

Felix-Ochieng_2024-M132-20790_Big-Data-Analytics-Project.R

HP

2026-02-11

```
# Load required libraries
```

```
library(data.table)
```

```
## Warning: package 'data.table' was built under R version 4.4.3
```

```
library(dplyr)
```

```
## Warning: package 'dplyr' was built under R version 4.4.3
```

```
##
```

```
## Attaching package: 'dplyr'
```

```
## The following objects are masked from 'package:data.table':
```

```
##
```

```
##      between, first, last
```

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

```
##
```

```
##      filter, lag
```

```
## The following objects are masked from 'package:base':
```

```
##
```

```
##      intersect, setdiff, setequal, union
```

```
library(ggplot2)
```

```
## Warning: package 'ggplot2' was built under R version 4.4.3
```

```
library(caret)
```

```
## Warning: package 'caret' was built under R version 4.4.3
```

```
## Loading required package: lattice
```

```
# PRACTICAL BIG DATA ANALYTICS SOLUTION USING R, UMU Case Study
# Step 1: Creating a Large Student Dataset (Volume). To simulate the university's large
#historical dataset to predict academic performance, a synthetic dataset of 20,000
#student records was generated. This dataset represents data collected from LMS usage,
#attendance tracking, and examination systems.
```

```
set.seed(2026)
```

```
student_data <- data.frame(
  Student_ID = 1:20000,
  Gender = sample(c("Male", "Female"), 20000, replace = TRUE),
  Attendance = sample(40:100, 20000, replace = TRUE),
  LMS_Engagement = sample(1:100, 20000, replace = TRUE),
  Assignment_Score = sample(10:40, 20000, replace = TRUE),
  Exam_Score = sample(20:60, 20000, replace = TRUE)
)
```

```
student_data$Final_Result <- ifelse(
  student_data$Attendance >= 70 &
  student_data$Exam_Score >= 40 &
  student_data$LMS_Engagement >= 50,
  "Pass",
  "Fail"
)
```

```
write.csv(student_data, "student_data_final_project.csv", row.names = FALSE)
```

```
# The above data shows realistic student academic indicators and clearly shows a
#Pass/Fail outcome for predictive modeling
```

```
#Step 2: Installing and Loading Required R Packages (Scalable analytics support)
install.packages(c("data.table", "dplyr", "ggplot2", "caret"))
```

```
## Warning: packages 'data.table', 'dplyr', 'ggplot2', 'caret' are in use and will
## not be installed
```

```
# Loading the packages
```

```
library(data.table)
library(dplyr)
library(ggplot2)
library(caret)
```

```
# Step 3: Loading the Big Dataset (Batch Processing)
```

```
student_data <- fread("student_data_final_project.csv")
```

```
#Step 4: Understanding the Data (Data understanding & validation)
```

```
# Checking structure
```

```
str(student_data)
```

```
## Classes 'data.table' and 'data.frame': 20000 obs. of 7 variables:
## $ Student_ID : int 1 2 3 4 5 6 7 8 9 10 ...
## $ Gender : chr "Male" "Male" "Male" "Female" ...
## $ Attendance : int 63 66 55 80 87 72 73 96 56 95 ...
```

```
## $ LMS_Engagement : int 6 63 7 97 94 55 88 31 48 66 ...
## $ Assignment_Score: int 26 20 23 22 28 37 24 15 14 25 ...
## $ Exam_Score      : int 34 36 39 34 33 48 60 45 43 40 ...
## $ Final_Result    : chr "Fail" "Fail" "Fail" "Fail" ...
## - attr(*, ".internal.selfref")=<externalptr>
```

```
#Checking the data size
nrow(student_data)
```

```
## [1] 20000
```

```
#Checking the data size
ncol(student_data)
```

```
## [1] 7
```

```
#Viewing first few rows
head(student_data)
```

```
##      Student_ID Gender Attendance LMS_Engagement Assignment_Score Exam_Score
##      <int> <char>      <int>          <int>          <int>          <int>
## 1:         1   Male         63             6             26             34
## 2:         2   Male         66             63            20             36
## 3:         3   Male         55             7             23             39
## 4:         4 Female         80             97            22             34
## 5:         5   Male         87             94            28             33
## 6:         6   Male         72             55            37             48
##      Final_Result
##      <char>
## 1:      Fail
## 2:      Fail
## 3:      Fail
## 4:      Fail
## 5:      Fail
## 6:      Pass
```

```
# Viewing data summary
summary(student_data)
```

```
##      Student_ID      Gender      Attendance      LMS_Engagement
## Min.   : 1 Length:20000 Min.   : 40.00 Min.   : 1.00
## 1st Qu.:5001 Class :character 1st Qu.: 55.00 1st Qu.: 26.00
## Median :10000 Mode  :character Median : 70.00 Median : 51.00
## Mean   :10000 Mean   : 70.22 Mean   : 50.58
## 3rd Qu.:15000 3rd Qu.: 86.00 3rd Qu.: 76.00
## Max.   :20000 Max.   :100.00 Max.   :100.00
## Assignment_Score Exam_Score Final_Result
## Min.   :10.00 Min.   :20.00 Length:20000
## 1st Qu.:17.00 1st Qu.:30.00 Class :character
## Median :25.00 Median :40.00 Mode  :character
## Mean   :25.16 Mean   :40.02
## 3rd Qu.:33.00 3rd Qu.:50.00
## Max.   :40.00 Max.   :60.00
```

```
# Step 5: Data Cleaning (Veracity Management. This ensures data accuracy and
#removes incomplete records
colSums(is.na(student_data))
```

```
##      Student_ID      Gender      Attendance      LMS_Engagement
##           0           0           0           0
## Assignment_Score      Exam_Score      Final_Result
##           0           0           0
```

```
student_data <- na.omit(student_data)
```

```
#Step 6: Efficient Data Processing (Scalable Analytics). This creates a total score
#and enables scalability for large data sets
```

```
student_data[, Total_Score := Assignment_Score + Exam_Score]
```

```
#Step 7: Descriptive Analytics (What Happened?). This gives a clear insight based on
#the overall pass vs fail rates, average performance levels and
#Gender-based performance trends
```

```
table(student_data$Final_Result)
```

```
##
## Fail Pass
## 17301 2699
```

```
mean(student_data$Exam_Score)
```

```
## [1] 40.0152
```

```
student_data[, .(Avg_Exam = mean(Exam_Score)), by = Gender]
```

```
##      Gender Avg_Exam
##      <char>   <num>
## 1:   Male 40.12949
## 2: Female 39.90230
```

```
#Step 8: Diagnostic Analytics (Why Did It Happen?). This gives an insight on the
#relationship between attendance and performance and the influence of
#LMS engagement on exam scores
```

```
cor(student_data$Attendance, student_data$Exam_Score)
```

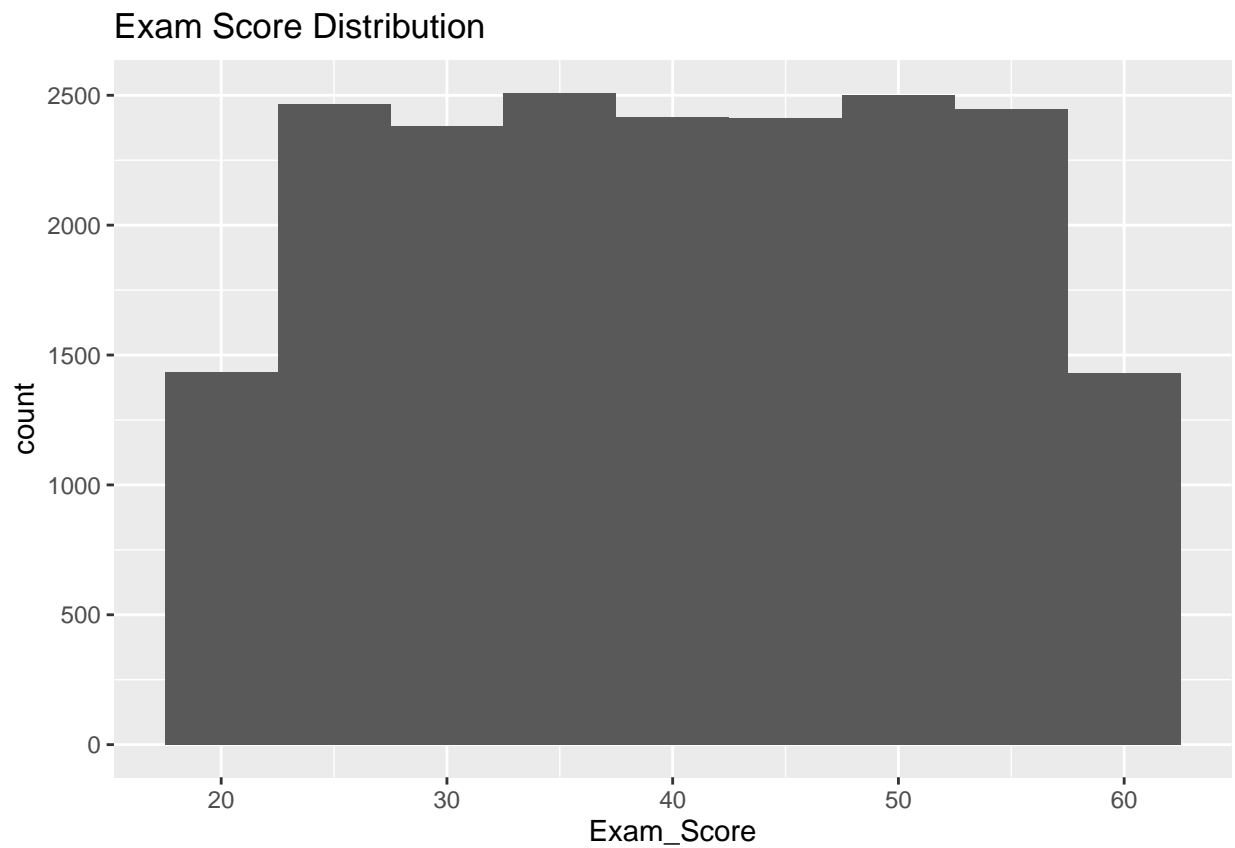
```
## [1] -0.007466529
```

```
cor(student_data$LMS_Engagement, student_data$Exam_Score)
```

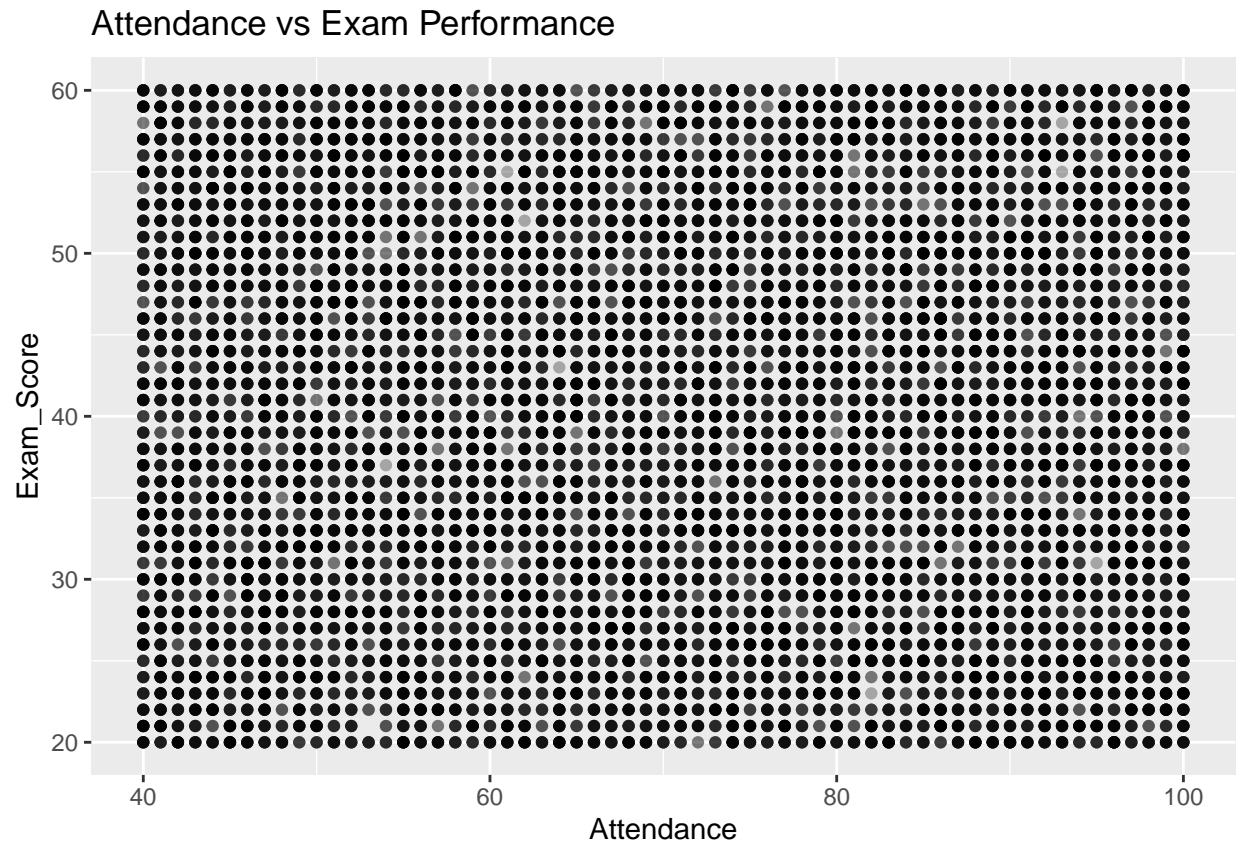
```
## [1] 0.001178173
```

*# From the above, the result shows a negative relationship between attendance and the
performance of students. This implies that as the student attendance
decreases, this negatively affects their performance in class.
However there is a positive relationship between LMS engagement and exam implying
that use of LMS engagement for learning increases the students' chances
of performing well in the exams*

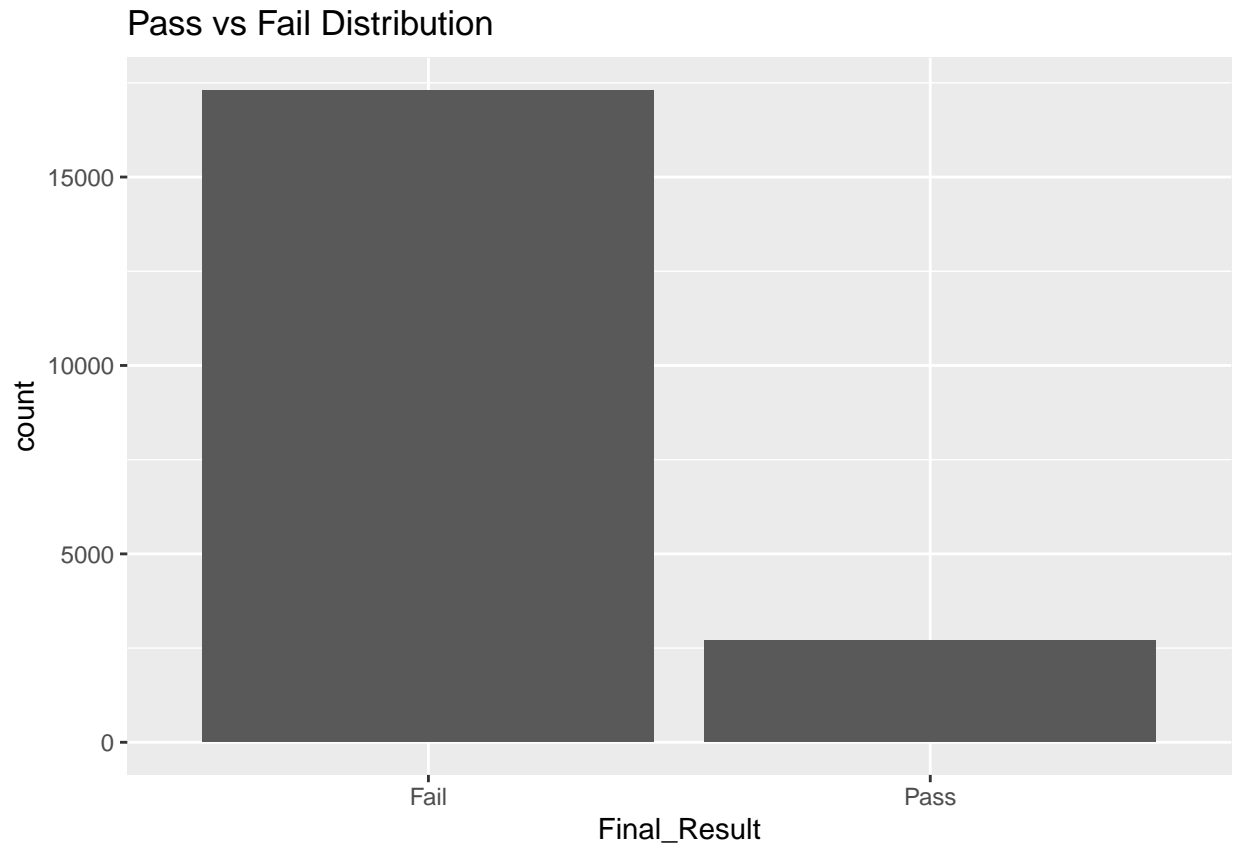
```
#Step 9: Big Data Visualization (Visual Analytics)  
ggplot(student_data, aes(x = Exam_Score)) +  
  geom_histogram(binwidth = 5) +  
  labs(title = "Exam Score Distribution")
```



```
ggplot(student_data, aes(x = Attendance, y = Exam_Score)) +  
  geom_point(alpha = 0.3) +  
  labs(title = "Attendance vs Exam Performance")
```



```
ggplot(student_data, aes(x = Final_Result)) +  
  geom_bar() +  
  labs(title = "Pass vs Fail Distribution")
```



```
#Step 10: Preparing Data for Prediction. The pass/fail has been coded as 1=Pass
#and 0=Fail
student_data$Result_Num <- ifelse(student_data$Final_Result == "Pass", 1, 0)

#Step 11: Splitting the Data (Predictive analytics preparation technique)
set.seed(2026)

train_index <- createDataPartition(
  student_data$Result_Num,
  p = 0.7,
  list = FALSE
)

train_data <- student_data[train_index]
test_data <- student_data[-train_index]

#Step 12: Building the Predictive Model. This uses predictive analytics techniques and
#logistic regression statistical method
model <- glm(
  Result_Num ~ Attendance + LMS_Engagement + Assignment_Score + Exam_Score,
  data = train_data,
  family = binomial
)

#Step 13: Understanding the Mode. This enables the data analyst to gain insights into the
#Significant predictors, direction of influence and
```

```
#model coefficients
summary(model)
```

```
##
## Call:
## glm(formula = Result_Num ~ Attendance + LMS_Engagement + Assignment_Score +
##      Exam_Score, family = binomial, data = train_data)
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)   -26.701973    0.612786 -43.575  <2e-16 ***
## Attendance      0.134623    0.003609  37.298  <2e-16 ***
## LMS_Engagement  0.081138    0.002203  36.827  <2e-16 ***
## Assignment_Score -0.004763    0.004394  -1.084    0.278
## Exam_Score      0.204461    0.005539  36.916  <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: 11048.4  on 13999  degrees of freedom
## Residual deviance:  4394.9  on 13995  degrees of freedom
## AIC: 4404.9
##
## Number of Fisher Scoring iterations: 8
```

```
#Step 14: Making Predictions
```

```
prob_predictions <- predict(model, test_data, type = "response")
class_predictions <- ifelse(prob_predictions >= 0.5, 1, 0)
```

```
#Step 15: Model Evaluation (Accuracy). This helps in model validation
```

```
confusionMatrix(
  as.factor(class_predictions),
  as.factor(test_data$Result_Num)
)
```

```
## Confusion Matrix and Statistics
##
##              Reference
## Prediction    0    1
##              0 5022 256
##              1  160 562
##
##              Accuracy : 0.9307
##              95% CI : (0.9239, 0.937)
##      No Information Rate : 0.8637
##      P-Value [Acc > NIR] : < 2.2e-16
##
##              Kappa : 0.6903
##
##      Mcnemar's Test P-Value : 3.197e-06
##
```

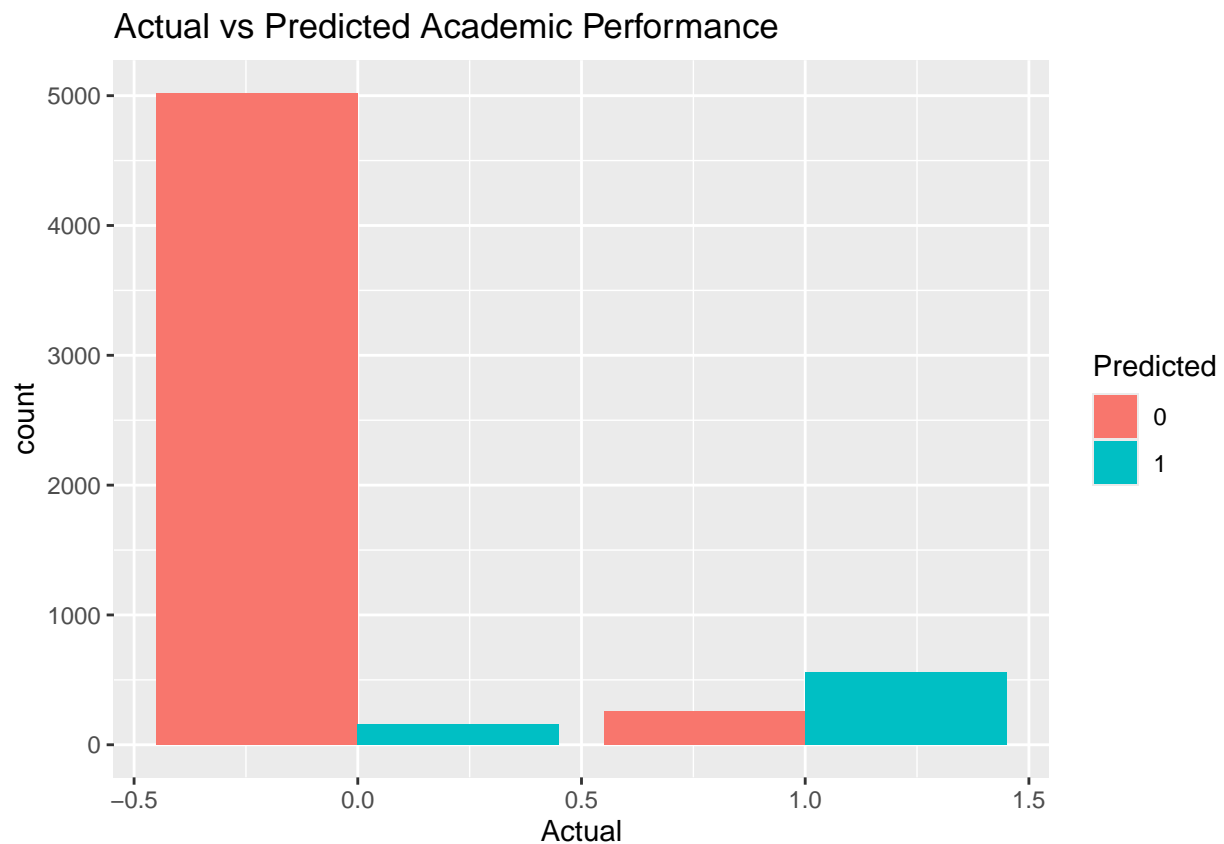


```
##           Sensitivity : 0.9691
##           Specificity : 0.6870
##           Pos Pred Value : 0.9515
##           Neg Pred Value : 0.7784
##           Prevalence : 0.8637
##           Detection Rate : 0.8370
##           Detection Prevalence : 0.8797
##           Balanced Accuracy : 0.8281
##
##           'Positive' Class : 0
##
```

#Step 16: Visualizing Prediction Results

```
prediction_results <- data.frame(
  Actual = test_data$Result_Num,
  Predicted = class_predictions
)

ggplot(prediction_results, aes(x = Actual, fill = as.factor(Predicted))) +
  geom_bar(position = "dodge") +
  labs(
    title = "Actual vs Predicted Academic Performance",
    fill = "Predicted"
  )
)
```



Interpretation

*#From the output, the predictive model achieved an accuracy of 93.07%, significantly
#outperforming baseline classification. WWith a sensitivity of 96.9%, the system is
#highly effective in identifying academically at-risk students, enabling timely
#intervention strategies. This demonstrates the value of Big Data Analytics in transforming
#historical educational data into actionable decision-support insights.*

#Conclusion

*# This practical implementation demonstrates how Big Data Analytics techniques can be
#applied using R to analyze large-scale historical student at UMU. By combining descriptive,
#diagnostic and predictive analytics, the solution enables early identification of
#academically at-risk students, student failure rates and supports data-driven academic
#decision-making. The approach aligns with Big Data principles of volume, scalability,
#prediction, validation and value extraction.*