# AMERICAN INTERNATIONAL UNIVERSITY-BANGLADESH

**FINAL-TERM PROJECT**

Course: INTRODUCTION TO DATA SCIENCE

Sec: A

<u>SUBMITTED TO</u>
Faculty: Tohedul Islam

<u>SUBMITTED BY</u>
Group: 15

| NAME | ID |
|------|----|
| 1. Abdullah Muhammad Hamja | 20 -43465-1 |
| 2. Sumaiya Ahmed Susmita | 21-45266-2 |

# DATA SCIENCE FINAL PROJECT
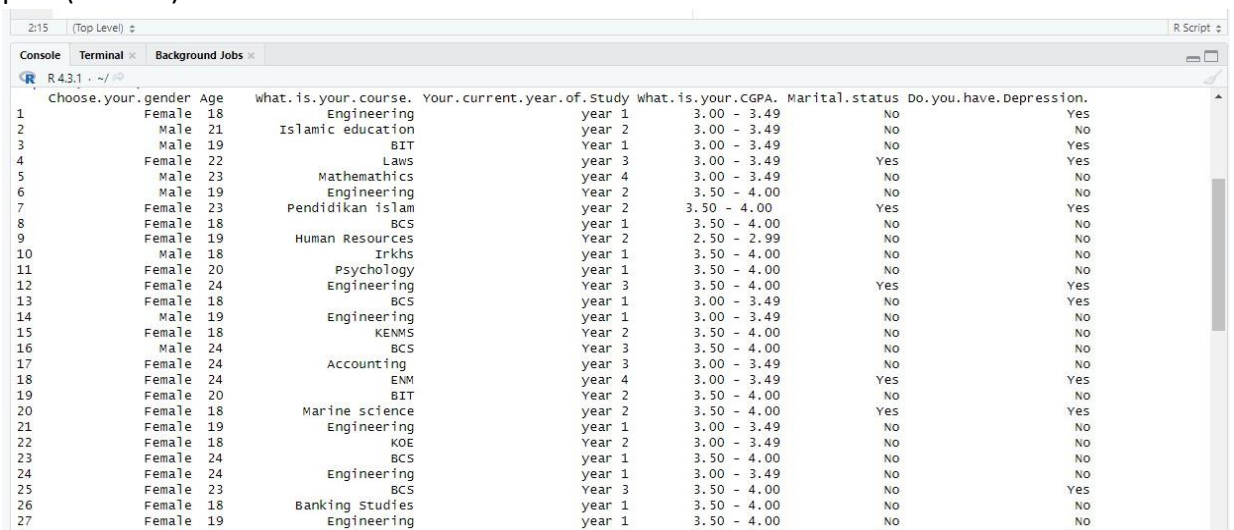## Student Mental Health

Dataset Description:

The 'Student Mental Health' dataset is a comprehensive collection aimed at exploring factors related to students' mental health. It includes key variables such as age, gender, CGPA (Cumulative Grade Point Average), course enrolment, and indicators for depression, anxiety, panic attacks, and seeking help. This dataset provides insights into demographic distributions, academic performance correlations with mental health, the prevalence of mental health conditions, and patterns of help-seeking behaviour. Potential use cases involve identifying risk factors, developing predictive models, and informing targeted interventions. Ethical considerations emphasize responsible data handling due to the sensitivity of mental health information. Overall, the dataset is a valuable resource for researchers, educators, and policymakers interested in addressing mental health challenges in the student population.

## CODES – CONSOLE – DETAILS:

1. Import CSV file.
   Code:
   dataset <- read.csv("D:/FALL2023/IntroToDataScience/Student_Mental_health.csv", header = TRUE, sep = ",")
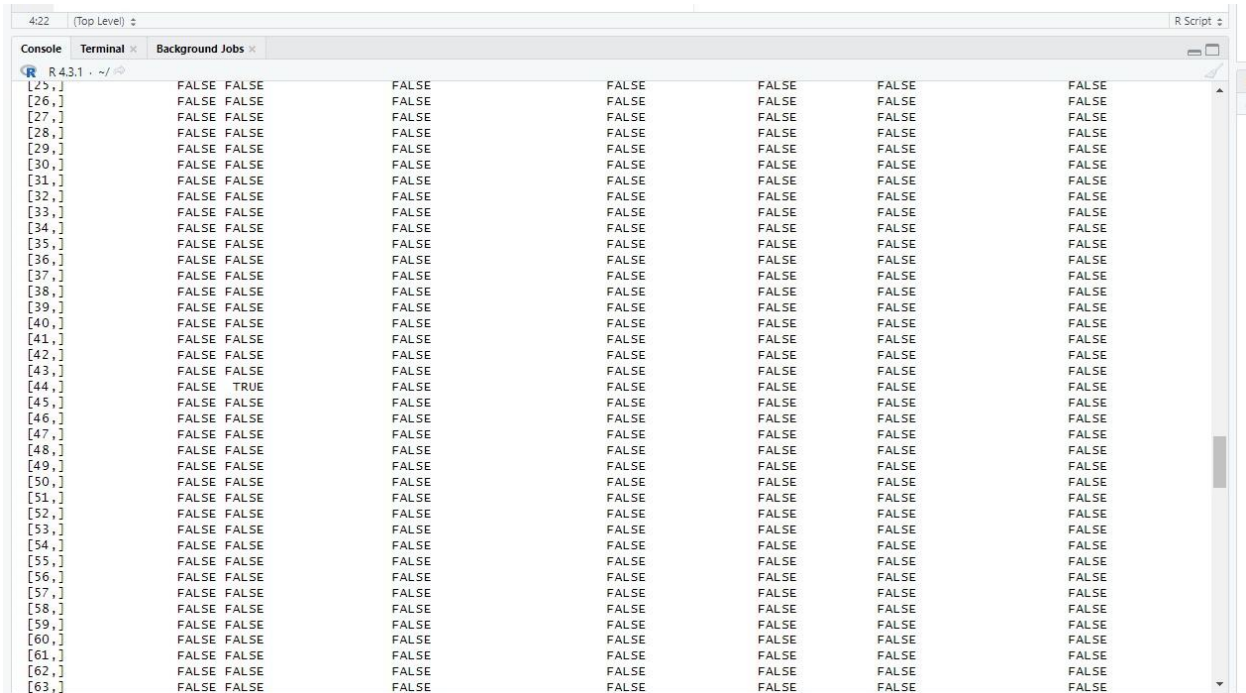   print(dataset)



*This code is used to import external Excel files (in CSV format) into R.*

2. Find missing values.
   Code:
   missing_values <- is.na(dataset)
   print(missing_values)



*This code is used to identify the missing values in a dataset. The one that has TRUE
written in it is a missing value.*

3. Discard missing values.
   Code:
   dataset <- na.omit(dataset)
   print(dataset)



*This code is used to remove the instance that had a missing value.*

4. <u>Converting Numeric Age values to Categorical Age values.</u>

Code:

```
catage <- "Age"
breaks <- c(0, 19, 25, Inf)
labels <- c("Teenager", "Young Adult", "Adult")
dataset$CatAge <- cut(dataset[[catage]], breaks = breaks, labels = labels,
include.lowest = TRUE)
print(dataset)
```

```
14:1    (Top Level) ≑

Console   Terminal   Background Jobs

R  R 4.3.1 · ~/
86          Female  18         psychology         year 1    3.50 - 4.00        No
87          Female  19         Fiqh fatwa         Year 3    3.00 - 3.49        No
88          Female  18         psychology         year 1    3.50 - 4.00        No
89            Male  24                BIT         year 1    3.00 - 3.49        No
90            Male  24        Engineering         Year 2    2.00 - 2.49        No
   Do.you.have.Anxiety. Do.you.have.Panic.attack. Did.you.seek.any.specialist.for.a.treatment.      CatAge
1                    No                       Yes                                            No     Teenager
2                   Yes                        No                                            No  Young Adult
3                   Yes                       Yes                                            No     Teenager
4                    No                        No                                            No  Young Adult
5                    No                        No                                            No  Young Adult
6                    No                       Yes                                            No     Teenager
7                    No                       Yes                                            No  Young Adult
8                   Yes                        No                                            No     Teenager
9                    No                        No                                            No     Teenager
10                  Yes                       Yes                                            No     Teenager
11                   No                        No                                            No  Young Adult
12                   No                        No                                            No  Young Adult
```

*This code is used to convert the categorize the age range to teenager, young adult, and adult. It is saved as a new column.*

5. <u>Pearson's Chi-squared Test</u>

Code:

```
dataset <- data.frame(
  "what.is.your.cgpa" = sample(c("0-1.99", "2-2.49", "2.5-2.99", "3-3.49", "3.5-4.00"),
100, replace = TRUE),
  "Do.you.have.anxiety" = sample(c("Yes", "No"), 100, replace = TRUE)
)
contingency_table <- table(dataset$What.is.your.cgpa.,
dataset$Do.you.have.anxiety.)
chi_squared_test <- chisq.test(contingency_table)
print(chi_squared_test)
```

```
> dataset <- data.frame(
+ "what.is.your.cgpa" = sample(c("0-1.99", "2-2.49", "2.5-2.99", "3-3.49", "3.5-4.00"), 100, replace = TRUE),
+ "Do.you.have.anxiety" = sample(c("Yes", "No"), 100, replace = TRUE)
+ )
> contingency_table <- table(dataset$what.is.your.cgpa, dataset$Do.you.have.anxiety)
> chi_squared_test <- chisq.test(contingency_table)
> print(chi_squared_test)

        Pearson's Chi-squared test

data:  contingency_table
X-squared = 2.9557, df = 4, p-value = 0.5653
```

*This code is used to find the significant attributes using Pearson's Chi-Squared Test.*

6. Install Package

   a) e1071
      Code:
      install.packages("e1071")
      library(e1071)

```
> library(e1071)
warning message:
package 'e1071' was built under R version 4.3.2
>
```

   *This code is used to install the necessary package for the 'Naïve Bayes'*
   *function.*

   b) caret
      Code:
      install.packages("caret")
      library(caret)

```
> library(caret)
Loading required package: ggplot2
Need help? Try Stackoverflow: https://stackoverflow.com/tags/ggplot2
Loading required package: lattice
warning messages:
1: package 'caret' was built under R version 4.3.2
2: package 'ggplot2' was built under R version 4.3.2
```

   *This code is used to install the necessary package for the 'Naïve Bayes'*
   *classification.*

7. Naïve Bayes
   Code:
   nb_model <- naiveBayes(What.is.your.course. ~ ., data = dataset)
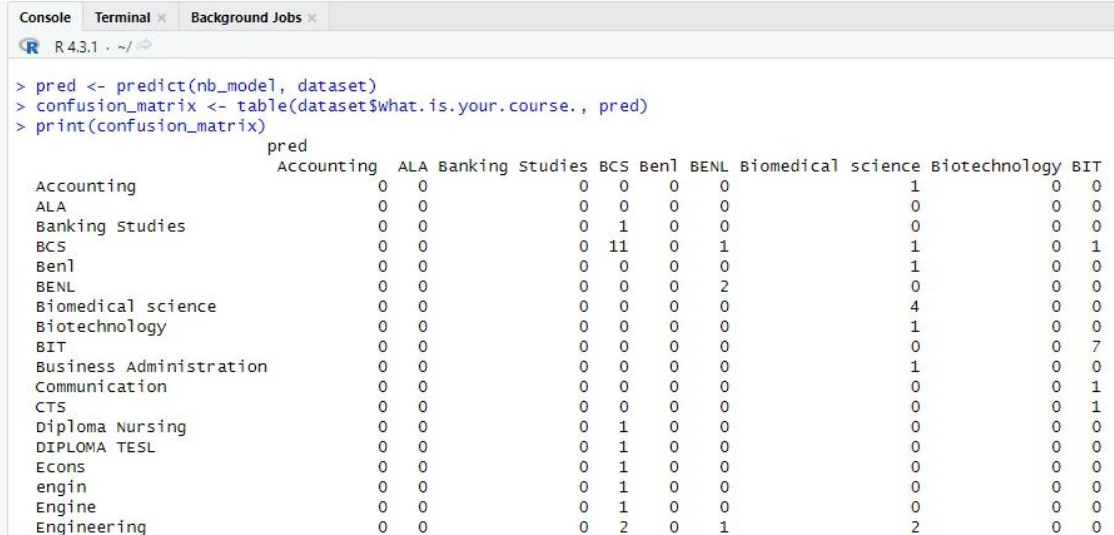   print(nb_model)



*This code is used to find the Naïve Bayes function to predict, based on the course of the student.*

## 8. Confusion Matrix

Code:
```
pred <- predict(nb_model, dataset)
confusion_matrix <- table(dataset$What.is.your.course., pred)
print(confusion_matrix)
```



```
> pred <- predict(nb_model, dataset)
> confusion_matrix <- table(dataset$what.is.your.course., pred)
> print(confusion_matrix)
                        pred
                         Accounting  ALA Banking Studies BCS Benl BENL Biomedical science Biotechnology BIT
  Accounting                      0    0               0   0    0    0                  1             0   0
  ALA                             0    0               0   0    0    0                  0             0   0
  Banking Studies                 0    0               0   1    0    0                  0             0   0
  BCS                             0    0               0  11    0    1                  1             0   1
  Benl                            0    0               0   0    0    0                  1             0   0
  BENL                            0    0               0   0    0    2                  0             0   0
  Biomedical science             0    0               0   0    0    0                  4             0   0
  Biotechnology                  0    0               0   0    0    0                  1             0   0
  BIT                            0    0               0   0    0    0                  0             0   7
  Business Administration        0    0               0   0    0    0                  1             0   0
  Communication                  0    0               0   0    0    0                  0             0   1
  CTS                            0    0               0   0    0    0                  0             0   1
  Diploma Nursing                0    0               0   1    0    0                  0             0   0
  DIPLOMA TESL                   0    0               0   1    0    0                  0             0   0
  Econs                          0    0               0   1    0    0                  0             0   0
  engin                          0    0               0   1    0    0                  0             0   0
  Engine                         0    0               0   1    0    0                  0             0   0
  Engineering                    0    0               0   2    0    1                  2             0   0
```
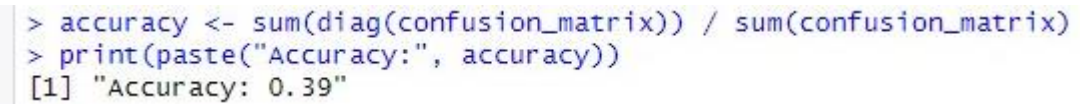
*This code has been used to find the confusion matrix on the student's courses.*

## 9. Accuracy

Code:
```
accuracy <- sum(diag(confusion_matrix)) / sum(confusion_matrix)
print(paste("Accuracy:", accuracy))
```



```
> accuracy <- sum(diag(confusion_matrix)) / sum(confusion_matrix)
> print(paste("Accuracy:", accuracy))
[1] "Accuracy: 0.39"
```

*This code is used to find the predictive accuracy of the Naïve Bayes classification.*

10. 10-Fold Cross Validation
    Code:
    set.seed(123)
    folds <- createFolds(dataset$Do.you.have.Depression., k = 10, list = TRUE, returnTrain = TRUE)

    for (i in 1:10) {
      train_fold <- dataset[unlist(folds[i]), ]
      test_fold <- dataset[-unlist(folds[i]), ]

      nb_model_fold <- naiveBayes(Do.you.have.Depression. ~ ., data = train_fold)

      predictions <- predict(nb_model_fold, test_fold)

      confusion_matrix <- table(test_fold$Do.you.have.Depression., predictions)
      accuracy <- sum(diag(confusion_matrix)) / sum(confusion_matrix)

      cat("Fold", i, "Accuracy:", accuracy, "\n")
    }

```
+
+   cat("Fold", i, "Accuracy:", accuracy, "\n")
+ }
Fold 1 Accuracy: 0.8
Fold 2 Accuracy: 0.9
Fold 3 Accuracy: 0.8181818
Fold 4 Accuracy: 0.7777778
Fold 5 Accuracy: 0.7
Fold 6 Accuracy: 0.8181818
Fold 7 Accuracy: 0.9
Fold 8 Accuracy: 0.8
Fold 9 Accuracy: 0.7777778
Fold 10 Accuracy: 0.8
>
```

*The code is used for 10 Fold Cross Validation and Accuracy.*

11. Train Set and Test Set
    Code:
    train_indices <- sample(1:nrow(dataset), 0.7 * nrow(dataset))
    train_data <- dataset[train_indices, ]
    test_data <- dataset[-train_indices, ]

```
> train_indices <- sample(1:nrow(dataset), 0.7 * nrow(dataset))
> train_data <- dataset[train_indices, ]
> test_data <- dataset[-train_indices, ]
```

   *This code is used to generate the train set and test set according to this dataset.*

## 12. Recall, Precision and F-measure value

Code:

```r
metrics <- data.frame(Recall = numeric(10), Precision = numeric(10), F_measure =
numeric(10))

for (i in 1:10) {
  train_fold <- dataset[unlist(folds[i]), ]
  test_fold <- dataset[-unlist(folds[i]), ]

  nb_model_fold <- naiveBayes(Do.you.have.Depression. ~ ., data = train_fold)

  predictions <- predict(nb_model_fold, test_fold)

  confusion_matrix <- table(test_fold$Do.you.have.Depression., predictions)


  recall <- confusion_matrix[2, 2] / sum(confusion_matrix[2, ])
  precision <- confusion_matrix[2, 2] / sum(confusion_matrix[, 2])
  f_measure <- 2 * (precision * recall) / (precision + recall)

  metrics[i, ] <- c(Recall = recall, Precision = precision, F_measure = f_measure)

  cat("Fold", i, "Recall:", recall, "Precision:", precision, "F-measure:", f_measure, "\n")
}


average_metrics <- colMeans(metrics)
cat("Average Recall:", average_metrics["Recall"], "Average Precision:",
average_metrics["Precision"], "Average F-measure:", average_metrics["F_measure"],
"\n")
```
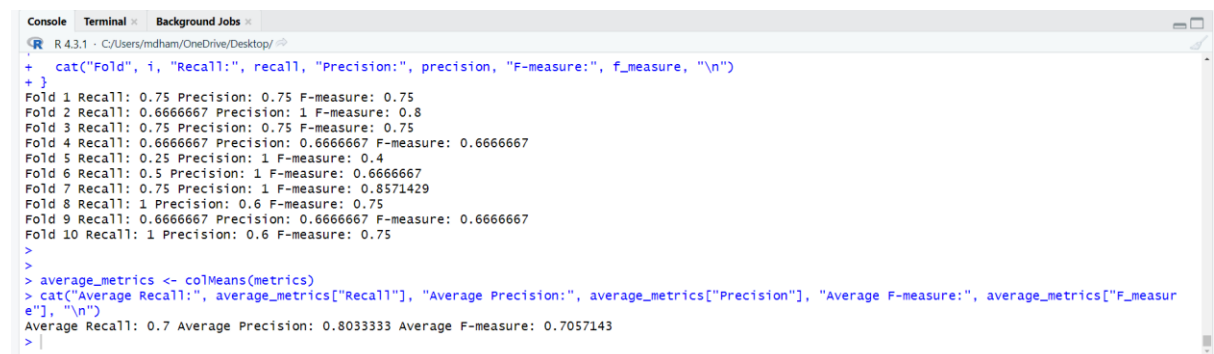


*The code shows 10 Fold (Recall, Precision, F-measure) value and it also shows the Average (Recall, Precision, F-measure) value.*