

CSCI946 Assignment

Yao Xiao
SID 2019180015

October 26, 2020

1 Task 1

```
1 df <- read.csv("income.csv")
2
3 head(df)
4
5 # The linear model formula can be written:
6 #      Income = 7.26299 + 0.99520 * Age + 1.75788 * Education -
7 #              0.93433 * Gender
8 modelA <- lm(Income ~ Age + Education + Gender, data=df)
9 summary(modelA)
10
11 # Call:
12 # lm(formula = Income ~ Age + Education + Gender, data = df)
13
14 # Residuals:
15 #      Min       1Q   Median       3Q      Max
16 # -37.340  -8.101   0.139   7.885  37.271
17
18 # Coefficients:
19 #              Estimate Std. Error t value Pr(>|t|)
20 # (Intercept)  7.26299     1.95575   3.714 0.000212 ***
21 # Age         0.99520     0.02057  48.373 < 2e-16 ***
22 # Education   1.75788     0.11581  15.179 < 2e-16 ***
23 # Gender     -0.93433     0.62388  -1.498 0.134443
24 # ---
25 # Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
26
27 # Residual standard error: 12.07 on 1496 degrees of freedom
28 # Multiple R-squared:  0.6364, Adjusted R-squared:  0.6357
29 # F-statistic: 873 on 3 and 1496 DF, p-value: < 2.2e-16
30 # modelB <- lm(Income ~ Age + Education, data=df)
31
32
33 # The linear model formula can be written:
34 #      Income = 6.75822 + 0.99603 * Age + 1.75860 * Education
```

```

35 modelB <- lm(Income ~ Age + Education, data=df)
36 summary(modelB)
37 # Call:
38 # lm(formula = Income ~ Age + Education, data = df)
39
40 # Residuals:
41 #      Min       1Q   Median       3Q      Max
42 # -36.889  -7.892   0.185   8.200  37.740
43
44 # Coefficients:
45 #              Estimate Std. Error t value Pr(>|t|)
46 # (Intercept)  6.75822     1.92728   3.507 0.000467 ***
47 # Age          0.99603     0.02057  48.412 < 2e-16 ***
48 # Education    1.75860     0.11586  15.179 < 2e-16 ***
49 # ---
50 # Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
51
52 # Residual standard error: 12.08 on 1497 degrees of freedom
53 # Multiple R-squared:  0.6359, Adjusted R-squared:  0.6354
54 # F-statistic: 1307 on 2 and 1497 DF, p-value: < 2.2e-16
55
56
57 # make prediction
58 prediction <- predict(modelB, item)
59 prediction
60 # 1
61 # 68.69884
62
63
64 # compute the confidence interval
65 ci <- predict(modelB, item, interval = "confidence")
66 ci
67 # fit      lwr      upr
68 # 1 68.69884 44.98867 92.40902
69
70
71 # compute the prediction interval
72 pi <- predict(modelB, item, interval = "prediction")
73 pi
74 # fit      lwr      upr
75 # 1 68.69884 44.98867 92.40902

```

2 Task 2

```

1 library(pROC)
2 df <- read.csv("churn.csv")
3
4 head(df)
5

```

```

6  modelA <- lm(Churned ~ Age + Married + Cust_years + Churned_
  contacts , data = df)
7  summary(modelA)
8
9
10 modelB <- lm(Churned ~ Age + Married + Churned_contacts , data = df)
11 summary(modelB)
12
13
14 # The linear model formula can be written:
15 #      Churned = 0.8226510 - 0.0163168 * Age + 0.0412362 * Churned_
  contacts
16 modelC <- lm(Churned ~ Age + Churned_contacts , data = df)
17 summary(modelC)
18
19 # Call:
20 # lm(formula = Churned ~ Age + Churned_contacts , data = df)
21
22 # Residuals:
23 #      Min       1Q   Median       3Q      Max
24 # -0.77637 -0.26017 -0.04805  0.15636  1.13144
25
26 # Coefficients:
27 #              Estimate Std. Error t value Pr(>|t|)
28 # (Intercept)    0.8226510   0.0133446   61.65  <2e-16 ***
29 # Age           -0.0163168   0.0002825  -57.77  <2e-16 ***
30 # Churned_contacts 0.0412362   0.0030280   13.62  <2e-16 ***
31 # ---
32 # Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
33
34 # Residual standard error: 0.3441 on 7997 degrees of freedom
35 # Multiple R-squared:  0.3054, Adjusted R-squared:  0.3052
36 # F-statistic: 1758 on 2 and 7997 DF, p-value: < 2.2e-16
37
38 pre <- predict(modelC, type='response')
39
40 # draw ROC
41 modelCroc <- roc(df$Churned, pre)
42
43 plot(modelCroc, print.auc=TRUE, auc.polygon=TRUE, grid=c(0.1, 0.2) ,
44      grid.col=c("blue", "red"), max.auc.polygon=TRUE,
45      auc.polygon.col="skyblue", print.thres=TRUE)
46
47 # put the predicted probability prob and the actual result y in a
  data frame
48 data <- data.frame(prob=pre, obs=df$Churned)
49
50 # sort by predicted probability from low to high
51 data <- data[order(data$prob),]
52 n <- nrow(data)
53 tpr <- fpr <- rep(0,n)

```

```

54
55 # calculate TPR and FPR according to different thresholds
56 for (i in 1:n) {
57   threshold <- data$prob[i]
58   tp <- sum(data$prob > threshold & data$obs == 1)
59   fp <- sum(data$prob > threshold & data$obs == 0)
60   tn <- sum(data$prob < threshold & data$obs == 0)
61   fn <- sum(data$prob < threshold & data$obs == 1)
62   tpr[i] <- tp/(tp+fn)
63   fpr[i] <- fp/(tn+fp)
64 }
65 plot(fpr , tpr , type='l')
66 abline(a=0,b=1)

```

Figure 1: ROC curve

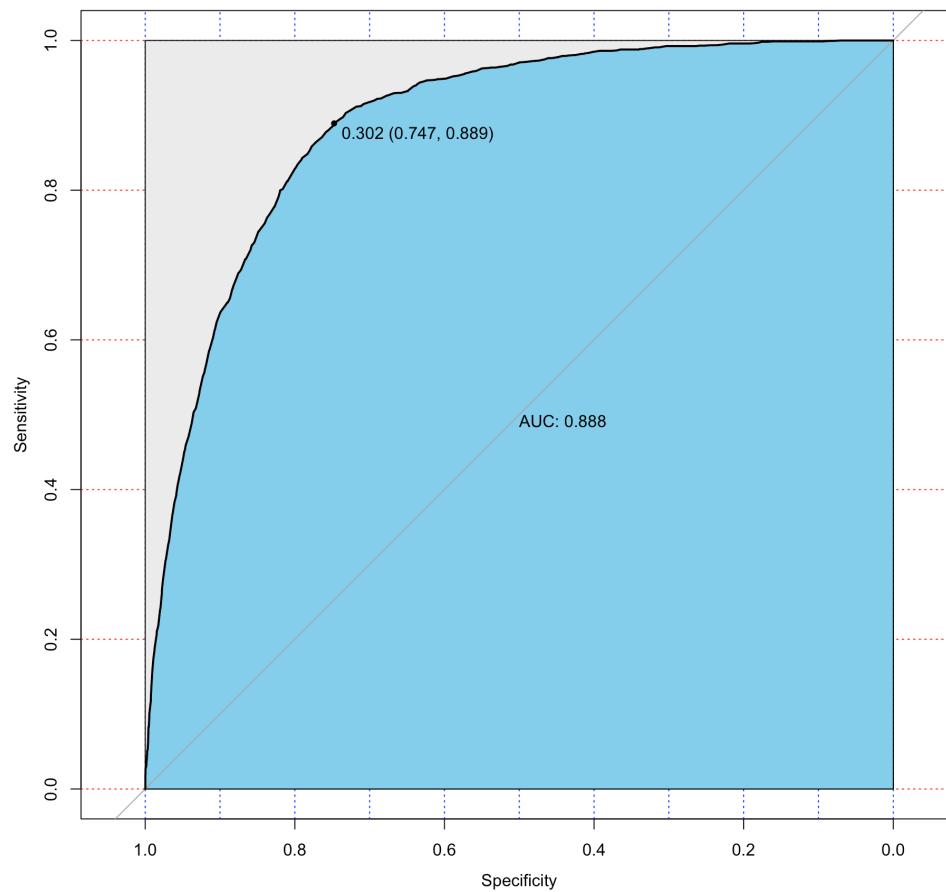


Figure 2: FPR and TPR

