



ILLINOIS TECH

HOMEWORK 2

EDUARDO GALEOTE ESCALERA

A20552496

CSP 571

Data Preparation and Analysis

September 2023

Contents

1	Recitation Exercises	3
1.1	Chapter 4	3
1.1.1	Excercise 4	3
1.1.2	Excercise 6	4
1.1.3	Excercise 7	4
1.1.4	Excercise 9	4
1.2	Chapter 5	5
1.2.1	Excercise 2	5
1.2.2	Excercise 3	6
2	Practicum Problems	6
2.1	Problem 1	6
2.2	Problem 2	8
2.3	Problem 3	9
2.4	Problem 4	11

1 Recitation Exercises

1.1 Chapter 4

1.1.1 Exercise 4

When the number of features p is large, there tends to be a deterioration in the performance of KNN and other local approaches that perform prediction using only observations that are near the test observation for which a prediction must be made. This phenomenon is known as the curse of dimensionality, and it ties into the fact that non-parametric approaches often perform poorly when p is large. We will now investigate this curse.

- a) Suppose that we have a set of observations, each with measurements on $p = 1$ feature, X . We assume that X is uniformly (evenly) distributed on $[0, 1]$. Associated with each observation is a response value. Suppose that we wish to predict a test observation's response using only observations that are within 10 % of the range of X closest to that test observation. For instance, in order to predict the response for a test observation with $X = 0.6$, we will use observations in the range $[0.55, 0.65]$. On average, what fraction of the available observations will we use to make the prediction?

Assuming $x \in [0.05, 0.95]$, then intervals: $[x - 0.05, x + 0.05]$, so $length = 0.1$. On average 10% of the observations would be available to make a prediction for the test observation.

- b) Now suppose that we have a set of observations, each with measurements on $p = 2$ features, X_1 and X_2 . We assume that (X_1, X_2) are uniformly distributed on $[0, 1] \times [0, 1]$. We wish to predict a test observation's response using only observations that are within 10 % of the range of X_1 and within 10 % of the range of X_2 closest to that test observation. For instance, in order to predict the response for a test observation with $X_1 = 0.6$ and $X_2 = 0.35$, we will use observations in the range $[0.55, 0.65]$ for X_1 and in the range $[0.3, 0.4]$ for X_2 . On average, what fraction of the available observations will we use to make the prediction?

Assuming $x \in [0.05, 0.95]$, $x1_{length} \times x2_{length} = 0.01$. Therefore, 1% of the available observations would be used to make a prediction.

- c) Now suppose that we have a set of observations on $p = 100$ features. Again the observations are uniformly distributed on each feature, and again each feature ranges in value from 0 to 1. We wish to predict a test observation's response using observations within the 10 % of each feature's range that is closest to that test observation. What fraction of the available observations will we use to make the prediction?

If $p=100$; $0.1^p \times 100 = 0.1^{100} \times 100\%$ of the observations are available.

- d) Using your answers to parts (a)–(c), argue that a drawback of KNN when p is large is that there are very few training observations “near” any given test observation.

As the number of predictors increase, the fraction of observations available to make a prediction is reduced exponentially.

- e) Now suppose that we wish to make a prediction for a test observation by creating a p -dimensional hypercube centered around the test observation that contains, on average, 10 % of the training observations. For $p = 1, 2$, and 100, what is the length of each side of the hypercube? Comment on your answer.

- $p = 1; d(length) = 0.1^{1/1} = 0.1$

- $p = 2; d(\text{length}) = 0.1^{1/2} = 0.32$
- $p = 100; d(\text{length}) = 0.1^{1/100} = 0.977$

As p increases the side length converges to 1, and this shows that the hypercube centered around the test observation with 10% of the test observation needs to be nearly the same size as the hypercube with all the observations.

1.1.2 Exercise 6

Suppose we collect data for a group of students in a statistics class with variables $X1 = \text{hours studied}$, $X2 = \text{undergrad GPA}$, and $Y = \text{receive an A}$. We fit a logistic regression and produce estimated coefficient, $\hat{\beta}_0 = 6$, $\hat{\beta}_1 = 0.05$, $\hat{\beta}_2 = 1$.

- a) Estimate the probability that a student who studies for 40 h and has an undergrad GPA of 3.5 gets an A in the class.

$$P(X) = \frac{\exp(\hat{\beta}_0 + \hat{\beta}_1 X_1 + \hat{\beta}_2 X_2)}{1 + \exp(\hat{\beta}_0 + \hat{\beta}_1 X_1 + \hat{\beta}_2 X_2)}$$

$$P(X) = \frac{\exp(-0.5)}{1 + \exp(-0.5)} = 0.38$$

- b) How many hours would the student in part (a) need to study to have a 50 % chance of getting an A in the class?

$$\log\left(\frac{P(X)}{1 - P(X)}\right) = \hat{\beta}_0 + \hat{\beta}_1 X_1 + \hat{\beta}_2 X_2$$

$$\log\left(\frac{0.5}{1 - 0.5}\right) = -6 + 0.05X_1 + 3.5$$

$$X_1 = 50 \text{ hours.}$$

1.1.3 Exercise 7

Suppose that we wish to predict whether a given stock will issue a dividend this year (“Yes” or “No”) based on X , last year’s percent profit. We examine a large number of companies and discover that the mean value of X for companies that issued a dividend was $\bar{X} = 10$, while the mean for those that didn’t was $\bar{X} = 0$. In addition, the variance of X for these two sets of companies was $\sigma^2 = 36$. Finally, 80 % of companies issued dividends. Assuming that X follows a normal distribution, predict the probability that a company will issue a dividend this year given that its percentage profit was $X = 4$ last year.

$$P(Y = \text{yes} | X = 4) = \frac{\pi_{\text{yes}} f_{\text{yes}}(x)}{\sum_{l=1}^K \pi_l f_l(x)} = \frac{\pi_{\text{yes}} \exp(-1/2\sigma^2(x - \mu_{\text{yes}})^2)}{\sum_{l=1}^K \pi_l \exp(-1/2\sigma^2(x - \mu_l)^2)}$$

$$P(Y = \text{yes} | X = 4) = \frac{0.8 \times \exp(-0.5)}{0.8 \times \exp(-0.5) + 0.2 \times \exp(-16/72)}$$

$$P(Y = \text{yes} | X = 4) = 0.75$$

1.1.4 Exercise 9

This problem has to do with odds.

- a) On average, what fraction of people with an odds of 0.37 of defaulting on their credit card payment will in fact default?

$$\text{Odds} = \frac{P(X)}{1 - P(X)}$$

$$P(X) = \frac{0.37}{1.37} = 0.27$$

27% of people with odds of 0.37 will default.

- b) Suppose that an individual has a 16% chance of defaulting on her credit card payment. What are the odds that she will default?

$$\text{Odds} = \frac{0.16}{1 - 0.16} = 0.19$$

1.2 Chapter 5

1.2.1 Exercise 2

We will now derive the probability that a given observation is part of a bootstrap sample. Suppose that we obtain a bootstrap sample from a set of n observations.

- a) What is the probability that the first bootstrap observation is not the j th observation from the original sample? Justify your answer.

When using random sampling with replacement to obtain a bootstrap sample from n observations: If $n=3$, with observations: 1,2,3 we have: Total permutations: $n^n = 3^3 = 27$. Total times an observation will appear in a particular order (say as the first observation): $n^{n-1} = 3^2 = 9$. So, the probability 1st bootstrap observation is the j th observation is: $1/n$. Therefore, the probability the 1st bootstrap observation is not the j th observation is: $1 - 1/n$.

- b) What is the probability that the second bootstrap observation is not the j th observation from the original sample?

$$1 - 1/n$$

- c) Argue that the probability that the j th observation is not in the bootstrap sample is $(1 - \frac{1}{n})^n$. Total times an observation will not appear in a bootstrap sample: $(n-1)^n$. Therefore, the probability that an observation is not in the bootstrap sample: $\frac{(n-1)^n}{n^n} = (1 - 1/n)^n$.

- d) When $n = 5$, what is the probability that the j th observation is in the bootstrap sample?

$$p = 1 - (1 - 1/5)^5 = 0.672$$

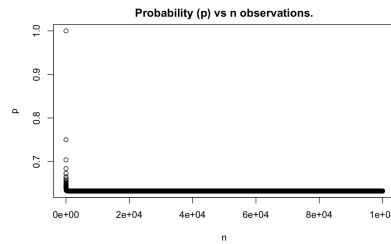
- e) When $n = 100$, what is the probability that the j th observation is in the bootstrap sample?

$$p = 1 - (1 - 1/100)^{100} = 0.634$$

- f) When $n = 10,000$, what is the probability that the j th observation is in the bootstrap sample?

$$p = 1 - (1 - 1/10000)^{10000} = 0.632$$

- g) Create a plot that displays, for each integer value of n from 1 to 100,000, the probability that the j th observation is in the bootstrap sample. Comment on what you observe.



Probability decreases rapidly to an asymptote around 0.63 as n increases.

- h) We will now investigate numerically the probability that a bootstrap sample of size $n = 100$ contains the j th observation. Here $j = 4$. We repeatedly create bootstrap samples, and each time we record whether or not the fourth observation is contained in the bootstrap sample. 0.6338. The result is similar to the value calculated using the formula in e.

1.2.2 Exercise 3

We now review k -fold cross-validation.

- a) Explain how k -fold cross-validation is implemented.
In k -fold cross validation, a training set is divided into k groups of equal size. The first group is treated as the validation set, and the model is fit on the remaining $k-1$ groups. The MSE is calculated using the validation set. This procedure is repeated k times and on each occasion the validation set and training sets will be different than the previous one. We then take the average of all the MSE's as the final MSE.
- b) What are the advantages and disadvantages of k -fold cross-validation relative to: i. The validation set approach? and ii. LOOCV?
- Less variance in test error estimate. A more accurate test error, as entire dataset is used.
 - Computational advantage, as less computing resources required than LOOCV. Only k models need to be fitted, unlike with LOOCV where n models need fitting. Higher bias than LOOCV, as fewer observations are used, but tends to have lower variance.

2 Practicum Problems

2.1 Problem 1

Load the abalone sample dataset from the UCI Machine Learning Repository (abalone.data) into R using a dataframe. Remove all observations in the Infant category, keeping the Male/Female classes. Using the caret package, use createDataPartition to perform an 80/20 test-train split (80% training and 20% testing). Fit a logistic regression using all feature variables via glm, and observe which predictors are relevant. Do the confidence intervals for the predictors contain 0 within the range? How does this relate to the null hypothesis? Use the confusionMatrix function in caret to observe testing results (use a 50% cutoff to tag Male/Female) - how does the accuracy compare to a random classifier ROC curve? Use the corrplot package to plot correlations between the predictors. How does this help explain the classifier performance?

```

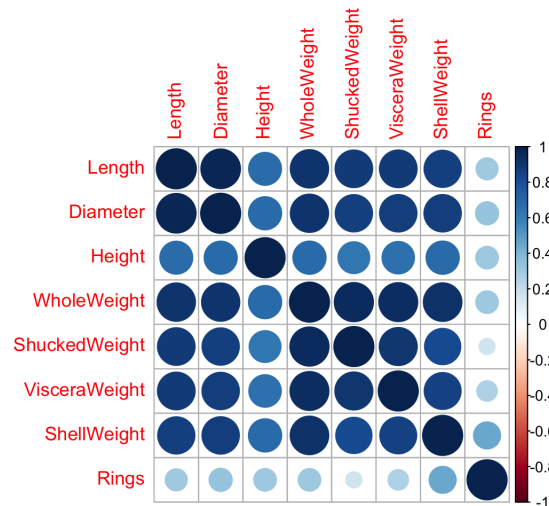
1 # Load required libraries
2 library(caret)
3 library(corrplot)
4 library(pROC)
5
6 # Load the abalone dataset from the UCI Machine Learning Repository
7 url <- "https://archive.ics.uci.edu/ml/machine-learning-databases/abalone/
  abalone.data"
8 abalone <- read.csv(url, header = FALSE)
9
10 # Add column names to the dataset
11 colnames(abalone) <- c("Sex", "Length", "Diameter", "Height", "WholeWeight", "
  ShuckedWeight", "VisceraWeight", "ShellWeight", "Rings")
12
13 # Filter out observations in the "Infant" category
14 abalone <- abalone[abalone$Sex %in% c("M", "F"), ]
15
16 # Recode "Sex" to a binary variable
17 abalone$Sex <- ifelse(abalone$Sex == "M", 1, 0)
18
19 # Perform an 80/20 test-train split
20 set.seed(123) # for reproducibility
21 trainIndex <- createDataPartition(abalone$Sex, p = 0.8, list = FALSE)
22 trainData <- abalone[trainIndex, ]
23 testData <- abalone[-trainIndex, ]
24
25 # Convert "Sex" to a factor with levels 0 and 1
26 trainData$Sex <- factor(trainData$Sex, levels = c(0, 1))
27 testData$Sex <- factor(testData$Sex, levels = c(0, 1))
28
29 # Fit a logistic regression model
30 model <- glm(Sex ~ ., data = trainData, family = binomial)
31
32 # Summary of the logistic regression model
33 summary(model)
34
35 # Confidence intervals for the predictors
36 confint(model)
37
38 # Check if confidence intervals contain 0
39 # If the interval contains 0, it suggests that the predictor is not
  statistically significant
40 # This relates to the null hypothesis, where the null hypothesis is that the
  predictor's coefficient is equal to 0
41 # If the interval contains 0, we fail to reject the null hypothesis, indicating
  the predictor may not be relevant.
42
43 # Predict on the test data using the logistic regression model
44 predictions <- predict(model, newdata = testData, type = "response")
45
46 # Convert probabilities to class predictions using a 50% cutoff
47 predicted_classes <- ifelse(predictions > 0.5, "1", "0")
48
49 predicted_classes <- factor(predicted_classes, levels = c(0, 1))
50
51 # Create a confusion matrix
52 confusion_matrix <- confusionMatrix(predicted_classes, testData$Sex)
53
54 # Print the confusion matrix
55 print(confusion_matrix)
56
57 # Calculate accuracy
58 accuracy <- confusion_matrix$overall["Accuracy"]

```

```

59 cat("Accuracy:", accuracy, "\n")
60
61 # ROC Curve
62 roc_curve <- roc(testData$Sex, predictions)
63 plot(roc_curve)
64
65 # Correlation plot
66 correlation_matrix <- cor(trainData[, -1]) # Exclude the "Sex" column
67 corplot(correlation_matrix, method = "circle")

```



High correlations between predictors indicate multicollinearity. It can be problematic for logistic regression models because it can make it challenging to isolate the individual effects of predictors on the response variable. Correlation plots can help you identify which predictors are strongly correlated with the response variable. Variables with strong correlations with the response are likely to be more informative and may have a greater impact on classifier performance.

Accuracy of the logistic regression model is 0.566 while the ROC curve's accuracy is 0.497.

2.2 Problem 2

Load the mushroom sample dataset from the UCI Machine Learning Repository (*agaricus-lepiota.data*) into R using a dataframe (Note: There are missing values with a ? character, you will have to explain your handling of these). Create a Naive Bayes classifier using the *e1071* package, using the sample function to split the data between 80% for training and 20% for testing. With the target class of interest being edible mushrooms, calculate the accuracy of the classifier both in-training and in-test. Use the table function to create a confusion matrix of predicted vs. actual classes - how many false positives did the model produce?

```

1 # Load the required packages
2 library(e1071)
3
4 # Load the mushroom dataset from the UCI Machine Learning Repository
5 # Specify the column names as per dataset documentation.
6 column_names <- c("class", "cap_shape", "cap_surface", "cap_color", "bruises", "
  odor", "gill_attachment", "gill_spacing", "gill_size", "gill_color", "stalk_
  shape", "stalk_root", "stalk_surface_above_ring", "stalk_surface_below_ring"
  , "stalk_color_above_ring", "stalk_color_below_ring", "veil_type", "veil_
  color", "ring_number", "ring_type", "spore_print_color", "population", "
  habitat")

```



```

7
8 # Load the dataset, treating "?" as missing values
9 mushroom_data <- read.csv("agaricus-lepiota.data", header = FALSE, col.names =
    column_names, na.strings = "?")
10
11 # Handle missing values by removing rows with any missing values
12 mushroom_data <- na.omit(mushroom_data)
13
14 # Split the data into training (80%) and testing (20%) sets
15 set.seed(123) # for reproducibility
16 train_indices <- sample(1:nrow(mushroom_data), 0.8 * nrow(mushroom_data))
17 train_data <- mushroom_data[train_indices, ]
18 test_data <- mushroom_data[-train_indices, ]
19
20 # Create a Naive Bayes classifier
21 nb_classifier <- naiveBayes(class ~ ., data = train_data)
22
23 # Predict the classes for training and testing sets
24 train_predictions <- predict(nb_classifier, train_data)
25 test_predictions <- predict(nb_classifier, test_data)
26
27 # train_predictions to char vector
28 train_predictions <- as.character(train_predictions)
29
30 # Calculate accuracy train set
31 train_accuracy <- sum(train_predictions == train_data$class) / nrow(train_data)
32
33
34 # test_predictions to char vector
35 test_predictions <- as.character(test_predictions)
36
37 # Calculate accuracy test set
38 test_accuracy <- sum(test_predictions == test_data$class) / nrow(test_data)
39
40 # Create a confusion matrix for the testing set
41 conf_matrix <- table(Predicted = test_predictions, Actual = test_data$class)
42
43 # Calculate the number of false positives
44 false_positives <- conf_matrix["e", "p"]
45
46 # Print the results
47 cat("Accuracy on the training set:", train_accuracy, "\n")
48 cat("Accuracy on the testing set:", test_accuracy, "\n")
49 cat("Number of false positives:", false_positives, "\n")

```

Accuracy on the training set: 0.9521595.

Accuracy on the testing set: 0.9583702.

Number of false positives: 45.

2.3 Problem 3

Load the Yacht Hydrodynamics sample dataset from the UCI Machine Learning Repository (*yacht hydrodynamics.data*) into R using a dataframe (Note: The feature labels need to be manually specified). Use the caret package to perform a 80/20 test-train split (via the *createDataPartition* function), and obtain a training fit for a linear model. (Hint: The model fit should use all available features with the residuary resistance as the target.). What are the training MSE/RMSE and R^2 results? Next, use the caret package to perform a bootstrap from the full sample dataset with $N=1000$ samples for fitting a linear model (via the *trainControl* method), resulting in a

training MSE/RMSE and R^2 for each resample. Plot a histogram of the RMSE values, and provide a mean RMSE and R^2 for the fit. How do these values compare to the basic model? How does the performance on the test set for the original and bootstrap model compare?

```

1 # Load required libraries
2 library(caret)
3 library(boot)
4
5 # Load the dataset and specify feature labels
6 url <- "https://archive.ics.uci.edu/ml/machine-learning-databases/00243/yacht_
  hydrodynamics.data"
7 yacht_data <- read.table(url, header = FALSE, col.names = c("Longitudinal_
  position", "Prismatic_coefficient", "Length_displacement_ratio", "Beam-
  draught_ratio", "Length-beam_ratio", "Froude_number", "Residuary_resistance"
  ))
8
9 # Perform an 80/20 train-test split
10 set.seed(123) # For reproducibility
11 splitIndex <- createDataPartition(yacht_data$Residuary_resistance, p = 0.8, list
  = FALSE)
12 train_data <- yacht_data[splitIndex, ]
13 test_data <- yacht_data[-splitIndex, ]
14
15 # Fit a linear model using all features
16 lm_model <- lm(Residuary_resistance ~ ., data = train_data)
17
18 # Calculate MSE and R-squared for the training set
19 train_pred <- predict(lm_model, newdata = train_data)
20 train_mse <- mean((train_data$Residuary_resistance - train_pred)^2)
21 train_rmse <- sqrt(train_mse)
22 train_r_squared <- 1 - (sum((train_data$Residuary_resistance - train_pred)^2) /
  sum((train_data$Residuary_resistance - mean(train_data$Residuary_resistance)
  )^2))
23
24 # Print results
25 cat("Training MSE:", train_mse, "\n")
26 cat("Training RMSE:", train_rmse, "\n")
27 cat("Training R-squared:", train_r_squared, "\n")

```

Training MSE: 79.91971.

Training RMSE: 8.939783.

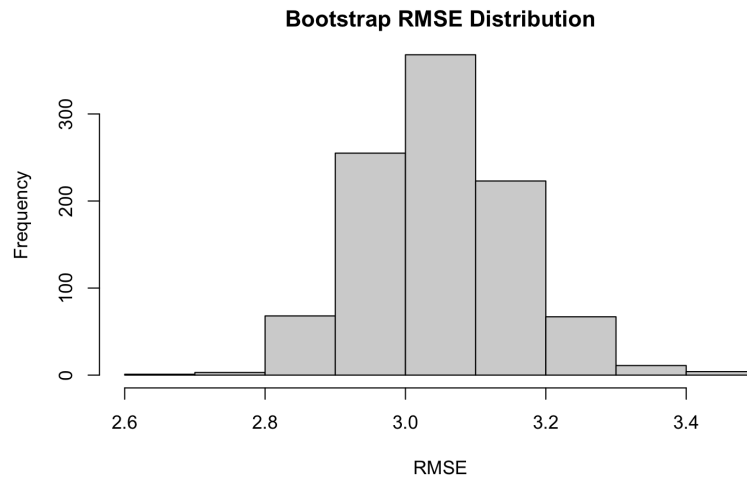
Training R-squared: 0.6606717.

```

1 # Set up bootstrap resampling
2 ctrl <- trainControl(method = "boot", number = 1000)
3
4 # Fit a linear model using bootstrap resampling
5 boot_results <- train(Residuary_resistance ~ ., data = train_data, method = "lm"
  , trControl = ctrl)
6
7 # Extract RMSE and R-squared for each resample
8 boot_rmse <- sqrt(boot_results$resample$RMSE)
9 boot_r_squared <- boot_results$resample$Rsquared
10
11 # Calculate mean RMSE and R-squared
12 mean_boot_rmse <- mean(boot_rmse)
13 mean_boot_r_squared <- mean(boot_r_squared)
14
15 # Plot histogram of RMSE values
16 hist(boot_rmse, main = "Bootstrap RMSE Distribution", xlab = "RMSE")
17
18 # Print mean RMSE and R-squared for bootstrap model
19 cat("Mean Bootstrap RMSE:", mean_boot_rmse, "\n")

```

```
20 cat("Mean Bootstrap R-squared:", mean_boot_r_squared, "\n")
```



Mean Bootstrap RMSE: 3.051137.

Mean Bootstrap R-squared: 0.6412813.

The Bootstrap RMSE is lower than the basic model, the R-squared is also a lower number in the Bootstrap model.

2.4 Problem 4

Load the German Credit Data sample dataset from the UCI Machine Learning Repository (*german.data-numeric*) into R using a dataframe (Note: The final column is the class variable coded as 1 or 2). Use the *caret* package to perform a 80/20 test-train split (via the *createDataPartition* function), and obtain a training fit for a logistic model via the *glm* package. (Hint: You may select a subset of the predictors based on exploratory analysis, or use all predictors for simplicity.). What are the training Precision/Recall and F_1 results? Next, use the *trainControl* and *train* functions to perform a $k=10$ fold cross-validation fit of the same model, and obtain cross-validated training Precision/Recall and F_1 values. How do these values compare to the original fit? How does the performance on the test set for the original and cross-validated model compare?

```
1 library(caret)
2 url <- "https://archive.ics.uci.edu/ml/machine-learning-databases/statlog/german/german.data-numeric"
3 data <- read.table(url, header = FALSE)
4 colnames(data)[ncol(data)] <- "Class"
5 set.seed(123) # For reproducibility
6 splitIndex <- createDataPartition(data$Class, p = 0.8, list = FALSE)
7 trainData <- data[splitIndex, ]
8 testData <- data[-splitIndex, ]
9
10 # Convert the "Class" variable to a factor with levels 1 and 2 in both train and test data
11 trainData$Class <- factor(trainData$Class, levels = c(1, 2))
12 testData$Class <- factor(testData$Class, levels = c(1, 2))
13
14 # Fit the logistic model using glm
15 model <- glm(Class ~ ., data = trainData, family = binomial)
16
17 # Make predictions on the training data
```

```

18 train_preds <- predict(model, newdata = trainData, type = "response")
19 train_preds <- ifelse(train_preds > 0.5, 2, 1) # Convert probabilities to class
      labels
20
21 train_preds <- factor(train_preds, levels = c(1, 2))
22
23 # Create a confusion matrix
24 conf_matrix <- confusionMatrix(data = train_preds, reference = trainData$Class)
25
26 # Extract Precision, Recall, and F1-score from the confusion matrix
27 precision <- conf_matrix$byClass["Precision"]
28 recall <- conf_matrix$byClass["Recall"]
29 f1 <- conf_matrix$byClass["F1"]
30
31 # Print the results
32 cat("Training Precision:", precision, "\n")
33 cat("Training Recall:", recall, "\n")
34 cat("Training F1 Score:", f1, "\n")

```

Training Precision: 0.8084772.

Training Recall: 0.9035088.

Training F1 Score: 0.8533554

```

1 # Convert the "Class" variable to a factor with valid R variable names
2 trainData$Class <- as.factor(trainData$Class)
3 testData$Class <- as.factor(testData$Class)
4
5 # Rename the levels to "Class_1" and "Class_2" (valid R variable names)
6 levels(trainData$Class) <- c("Class_1", "Class_2")
7 levels(testData$Class) <- c("Class_1", "Class_2")
8
9 # Define the training control with 10-fold cross-validation
10 ctrl <- trainControl(method = "cv", number = 10, classProbs = TRUE,
      summaryFunction = twoClassSummary)
11
12 # Fit the logistic model using cross-validation
13 cv_model <- train(Class ~ ., data = trainData, method = "glm", trControl = ctrl,
      family = binomial)
14
15 # Extract Precision, Recall, and F1-score from the cross-validated results
16 cv_precision <- cv_model$results$Precision
17 cv_recall <- cv_model$results$Recall
18 cv_f1 <- cv_model$results$F1

```