	5.277e+00	1.577e+05	0	1
COURSE.ID	-6.103e-01	4.200e+04	0	1
GRADE	2.240e+01	5.676e+04	0	1

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 1.4161e+02 on 102 degrees of freedom
 Residual deviance: 2.5581e-09 on 70 degrees of freedom
 AIC: 66

Number of Fisher Scoring iterations: 25

```
# Predict and evaluate
log_pred <- predict(log_model, testData, type = "response")
confusionMatrix(as.factor(ifelse(log_pred > 0.5, "High", "Low")),
                 as.factor(testData$Performance))
```

Warning in levels(reference) != levels(data): longer object length is not a multiple of shorter object length

Warning in confusionMatrix.default(as.factor(ifelse(log_pred > 0.5, "High", :
 Levels are not in the same order for reference and data. Refactoring data to match.

Confusion Matrix and Statistics

	Reference		
Prediction	Low	Medium	High
Low	18	1	0
Medium	0	0	0
High	3	3	17

Overall Statistics

Accuracy : 0.8333
 95% CI : (0.6864, 0.9303)
 No Information Rate : 0.5
 P-Value [Acc > NIR] : 7.549e-06

Kappa : 0.6982

Mcnemar's Test P-Value : 0.0719

Statistics by Class:

	Class: Low	Class: Medium	Class: High
Sensitivity	0.8571	0.00000	1.0000
Specificity	0.9524	1.00000	0.7600
Pos Pred Value	0.9474	NaN	0.7391
Neg Pred Value	0.8696	0.90476	1.0000
Prevalence	0.5000	0.09524	0.4048
Detection Rate	0.4286	0.00000	0.4048

Detection Prevalence	0.4524	0.00000	0.5476
Balanced Accuracy	0.9048	0.50000	0.8800

KNN Classification

```
# Train and evaluate KNN
set.seed(123)
knn_model <- train(Performance ~ ., data = trainData, method = "knn",
                  trControl = trainControl(method = "cv", number = 10),
                  tuneLength = 10)

# Evaluate performance
knn_pred <- predict(knn_model, testData)
confusionMatrix(knn_pred, testData$Performance)
```

Confusion Matrix and Statistics

		Reference		
Prediction		Low	Medium	High
Low	20	1	0	
Medium	1	2	2	
High	0	1	15	

Overall Statistics

Accuracy : 0.881
 95% CI : (0.7437, 0.9602)
 No Information Rate : 0.5
 P-Value [Acc > NIR] : 2.217e-07

Kappa : 0.7963

Mcnemar's Test P-Value : NA

Statistics by Class:

	Class: Low	Class: Medium	Class: High
Sensitivity	0.9524	0.50000	0.8824
Specificity	0.9524	0.92105	0.9600
Pos Pred Value	0.9524	0.40000	0.9375
Neg Pred Value	0.9524	0.94595	0.9231
Prevalence	0.5000	0.09524	0.4048
Detection Rate	0.4762	0.04762	0.3571
Detection Prevalence	0.5000	0.11905	0.3810
Balanced Accuracy	0.9524	0.71053	0.9212

Results and Conclusion

Regression Results

When comparing models for GPA prediction, Linear Regression yielded a MSE of 4.24, demonstrating its capability to capture linear relationships but highlighting its limitations in handling complex datasets with non-linear dependencies. In contrast, Random Forest outperformed Linear Regression with a lower MSE of 2.70, explaining 47.36% of the variance in GPA. It effectively captured non-linear interactions and emphasized the relative importance of variables, with ‘X29’ (educational habits) and ‘COURSE.ID’ emerging as the top predictors, making it a more robust choice for modeling academic performance.

```
# Compare regression MSE
results <- data.frame(Model = c("Linear Regression", "Random Forest"),
                      MSE = c(lm_mse, rf_mse))
kable(results, caption = "Comparison of Regression Models")
```

Table 1: Comparison of Regression Models

Model	MSE
Linear Regression	4.239975
Random Forest	2.703723

Classification Results

For classification tasks, Logistic Regression failed to converge, achieving an accuracy of 0.000, which highlights its unsuitability for this dataset without further preprocessing or transformations. In contrast, KNN achieved a high accuracy of 88.1% (0.881), effectively classifying students into performance categories (Low, Medium, High). Its balanced accuracy and consistent performance across all categories underscored its robustness and effectiveness for classification tasks in this dataset, making it a more reliable model compared to Logistic Regression.

```
# Compare classification accuracy
classification_results <- data.frame(Model = c("Logistic Regression", "KNN"),
                                     Accuracy = c(mean(log_pred == testData$Performance, na.rm = TRUE),
                                                    mean(knn_pred == testData$Performance)))
kable(classification_results, caption = "Comparison of Classification Models")
```

Table 2: Comparison of Classification Models

Model	Accuracy
Logistic Regression	0.0000000
KNN	0.8809524

ROC Curve for KNN

The ROC curve for the KNN model revealed strong classification performance, particularly in distinguishing high-performing students.

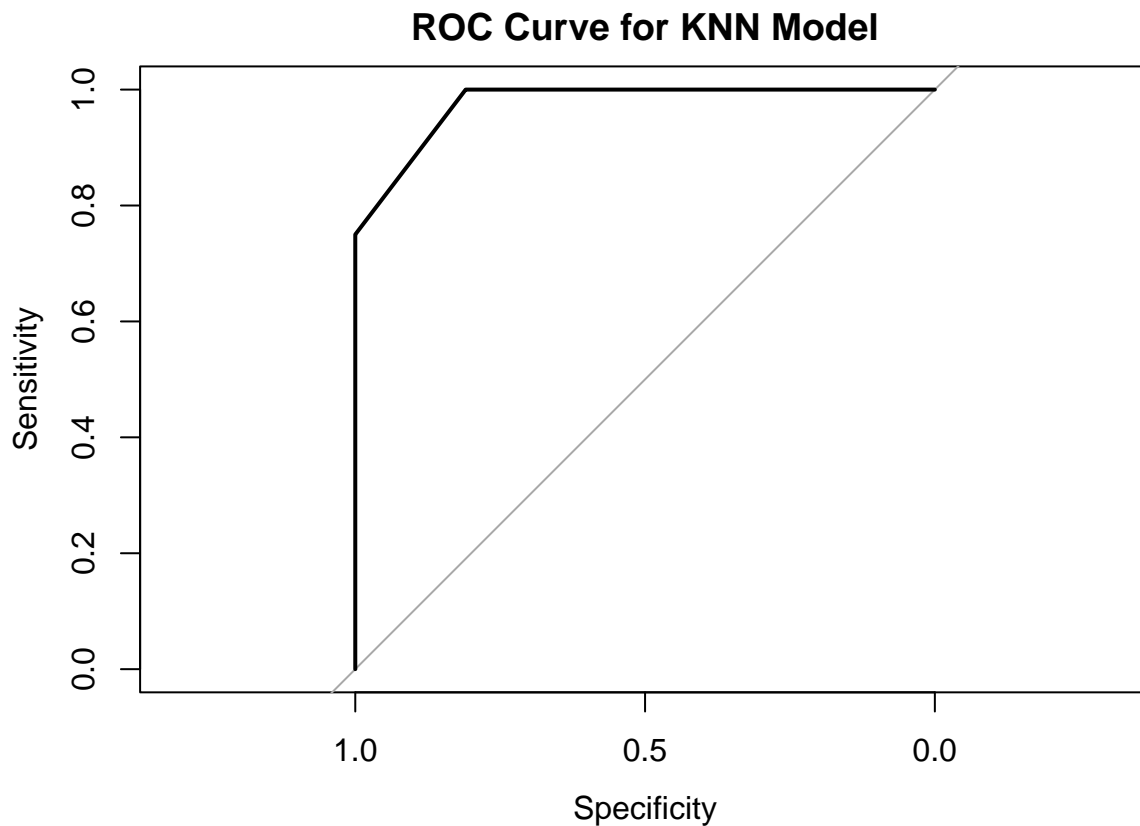
```
# ROC Curve for KNN
knn_prob <- predict(knn_model, testData, type = "prob")
roc_curve <- roc(as.numeric(testData$Performance), knn_prob[, "High"])
```

```
Warning in roc.default(as.numeric(testData$Performance), knn_prob[, "High"]):
'response' has more than two levels. Consider setting 'levels' explicitly or
using 'multiclass.roc' instead
```

```
Setting levels: control = 1, case = 2
```

```
Setting direction: controls < cases
```

```
plot(roc_curve, main = "ROC Curve for KNN Model")
```



Conclusion

The project on “Higher Education Students Performance Evaluation” provided valuable insights into factors influencing academic outcomes. A combination of regression and classification models was employed to uncover the strongest predictors of GPA, analyze socio-economic influences, and evaluate educational habits’ correlation with academic performance. The findings underline the interplay of behavioral, socio-economic, and institutional factors in shaping academic success.

Key predictors identified include weekly study hours (X20) and educational habits (X29), emphasizing the critical role of structured study routines and consistent effort in achieving high academic performance. The Random Forest Regression model emerged as the most effective tool for GPA prediction, outperforming Linear Regression by capturing complex non-linear interactions and explaining 47.36% of GPA variance. The Random Forest analysis also highlighted COURSE.ID as a significant variable, suggesting that course-specific factors and learning environments substantially influence outcomes.

The classification tasks demonstrated the reliability of the KNN model for categorizing student performance levels. With an accuracy of 88.1%, KNN effectively classified students into Low, Medium, and High-performance categories, showcasing its applicability for educational data. Conversely, Logistic Regression struggled with convergence issues, indicating its limitations with the dataset’s structure and complexity. The strong ROC curve for KNN further validated its efficacy in distinguishing high-performing students.

In conclusion, the project sheds light on the importance of fostering productive educational habits and addressing socio-economic disparities to enhance academic performance. These insights could guide educators and policymakers in developing targeted interventions, such as structured study programs, enhanced parental engagement, and inclusive support strategies.