

Data Analytics & Predictive Modeling

Lab#4, Fall 2023

The dataset for this lab ("GA_Pred.csv") contains several parameters that are important for graduate admission in the US. Your objective is to use this dataset to analyze and predict the relationship between these parameters and the chance of admission.

1. (10 points) Name this dataset as "GA_Pred_C". Pick a method to identify outliers for each numerical variable, and explain this method. Document your findings in the report.

*For the following questions use "GA_Pred_C" unless otherwise mentioned.

In order to find outliers, I am first validating if the data has Nan values.

R code:

```
32 # total number of missing values
33 count_na <- sum(is.na(GA_Pred_C))
34 print(paste("Total Number of missing values in the dataframe: ", count_na))
35
36 # total number of missing values in each column
37 sapply(GA_Pred_C, function(col) sum(is.na(col)))
```

Output

```
> print(paste("Total Number of missing values in the dataframe: ", count_na))
[1] "Total Number of missing values in the dataframe: 4"
>
> sapply(GA_Pred_C, function(col) sum(is.na(col)))
      Serial.No.      GRE.Score      TOEFL.Score University.Rating      SOP      LOR      CGPA
      0           0           0           1           1           1           0
      Research    Chance.of.Admit
      1           0
```

We can see from the above output that 4 columns contains nan values. Out of 4 columns, 2 are numerical columns and 2 are categorical columns. I am imputing nan in numerical columns with mean and nan in categorical columns with mode.

R code:

```
39 numeric_columns <- c("GRE.Score", "TOEFL.Score", "SOP", "LOR", "CGPA", "Chance.of.Admit")
40 categorical_columns <- c("University.Rating", "Research")
41
42 na_columns <- colSums(is.na(GA_Pred_C[, numeric_columns]))
43 columns_with_na <- names(na_columns[na_columns > 0])
44 cat("Numeric Columns with NA values:", toString(columns_with_na), "\n")
45 for(col in columns_with_na){
46   # replace the numeric column with mean
47   GA_Pred_C[[col]][is.na(GA_Pred_C[[col]])] <- mean(GA_Pred_C[[col]], na.rm=TRUE)
48 }
49
50 na_columns <- colSums(is.na(GA_Pred_C[, categorical_columns]))
51 columns_with_na <- names(na_columns[na_columns > 0])
52 cat("Categorical Columns with NA values:", toString(columns_with_na), "\n")
53 for(col in columns_with_na){
54   # replace the numeric column with mode
55   GA_Pred_C[[col]][is.na(GA_Pred_C[[col]])] <- mode(GA_Pred_C[[col]])
56 }
57
58 sapply(GA_Pred_C, function(col) sum(is.na(col)))
```

Output:

```
> sapply(GA_Pred_C, function(col) sum(is.na(col)))
      Serial.No.      GRE.Score      TOEFL.Score University.Rating      SOP      LOR      CGPA
      0          0          0          0          0          0          0
      Research    Chance.of.Admit
      0          0
```

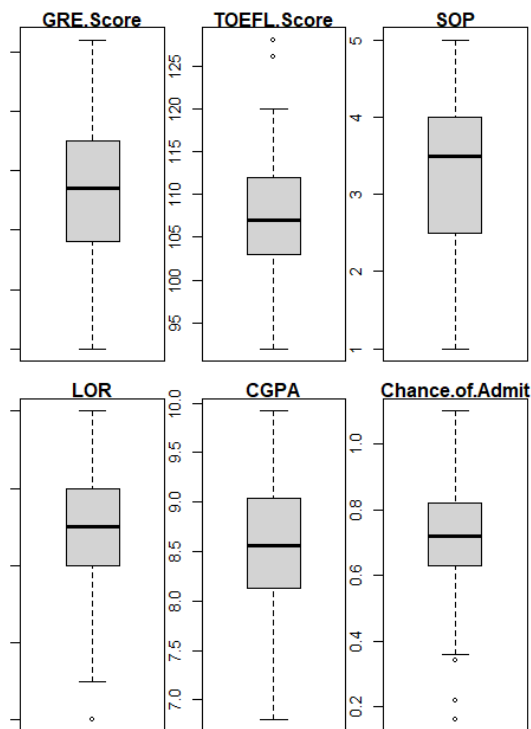
About output presents that there are no nan values in the data.

Box Plots: Box plots give an illustration of the data distribution. Points outside of the box plot's "whiskers" are known as outliers.

R code:

```
63 #Q1
64
65 par(mfrow = c(2, 3), mar=c(1,1,1,1))
66 for (col in numeric_columns) {
67   boxplot(GA_Pred_C[[col]], main = col, ylab = col)
68 }
```

Output:



Z-Score: The Z-score calculates a data point's deviation from the mean in standard deviations. An outlier is usually indicated by a z-score larger than a 3 standard deviation as threshold.

$$Z = (X - \mu) / \sigma$$

R Code:

```
70 # 3 standard deviation as threshold
71 outlier_threshold <- 3
72
73 # store outliers for each variable
74 outliers_list <- list()
75
76 for (col in numeric_columns) {
77   # z-scores
78   z_scores <- scale(GA_Pred_C[[col]])
79
80   # upper and lower thresholds
81   upper_threshold <- mean(z_scores) + outlier_threshold * sd(z_scores)
82   lower_threshold <- mean(z_scores) - outlier_threshold * sd(z_scores)
83
84   # outliers
85   outliers <- which(abs(z_scores) > outlier_threshold)
86
87   outliers_list[[col]] <- GA_Pred_C[outliers, ]
88 }
89
90 for (col in numeric_columns) {
91   cat("Outliers for", col, ":\n")
92   print(outliers_list[[col]])
93   cat("\n")
94 }
```

Output:

```
Outliers for GRE.Score :
[1] Serial.No.      GRE.Score      TOEFL.Score      University.Rating SOP      LOR      CGPA
[8] Research      Chance.of.Admit
<0 rows> (or 0-length row.names)

Outliers for TOEFL.Score :
  Serial.No.  GRE.Score  TOEFL.Score  University.Rating SOP  LOR  CGPA  Research  Chance.of.Admit
25         25        336         128           5 4.0 3.5 9.80         1         0.97
96         96        304         126           4 1.5 2.5 7.84         0         0.42

Outliers for SOP :
[1] Serial.No.      GRE.Score      TOEFL.Score      University.Rating SOP      LOR      CGPA
[8] Research      Chance.of.Admit
<0 rows> (or 0-length row.names)

Outliers for LOR :
[1] Serial.No.      GRE.Score      TOEFL.Score      University.Rating SOP      LOR      CGPA
[8] Research      Chance.of.Admit
<0 rows> (or 0-length row.names)

Outliers for CGPA :
[1] Serial.No.      GRE.Score      TOEFL.Score      University.Rating SOP      LOR      CGPA
[8] Research      Chance.of.Admit
<0 rows> (or 0-length row.names)

Outliers for Chance.of.Admit :
  Serial.No.  GRE.Score  TOEFL.Score  University.Rating SOP  LOR  CGPA  Research  Chance.of.Admit
235         235        330         113           5 5 4.0 9.31         1         0.16
274         274        312          99           1 1 1.5 8.01         1         0.22
```

From the above output, we can see that there are outliers for TOEFL score and chance of admit.

2. (10 points) Apply polynomial regression with d (degree) chosen by cross-validation in order to predict the chance of admission using GRE score. Plot the resulting polynomial fit to the data. Comment on the selected d and the plot.

In polynomial regression, an n -th degree polynomial is used to model the relationship between the independent variable (x) and the dependent variable (y).

It applies a polynomial equation to the data in order to capture more intricate relationships.

One of the most important aspects of polynomial regression is choosing the degree of the polynomial (n), which is typically done using methods like cross-validation to identify the model that strikes the best balance between complexity and data fit.

Step 1: I am splitting the data into 80-20 ratio. 80% of data is train set and 20% of the data is test set.

```
97 #Q2
98 # Sample split for training and testing
99 train_ratio <- 0.8
100 split <- sample.split(GA_Pred_C$Chance.of.Admit, SplitRatio = train_ratio)
101
102 # Create training and testing datasets
103 train_data <- subset(GA_Pred_C, split == TRUE)
104 test_data <- subset(GA_Pred_C, split == FALSE)
105
106 print(paste("Train data size: ", nrow(train_data), ncol(train_data)))
107 print(paste("Test data size: ", nrow(test_data), ncol(test_data)))
```

Output:

```
> print(paste("Train data size: ", nrow(train_data), ncol(train_data)))
[1] "Train data size: 400 9"
> print(paste("Test data size: ", nrow(test_data), ncol(test_data)))
[1] "Test data size: 100 9"
```

Step 2: I am performing cross validation to find the optimal degree for polynomial regression.

In cross-validation, the dataset is divided into several subsets, the model is trained on some of these subsets, and its performance is assessed on the remaining data. To determine the ideal degree, this process is repeated, averaging performance over various subsets.

The optimal degree of freedom is 2.

R code:

```

110 # Cross-validation to find the optimal degree for polynomial regression
111 max_degree <- 10
112 cv_error <- rep(0, max_degree)
113
114 for (degree in 1:max_degree) {
115   model_formula <- as.formula(paste("Chance.of.Admit ~ poly(GRE.Score, ", degree, ")", sep = ""))
116   cv_error[degree] <- cv.glm(train_data, glm(model_formula, data = train_data), K = 10)$delta[1]
117 }
118
119 optimal_degree <- which.min(cv_error)
120 cat("The optimal degree for the polynomial regression model is:", optimal_degree, "\n")
121

```

Output:

```

> cat("The optimal degree for the polynomial regression model is:", optimal_degree, "\n")
The optimal degree for the polynomial regression model is: 2

```

Step 3: I am training the polynomial regression model with the optimal degree of freedom and then making the predictions on the test data. Evaluating the model's performance with RMSE.

R code

```

122 # Fit the tuned polynomial regression model with the optimal degree
123 tuned_model <- lm(Chance.of.Admit ~ poly(GRE.Score, optimal_degree), data = train_data)
124
125 # Make predictions on both training and test sets
126 train_pred <- predict(tuned_model, newdata = list(GRE.Score = train_data$GRE.Score))
127 test_pred <- predict(tuned_model, newdata = list(GRE.Score = test_data$GRE.Score))
128
129 # Compute the root mean squared error (RMSE) of the predictions
130 train_rmse <- sqrt(mean((train_pred - train_data$Chance.of.Admit)^2))
131 test_rmse <- sqrt(mean((test_pred - test_data$Chance.of.Admit)^2))
132
133 cat("Train RMSE with the optimal degree for the polynomial regression model is:", train_rmse, "\n")
134 cat("Test RMSE with the optimal degree for the polynomial regression model is:", test_rmse, "\n")
135

```

Output:

```

> cat("Train RMSE with the optimal degree for the polynomial regression model is:", train_rmse, "\n")
Train RMSE with the optimal degree for the polynomial regression model is: 0.09840323
> cat("Test RMSE with the optimal degree for the polynomial regression model is:", test_rmse, "\n")
Test RMSE with the optimal degree for the polynomial regression model is: 0.07121853

```

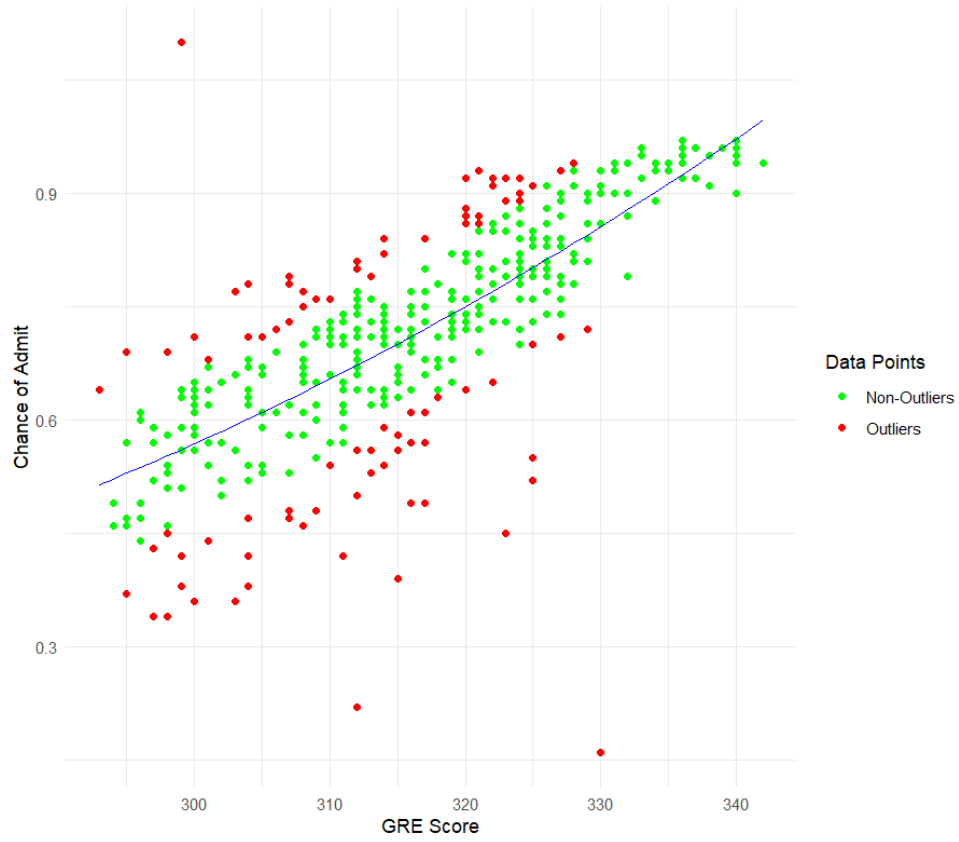
Step 4: Plotting the model fit on train data and the predictions.

The polynomial degree of two selected for the plot is suitable for the given data. The graph clearly demonstrates a parabolic relationship between the likelihood of admission and the GRE score. This relationship can be effectively represented by a quadratic polynomial regression model.

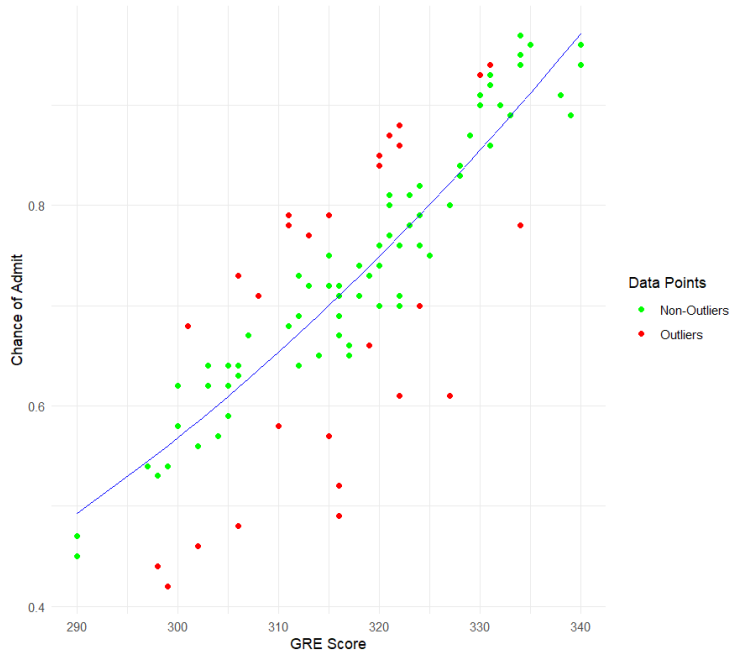
The following are some benefits of utilising a regression model with quadratic polynomials:

- Non-linear relationships between the variables can be captured by it.
- It is not too difficult to understand and use.
- Computationally, it is effective.

Training Set: Tuned Polynomial Regression Fit



Test Set: Tuned Polynomial Regression Fit



3. (20 points) Smoothing splines:

I. Use a smoothing spline with df and λ chosen by cross-validation in order to fit the chance of admission to the TOEFL scores. Plot the fitted spline on the data, and comment on the plot. (You may use `splines` package and `cv=TRUE` for cross-validation)

The non-linear relationship between TOEFL scores and the chance of admission is indicated by the fitted smoothing spline displayed in the image. The initial concavity of the curve suggests that the likelihood of admission rises quickly with a higher TOEFL score. At higher TOEFL scores, the curve starts to flatten out, indicating that the likelihood of admission reaches a plateau.

The fitted smoothing spline offers a decent fit for the data. In addition to accounting for data outliers, it captures the non-linear relationship between TOEFL scores and admission chance.

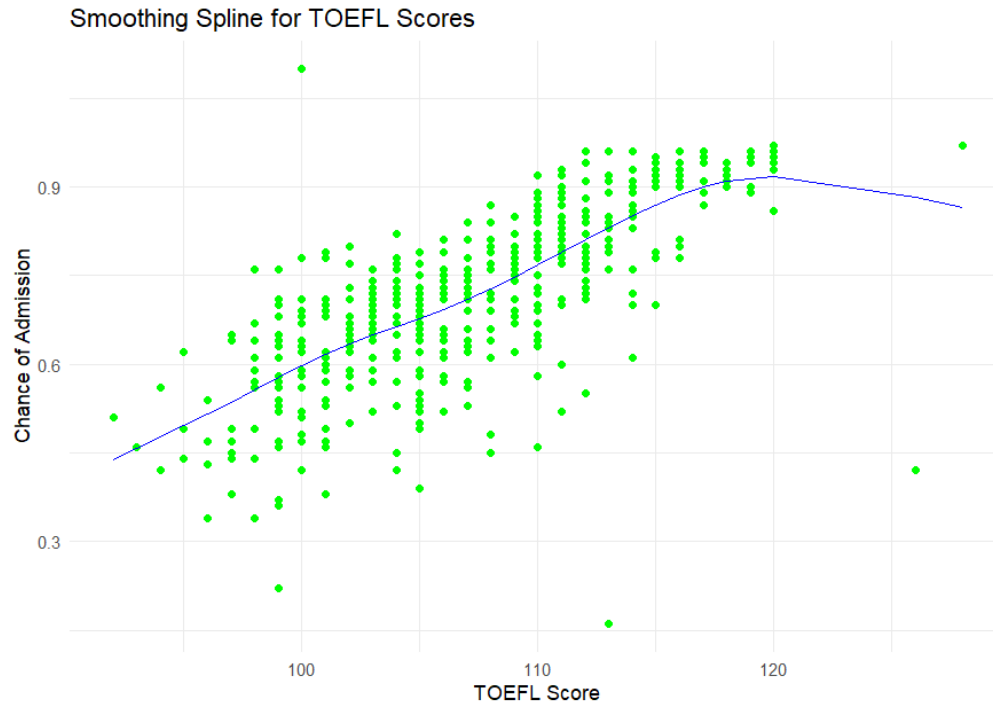
The above graph's smoothing parameter, 0.0126916, is a reasonable choice. Without overfitting, the curve is smooth and provides a good fit to the data.

The trade-off between smoothing out the noise and fitting the data is managed by the smoothing parameter. A smoother curve will be produced by increasing the smoothing parameter, but the data may lose some significant details as a result. A wigglier curve will be produced by a smaller smoothing parameter, but this may be overfitting the data's noise.

R code:

```
167 #Q3
168
169 # applying smooth spline
170 smoothend_spline <- smooth.spline(GA_Pred_C$TOEFL.Score,
171                                   GA_Pred_C$Chance.of.Admit, cv = TRUE)
172
173 # Create a data frame for plotting
174 plot_data <- data.frame(TOEFL.Score = GA_Pred_C$TOEFL.Score,
175                          Chance.of.Admit = GA_Pred_C$Chance.of.Admit)
176
177 # Plot the data and the fitted spline using ggplot2
178 ggplot(plot_data, aes(x = TOEFL.Score, y = Chance.of.Admit)) +
179   geom_point(col = "green") +
180   geom_line(aes(y = fitted(smoothend_spline)), color = "blue") +
181   labs(x = "TOEFL Score", y = "Chance of Admission") +
182   ggtitle("Smoothing Spline for TOEFL Scores") +
183   theme_minimal() +
184   theme(legend.position = "topright") +
185   guides(col = guide_legend(title = "Smoothing Spline"))
186
187 # Print the chosen values of df and lambda
188 cat("Degrees of Freedom (df):", smoothend_spline$df, "\n")
189 cat("Smoothing Parameter (lambda):", smoothend_spline$lambda, "\n")
190
```

Output:



```
> cat("Degrees of Freedom (df):", smoothend_spline$df, "\n")
Degrees of Freedom (df): 5.32607
> cat("Smoothing Parameter (lambda):", smoothend_spline$lambda, "\n")
Smoothing Parameter (lambda): 0.0126916
```

II. Discuss how the values of λ and df affect smoothing splines.

The main function of λ is to manage the trade-off between obtaining a smooth curve and closely fitting the data. A smoother curve is associated with a larger λ because it penalizes deviations from a simple function or straight line.

Lower values of λ enable the spline to closely track the data points, which may cause noise to be captured and overfitting to occur. Greater generalization is facilitated by larger values of λ , which lessen the influence of individual data points and prevent overfitting.

A smoothing spline's degrees of freedom show how flexible the model is. The spline can be more flexible and capture minute details in the data with a higher df . On the other hand, a smoother, less flexible curve is produced by a lower df .

Increasing df usually makes the model more accurate at fitting the data, but it can also cause overfitting. Achieving a balance between preserving simplicity and identifying the underlying trend is crucial.

Selecting the right value for λ and df requires striking a compromise between preventing needless complexity and providing a good fit for the data. The best values can be found with the aid of model selection methods such as cross-validation.

A small and a high df could cause overfitting and noise capture, while a large λ and a low df could lead to an oversimplified model that misses significant patterns.

By resolving a penalized least squares issue, smoothing splines are fitted. Higher df values result in an increase in computational cost, while larger λ values can simplify the optimization problem and increase its computational efficiency.

4. (20 points) Assume that we are interested in predicting whether a student's chance of acceptance is above 0.72. Split the dataset into 80% train and 20% test. Use the train set and fit a polynomial logistic regression model using GRE score with the optimal d obtained in question 2 and CGPA with d=1. Apply the fitted model to the test dataset to predict the chance of admission. Plot the predicted probabilities versus actual probabilities and discuss the performance of the fitted model. (Use type="response" in predict()).

Step 1: I am using the optimal degree for GRE as 2 (obtained in q2) and for CGPA is 1. I am fitting the logistic regression model on train data and testing on test data. Chance of acceptance above 0.72 is 1 and below or equal to 0.72 is 0.

R code:

```
192 #Q4
193 optimal_degree_GRE <- 2 #optimal degree from Q2
194 optimal_degree_CGPA <- 1
195
196 # Create formula for logistic regression
197 formula <- as.formula(paste("Chance.of.Admit > 0.72 ~ poly(GRE.Score, ",
198                               optimal_degree_GRE, ") + poly(CGPA, ", optimal_degree_CGPA, ")"))
199
200 # Fit logistic regression model
201 logistic_model <- glm(formula, data = train_data, family = binomial)
202
203 # Predict on the test set
204 predicted_probabilities <- predict(logistic_model, newdata = test_data, type = "response")
205
206 test_rmse <- sqrt(mean((predicted_probabilities - test_data$Chance.of.Admit)^2))
207 cat(test_rmse)
208
```

Output

```
> cat(test_rmse)
0.3456252
```

Step 2: To evaluate the model, I am first using 0.50 as the threshold on predicted probabilities which is the standard way of doing.

R code:

```

215 # Evaluate model performance
216 threshold <- 0.5
217 predicted_labels <- ifelse(predicted_probabilities > threshold, 1, 0)
218 actual_labels <- ifelse(test_data$Chance.of.Admit > 0.72, 1, 0)
219
220 # Confusion matrix
221 conf_matrix <- table(Actual = actual_labels, Predicted = predicted_labels)
222 print("Confusion Matrix:")
223 print(conf_matrix)
224
225 # Calculate accuracy
226 accuracy <- sum(diag(conf_matrix)) / sum(conf_matrix)
227 cat("Accuracy:", accuracy, "\n")
228

```

Output:

```

> print("Confusion Matrix:")
[1] "Confusion Matrix:"
> print(conf_matrix)
      Predicted
Actual 0  1
      0 39  9
      1  9 43

> cat("Accuracy:", accuracy, "\n")
Accuracy: 0.82

```

Step 3: I am using 0.72 as the threshold on predicted probabilities.

R code:

```

229 # Evaluate model performance
230 threshold <- 0.72
231 predicted_labels <- ifelse(predicted_probabilities > threshold, 1, 0)
232 actual_labels <- ifelse(test_data$Chance.of.Admit > 0.72, 1, 0)
233
234 # Confusion matrix
235 conf_matrix <- table(Actual = actual_labels, Predicted = predicted_labels)
236 print("Confusion Matrix:")
237 print(conf_matrix)
238
239 # Calculate accuracy
240 accuracy <- sum(diag(conf_matrix)) / sum(conf_matrix)
241 cat("Accuracy:", accuracy, "\n")
242

```

Output:

```

[1] "Confusion Matrix:"
> print(conf_matrix)
      Predicted
Actual 0  1
      0 45  3
      1 12 40

```

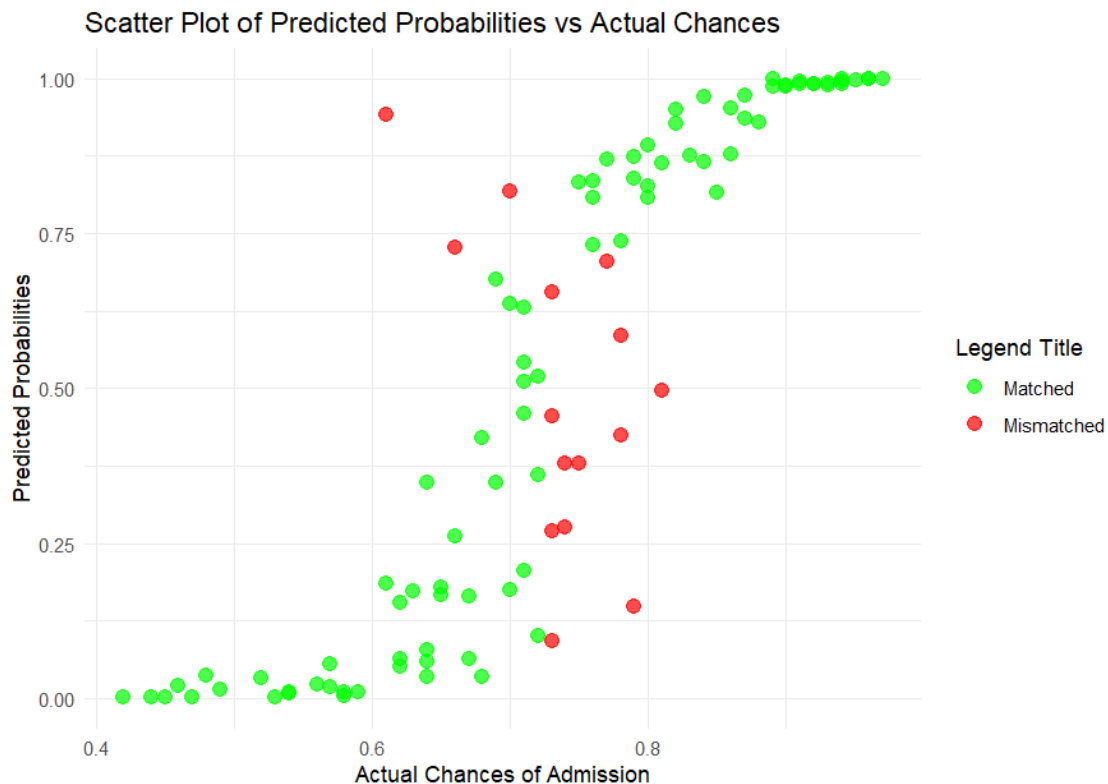
```
> cat("Accuracy:", accuracy, "\n")
Accuracy: 0.85
```

Step 4: Plotting the fitted model.

R code:

```
243 # Create a data frame for plotting
244 plot_data <- data.frame(
245   Actual_Chance = test_data$Chance.of.Admit,
246   Predicted_Probabilities = predicted_probabilities,
247   Color = ifelse((predicted_probabilities > 0.72 & test_data$Chance.of.Admit <= 0.72) |
248                 (predicted_probabilities <= 0.72 & test_data$Chance.of.Admit > 0.72), "red", "green")
249 )
250
251 # Plot the scatter plot with color-coded points for mismatched probabilities
252 library(ggplot2)
253
254 ggplot(plot_data, aes(x = Actual_Chance, y = Predicted_Probabilities, color = Color)) +
255   geom_point(alpha = 0.7, size = 3) + # Fix: alpha should be between 0 and 1
256   ggtitle("Scatter Plot of Predicted Probabilities vs Actual Chances") +
257   xlab("Actual Chances of Admission") +
258   ylab("Predicted Probabilities") +
259   scale_color_manual(values = c("green" = "green", "red" = "red"), name = "Legend Title",
260                     labels = c("Matched", "Mismatched")) +
261   theme_minimal()
```

Output:



The graph displays a logistic regression model's scatter plot of expected probabilities versus actual admission chances. The diagonal line occupies the majority of the scatter plot, suggesting that the model is well-calibrated. This indicates that the expected and actual chances of admission are fairly close to each other.

A small number of the data points are situated below the diagonal. This indicates that these students' actual chances of admission were higher than the model's predicted probability of admission. Several things could be the cause of this, including:

- It's possible that the model did not account for every significant factor influencing admissions decisions.
- It's possible that the model is overfitting the training set, which would explain why it performs poorly when applied to fresh data.
- Predictions may be off if there is noise in the data.

5. (20 points) Split the dataset into 70% train and 30% test. Fit a MARS on the training and use that model to predict the test data set. Plot the predicted probabilities versus actual probabilities and discuss the performance of the fitted model. (Use “mda” or “earth” package)

Step 1: I am splitting the data with 70-30 % ratio.

R code:

```
262 #Q5
263
264 # Sample split for training and testing
265 train_ratio <- 0.8
266 split <- sample.split(GA_Pred_C$Chance.of.Admit, SplitRatio = train_ratio)
267
268 # Create training and testing datasets
269 train_data_Q5 <- subset(GA_Pred_C, split == TRUE)
270 test_data_Q5 <- subset(GA_Pred_C, split == FALSE)
271
272 print(paste("Train data size: ", nrow(train_data_Q5), ncol(train_data_Q5)))
273 print(paste("Test data size: ", nrow(test_data_Q5), ncol(test_data_Q5)))
274
```

Output:

```
> print(paste("Train data size: ", nrow(train_data_Q5), ncol(train_data_Q5)))
[1] "Train data size: 400 9"
> print(paste("Test data size: ", nrow(test_data_Q5), ncol(test_data_Q5)))
[1] "Test data size: 100 9"
```

Step 2: Fitting the MARS on the train data and predicting on test data. Evaluating the model's performance by considering 0.50 as threshold.

R code:

```

275 mars_model <- earth(Chance.of.Admit ~ GRE.Score + TOEFL.Score + CGPA, data = train_data_Q5)
276
277 # Predict on the test set
278 predicted_probabilities <- predict(mars_model, newdata = test_data_Q5)
279
280 # Evaluate model performance
281 threshold <- 0.5
282 predicted_labels <- ifelse(predicted_probabilities > threshold, 1, 0)
283 actual_labels <- ifelse(test_data_Q5$Chance.of.Admit > 0.72, 1, 0)
284
285 # Confusion matrix
286 conf_matrix <- table(Actual = actual_labels, Predicted = predicted_labels)
287 print("Confusion Matrix:")
288 print(conf_matrix)
289
290 # Calculate accuracy
291 accuracy <- sum(diag(conf_matrix)) / sum(conf_matrix)
292 cat("Accuracy:", accuracy, "\n")

```

Output:

```

> print(conf_matrix)
      Predicted
Actual 0 1
      0 5 43
      1 0 52
\ |

> cat("Accuracy:", accuracy, "\n")
Accuracy: 0.57

```

Step 3: Evaluating the model's performance by considering 0.72 as threshold.

R code:

```

294 threshold <- 0.72
295 predicted_labels <- ifelse(predicted_probabilities > threshold, 1, 0)
296 actual_labels <- ifelse(test_data_Q5$Chance.of.Admit > 0.72, 1, 0)
297
298 # Confusion matrix
299 conf_matrix <- table(Actual = actual_labels, Predicted = predicted_labels)
300 print("Confusion Matrix:")
301 print(conf_matrix)
302
303 # Calculate accuracy
304 accuracy <- sum(diag(conf_matrix)) / sum(conf_matrix)
305 cat("Accuracy:", accuracy, "\n")

```

Output:

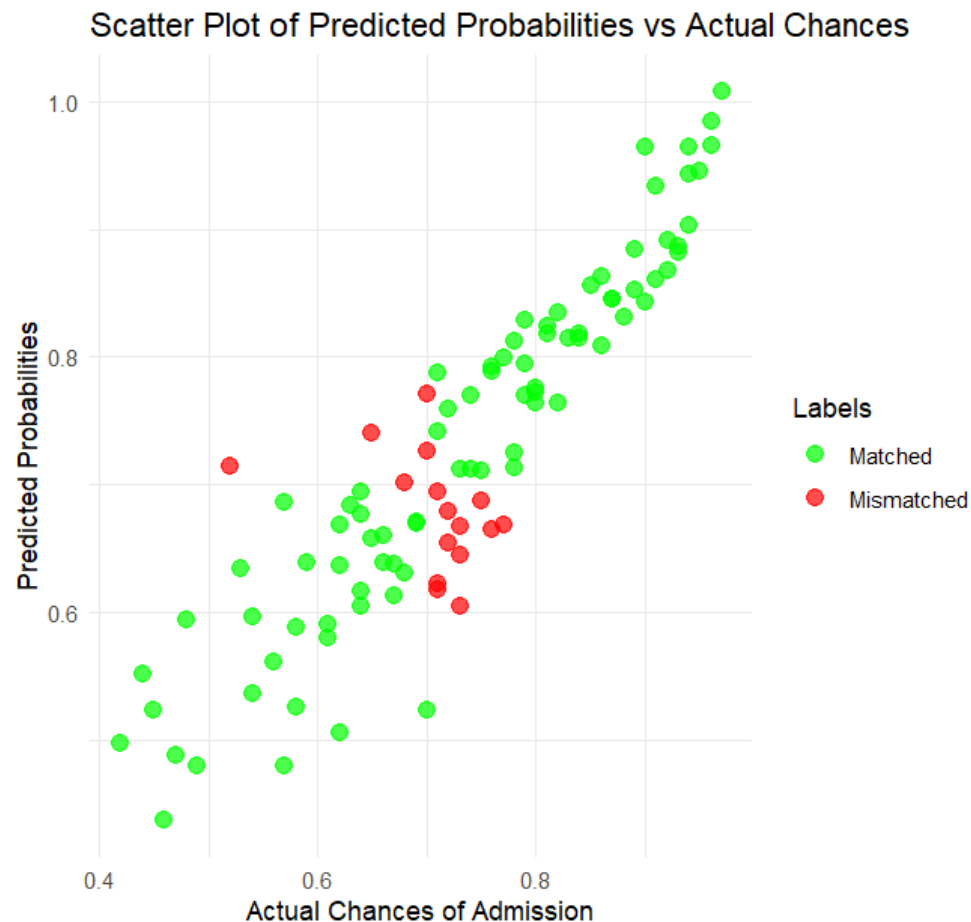
```

- -
> print(conf_matrix)
      Predicted
Actual 0 1
      0 42 6
      1 10 42

```

```
> cat("Accuracy:", accuracy, "\n")
Accuracy: 0.84
```

Step 4: Plotting the predicted probabilities versus actual probabilities



The data seem to fit well with the MARS model. Its accuracy is relatively high, and it is well-calibrated.

A small number of the data points are situated below the diagonal. This indicates that these students' actual chances of admission were higher than the model's predicted probability of admission.

6. (20 points) Open question: Split the dataset into train and test. Fit a GAM on the training dataset and use the fitted model to predict the response variable on the test set. Report and interpret the results, and provide plots in order to discuss the performance of the model and its outputs.

Step 1: Fitting the GAM model and making predictions.

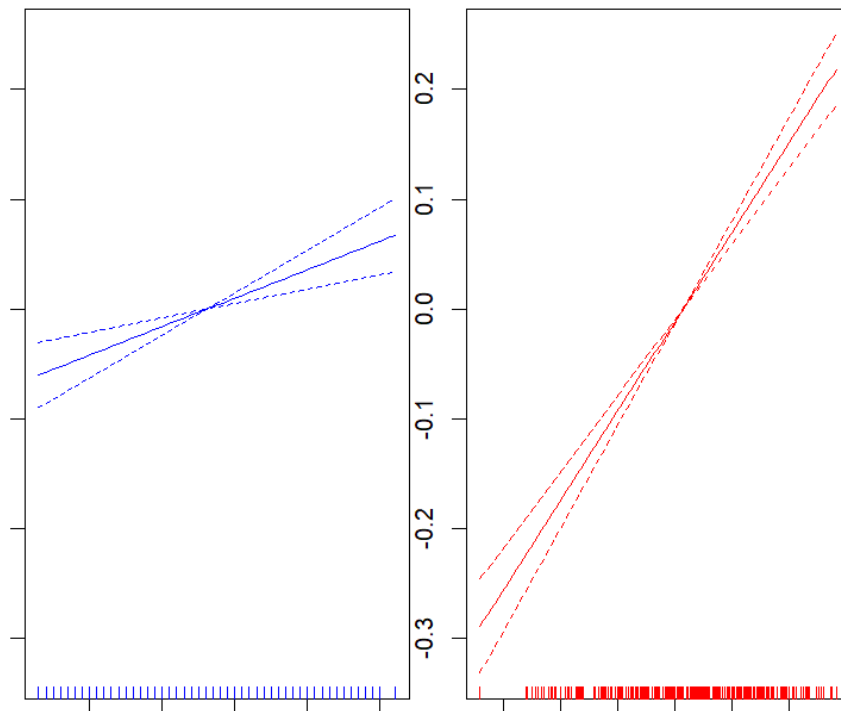
R code:

```

328 # Fitting GAM model on training data
329 gam_model <- gam(Chance.of.Admit ~ s(GRE.Score) + s(CGPA), data = train_data)
330
331 # Predict on the test set
332 predicted_probabilities <- predict(gam_model, newdata = test_data)
333
334 par(mfrow = c(1, 2))
335 plot(gam_model)
336
337 names(plot_data) <- c("Actual_Chance", "Predicted_Probabilities", "Color")
338

```

Output:



Blue is the plot for GRE.Score and red graph is the plot for CGPA.

Step 2: Evaluating the model's performance with threshold 0.50.

R code:

```

339 # Evaluate model performance
340 threshold <- 0.5
341 predicted_labels <- ifelse(predicted_probabilities > threshold, 1, 0)
342 actual_labels <- ifelse(test_data$Chance.of.Admit > 0.72, 1, 0)
343
344 # Confusion matrix
345 conf_matrix <- table(Actual = actual_labels, Predicted = predicted_labels)
346 print("Confusion Matrix:")
347 print(conf_matrix)
348
349 # Calculate accuracy
350 accuracy <- sum(diag(conf_matrix)) / sum(conf_matrix)
351 cat("Accuracy:", accuracy, "\n")
352

```

Output:

```

[1] "Confusion Matrix:"
> print(conf_matrix)
      Predicted
Actual 0 1
      0 6 42
      1 0 52
>
> # Calculate accuracy
> accuracy <- sum(diag(conf_matrix)) / sum(conf_matrix)
> cat("Accuracy:", accuracy, "\n")
Accuracy: 0.58
> |

```

Step 3: Evaluating the model's performance with threshold 0.72.

R code:

```

353 threshold <- 0.72
354 predicted_labels <- ifelse(predicted_probabilities > threshold, 1, 0)
355 actual_labels <- ifelse(test_data$Chance.of.Admit > 0.72, 1, 0)
356
357 # Confusion matrix
358 conf_matrix <- table(Actual = actual_labels, Predicted = predicted_labels)
359 print("Confusion Matrix:")
360 print(conf_matrix)
361
362 # Calculate accuracy
363 accuracy <- sum(diag(conf_matrix)) / sum(conf_matrix)
364 cat("Accuracy:", accuracy, "\n")
365

```

Output:


```

> # Confusion matrix
> conf_matrix <- table(Actual = actual_labels, Predicted = predicted_labels)
> print("Confusion Matrix:")
[1] "Confusion Matrix:"
> print(conf_matrix)
      Predicted
Actual 0 1
      0 41 7
      1 9 43
>
> # Calculate accuracy
> accuracy <- sum(diag(conf_matrix)) / sum(conf_matrix)
> cat("Accuracy:", accuracy, "\n")
Accuracy: 0.84

```

Step 4: Plotting the graph.

```

366 plot_data <- data.frame(
367   Actual_Chance = test_data$Chance.of.Admit,
368   Predicted_Probabilities = predicted_probabilities,
369   Color = ifelse((predicted_probabilities > 0.7 & test_data$Chance.of.Admit <= 0.7) |
370                 (predicted_probabilities <= 0.7 & test_data$Chance.of.Admit > 0.7), "red", "green")
371 )
372 # Plot the scatter plot with color-coded points
373 ggplot(plot_data, aes(x = Actual_Chance, y = Predicted_Probabilities, color = Color)) +
374   geom_point(alpha = 0.7, size = 3) +
375   ggtitle("Scatter Plot of Predicted Probabilities vs Actual Chances") +
376   xlab("Actual Chances of Admission") +
377   ylab("Predicted Probabilities") +
378   scale_color_manual(values = c("green", "red"), name = "Labels",
379                     labels = c("Matched", "Mismatched")) +
380   theme_minimal()
381

```

