

Appendix –

R-codes

```

1  install.packages("MASS")
2  library(MASS)
3  library(leaps)
4  library(glm)
5  library(boot)
6  library(caret)
7  data(Pima.tr)
8
9  dim(Pima.tr)
10 data(Pima.te)
11 head(Pima.te)
12 data(Pima.te)
13 dim(Pima.te)
14 head(Pima.te)
15
16 #Cleaning the data
17 comb_data <- rbind(Pima.tr, Pima.te)
18 dim(comb_data)
19
20 any(is.na(comb_data))
21 colSums(is.na(comb_data))
22 comb_data$type <- as.factor(comb_data$type)
23
24 comb_data$type <- ifelse(comb_data$type == "Yes", 1, 0)
25 head(comb_data)
26
27 #####EDA#####
28
29 summary(comb_data)
30
31 #Histograms
32
33 par(mfrow=c(3,3))
34 hist(comb_data$npreg, main="Distribution of Number of Pregnancies", xlab="Number of Pregnancies", col="lightblue", border="black")
35 hist(comb_data$glu, main="Distribution of Glucose Levels", xlab="Glucose Levels", col="lightblue", border="black")
36 hist(comb_data$bp, main="Distribution of Blood Pressure", xlab="Blood Pressure", col="lightblue", border="black")
37 hist(comb_data$skin, main="Distribution of Skin Fold Thickness", xlab="Skin Fold Thickness", col="lightblue", border="black")
38 hist(comb_data$bmi, main="Distribution of BMI", xlab="BMI", col="lightblue", border="black")
39 hist(comb_data$ped, main="Distribution of Pedigree Function", xlab="Pedigree Function", col="lightblue", border="black") # Fixed
40 hist(comb_data$age, main="Distribution of Age", xlab="Age", col="lightblue", border="black")
41
42 #counting the the type variable
43 table(comb_data$type)
44
45
46
47 #counting the the type variable
48 table(comb_data$type)
49
50 #plotting the scatterplot matrix
51
52 # Exclude the 'type' column as it's a binary categorical variable
53 pairs(comb_data[, ], pch = 19, col = ifelse(comb_data$type == 1, "red", "blue"))
54
55 #correlation chart
56
57 cor(comb_data)
58
59 #checking the mean values for the independent variables in every column
60
61 mean_npreg_by_type <- tapply(comb_data$npreg, comb_data$type, mean)
62 mean_glu_by_type <- tapply(comb_data$glu, comb_data$type, mean)
63 mean_bp_by_type <- tapply(comb_data$bp, comb_data$type, mean)
64 mean_skin_by_type <- tapply(comb_data$skin, comb_data$type, mean)
65 mean_bmi_by_type <- tapply(comb_data$bmi, comb_data$type, mean)
66 mean_ped_by_type <- tapply(comb_data$ped, comb_data$type, mean)
67 mean_age_by_type <- tapply(comb_data$age, comb_data$type, mean)
68
69 cat("Mean of Number of Pregnancies by Type:", "\n")
70 print(mean_npreg_by_type)
71
72 cat("\nMean of Glucose Levels by Type:", "\n")
73 print(mean_glu_by_type)
74
75 cat("\nMean of Blood Pressure by Type:", "\n")
76 print(mean_bp_by_type)
77
78 cat("\nMean of Skin Fold Thickness by Type:", "\n")
79 print(mean_skin_by_type)
80
81 cat("\nMean of BMI by Type:", "\n")
82 print(mean_bmi_by_type)
83
84 cat("\nMean of Diabetes Pedigree Function by Type:", "\n")
85 print(mean_ped_by_type)
86
87 cat("\nMean of Age by Type:", "\n")
88 print(mean_age_by_type)
89

```

```

82 cat("\nMean of Age by Type:", "\n")
83 print(mean_age_by_type)
84
85 #####Building our models#####
86
87 # Constructing logistic regression model
88
89 lr_model <- glm(type ~ ., data=comb_data, family=binomial)
90 summary(lr_model)
91
92
93 #####Best Subset #####
94
95 full.model <- glm(type ~ ., data = comb_data, family = binomial)
96 subsets <- regsubsets(type ~ ., data = comb_data, method="exhaustive")
97
98 results <- summary(subsets)
99
100 bic = results$bic
101 which.min(bic)
102 aic = results$cp
103 which.min(aic)
104
105 coef(subsets, which.min(bic))
106
107 coef(subsets, which.min(aic))
108
109 aic_values = best_models$cp
110 bic_values = best_models$bic
111
112 #K-cross
113
114 k_values <- 1:length(aic_values)
115
116 optimal_k_bic <- 4
117 optimal_k_aic <- 5
118
119
120 cv.error <- function(data, indices, size) {
121   train <- data[indices, ]
122   test <- data[-indices, ]
123   fit <- glm(type ~ ., data=train, family=binomial, subset = size)
124   preds <- predict(fit, newdata=test, type="response")
125   class_preds <- ifelse(preds > 0.5, 1, 0)
126   return(mean(class_preds != test$type))
127 }

129 set.seed(123)
130 k <- 10 #setting the number of folds
131 folds <- sample(1:k, nrow(comb_data), replace=TRUE)
132 cv.errors <- matrix(NA, k, ncol(comb_data)-1)
133
134 for(j in 1:k){
135   for(i in 1:(ncol(comb_data)-1)){
136     cv.errors[j, i] <- cv.error(comb_data, which(folds != j), i)
137   }
138 }
139 mean.cv.errors <- apply(cv.errors, 2, mean)
140 optimal.number <- which.min(mean.cv.errors)
141 cat("Optimal number of predictors:", optimal.number, "\n")
142
143 best.subset <- regsubsets(type ~ ., data=comb_data, nvmax=ncol(comb_data)-1)
144 coef(best.subset, optimal.number)
145
146
147 mean.cv.errors <- apply(cv.errors, 2, mean)
148
149 # Plotting
150
151 par(mfrow=c(3,1))# 2 rows, 1 column
152
153 # Plot AIC
154 plot(k_values, aic_values, type="b", col="blue", ylab="AIC Value", xlab="k (Number of Variables)", main="AIC vs. k")
155 points(optimal_k_aic, aic_values[optimal_k_aic], col="red", pch=19, cex=2)
156 text(optimal_k_aic, aic_values[optimal_k_aic], labels="Optimal k", pos=3, col="blue")
157
158 # Plot BIC
159 plot(k_values, bic_values, type="b", col="green", ylab="BIC Value", xlab="k (Number of Variables)", main="BIC vs. k")
160 points(optimal_k_bic, bic_values[optimal_k_bic], col="red", pch=19, cex=2)
161 text(optimal_k_bic, bic_values[optimal_k_bic], labels="Optimal k", pos=4, col="green")
162 plot(1:length(mean.cv.errors), mean.cv.errors, type="b",
163      xlab="Number of Predictors", ylab="Mean CV Error",
164      main="Cross-validation Error vs. Subset Size", pch=19, col="red")
165
166 points(optimal.number, mean.cv.errors[optimal.number], col="blue", pch=4, lwd=3)
167
168 #####LDA model#####
169
170 predictors <- c("npreg", "glu", "bmi", "ped")
171
172 lda.model <- lda(type ~ npreg + glu + bmi + ped, data=comb_data)
173 print(summary(lda.model))
174 print(lda.model$means)

```

```

169 ##### LDA MODEL #####
170 predictors <- c("npreg", "glu", "bmi", "ped")
171
172 lda.model <- lda(type ~ npreg + glu + bmi + ped, data=comb_data)
173 print(summary(lda.model))
174 print(lda.model$means)
175
176 ##### QDA MODEL #####
177 qda.model <- qda(type ~ npreg + glu + bmi + ped, data = comb_data)
178
179 print(summary(qda.model))
180 print(qda.model$means)
181
182
183 #####Plotting confusion matrix#####
184
185 comb_data$type <- as.factor(comb_data$type)
186
187 #Train Models:
188 control <- trainControl(method="cv", number=10, savePredictions="final")
189
190
191
192 #Logistic Regression:
193 logistic_model <- train(type ~ npreg + glu + bmi + ped, data=comb_data, method="glm", family="binomial", trControl=control)
194
195 lda_model <- train(type ~ npreg + glu + bmi + ped, data=comb_data, method="lda", trControl=control)
196
197 qda_model <- train(type ~ npreg + glu + bmi + ped, data=comb_data, method="qda", trControl=control)
198
199 # plot confusion matrix
200
201 confusionMatrix(logistic_model$pred$pred, logistic_model$pred$obs)
202 confusionMatrix(lda_model$pred$pred, lda_model$pred$obs)
203 confusionMatrix(qda_model$pred$pred, qda_model$pred$obs)
204
205 ###Training error###
206 # For Logistic Regression:
207 logistic_preds <- predict(logistic_model, newdata = comb_data)
208 logistic_error <- mean(logistic_preds != comb_data$type)
209
210 # For LDA:
211 lda_preds <- predict(lda_model, newdata = comb_data)
212 lda_error <- mean(lda_preds != comb_data$type)
213
214 # For QDA:
215 qda_preds <- predict(qda_model, newdata = comb_data)
216 qda_error <- mean(qda_preds != comb_data$type)
217
218 # Print the training errors:
219 cat("Logistic Regression Training Error:", logistic_error, "\n")
220 cat("LDA Training Error:", lda_error, "\n")
221 cat("QDA Training Error:", qda_error, "\n")
222

```

Results –

Best subset Analysis AIC and BIC –

```

> bic = results$bic
> which.min(bic)
[1] 4
> aic = results$cp
> which.min(aic)
[1] 5
>
> coef(subsets, which.min(bic))
(Intercept)      npreg      glu      bmi      ped
-0.018360784  0.028246117  0.006275351  0.012108102  0.186859358
>
> coef(subsets, which.min(aic))
(Intercept)      npreg      glu      bmi      ped
-0.079262888  0.020206720  0.006007237  0.012009726  0.182745138
      age
0.004014933

```



```
> confusionMatrix(qda_model$pred$pred, qda_model$pred$obs)
Confusion Matrix and Statistics
```

	Reference	
Prediction	0	1
0	307	71
1	48	106

```

          Accuracy : 0.7763
          95% CI : (0.7385, 0.811)
    No Information Rate : 0.6673
    P-Value [Acc > NIR] : 2.283e-08
```

```

          Kappa : 0.4793
```

Training errors -

```

> # Print the training errors:
> cat("Logistic Regression Training Error:", logistic_error, "\n")
Logistic Regression Training Error: 0.2011278
> cat("LDA Training Error:", lda_error, "\n")
LDA Training Error: 0.2105263
> cat("QDA Training Error:", qda_error, "\n")
QDA Training Error: 0.2218045
```