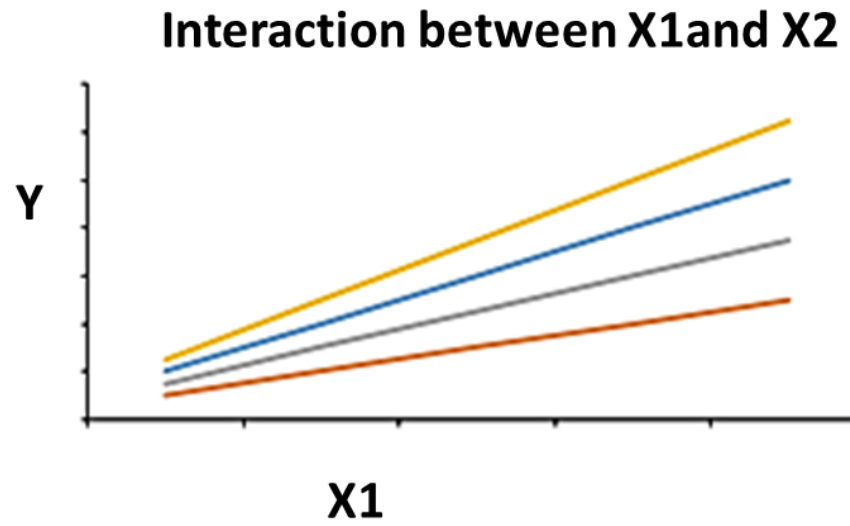# Data 603:Statistical Modelling with Data

Logistic Regression

Part III : Model Building in Multiple Logistic Regression Model and Assumptions

## Model building in Multiple Regression (An Interaction Model with both Quantitative and Qualitative variables)

An interaction occurs if the relation between one predictor, $X_1$, and the outcome (response) variable, $Y$, depends on the value of another independent variable, $X_2$. The regression coefficient for the product term represents the degree to which there is an interaction between the two variables. The effect of $X_1$ on Y is not the same for all values of $X_2$, which, in linear regression, is graphically represented by non-parallel slopes.

## Interaction between X1and X2



Non-parallel slopes represent interation terms between X1 and X2

If slopes are parallel, the effect of $X_1$ on $Y$ is the same at all levels of $X_2$, and there is no interaction.

**Variable X1 and X2 may be binary or continuous** . Interactions are similarly specified in logistic regression if the response is binary. The right hand side of the logit equation includes coefficients for the predictors, X1,X2, and X1*X2.

If the interaction coefficient $\beta_3$ is significant, we conclude that the association between $X_1$ and the probability that $Y = 1$ depends on the values of X2, X1 and X2 may be binary or continuous.

The test of the interaction may be conducted with **the Wald chi square test** or **a likelihood ratio test** comparing models with and without the interaction term.

For example, using Default data to predict the probability of default, test the interaction term for the logistic regression model

```
library(ISLR)
library(lmtest)
```

```
## Loading required package: zoo
```

```
##
## Attaching package: 'zoo'
```

```
## The following objects are masked from 'package:base':
##
##     as.Date, as.Date.numeric
```

```
mylogit <- glm(default ~ balance+income, data = Default, family = "binomial")
summary(mylogit)#Wald z test
```

```
##
## Call:
## glm(formula = default ~ balance + income, family = "binomial",
##     data = Default)
##
## Deviance Residuals:
##     Min       1Q    Median       3Q      Max
## -2.4725  -0.1444   -0.0574  -0.0211   3.7245
##
## Coefficients:
##               Estimate Std. Error z value Pr(>|z|)
## (Intercept) -1.154e+01  4.348e-01 -26.545  < 2e-16 ***
## balance      5.647e-03  2.274e-04  24.836  < 2e-16 ***
## income       2.081e-05  4.985e-06   4.174 2.99e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##     Null deviance: 2920.6  on 9999  degrees of freedom
## Residual deviance: 1579.0  on 9997  degrees of freedom
## AIC: 1585
##
## Number of Fisher Scoring iterations: 8
```

```
interlogit <- glm(default ~ balance+income+balance*income, data = Default, family = "binomial")
summary(interlogit)
```

```
##
## Call:
## glm(formula = default ~ balance + income + balance * income,
##     family = "binomial", data = Default)
##
## Deviance Residuals:
##     Min       1Q   Median       3Q      Max
## -2.5415  -0.1441  -0.0570  -0.0207   3.7546
##
## Coefficients:
##                  Estimate Std. Error z value Pr(>|z|)
## (Intercept)    -1.092e+01  9.489e-01 -11.504   <2e-16 ***
## balance         5.265e-03  5.648e-04   9.323   <2e-16 ***
## income          1.600e-06  2.683e-05   0.060    0.952
## balance:income  1.193e-08  1.638e-08   0.728    0.466
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##     Null deviance: 2920.6  on 9999  degrees of freedom
## Residual deviance: 1578.4  on 9996  degrees of freedom
## AIC: 1586.4
##
## Number of Fisher Scoring iterations: 8
```

```
anova(mylogit,interlogit,test="Chisq")
```

| Resid. Df | Resid. Dev | Df | Deviance | Pr(>Chi) |
| --- | --- | --- | --- | --- |
| <dbl> | <dbl> | <dbl> | <dbl> | <dbl> |
| 1 | 9997 | 1578.966 | NA | NA | NA |
| 2 | 9996 | 1578.431 | 1 | 0.5353864 | 0.4643511 |

2 rows

```
#likelihood ratio test
lrtest(mylogit,interlogit)
```

| | #Df | LogLik | Df | Chisq | Pr(>Chisq) |
|---|---|---|---|---|---|
| | <dbl> | <dbl> | <dbl> | <dbl> | <dbl> |
| 1 | 3 | -789.4831 | NA | NA | NA |
| 2 | 4 | -789.2154 | 1 | 0.5353864 | 0.4643511 |

2 rows

From the output, by using **the Wald Z test** and **Likelihood ratio test** ,we see that the p-value =0.466>005 (from the Wald Z test)>0.05 and the p-value =0.4644>0.05 (from Likelihood Ratio test). Therefore, the interaction term is not significant. We should drop this term out of the model.

For Default data, consider the Multiple Logistic model with both Qualitative and Quantitative variables with interation terms. We add a Student predictor (qualitative variable) into the logistic model and also add all interaction terms.

```
library(ISLR)
library(lmtest)
mylogit<- glm(default ~ balance+income+factor(student), data = Default, family = "binomial")

#Wald z test for testing individual predictors
summary(mylogit)
```

```
## 
## Call:
## glm(formula = default ~ balance + income + factor(student), family = "binomial",
##     data = Default)
## 
## Deviance Residuals:
##     Min       1Q   Median       3Q      Max
## -2.4691  -0.1418  -0.0557  -0.0203   3.7383
## 
## Coefficients:
##                    Estimate Std. Error z value Pr(>|z|)
## (Intercept)      -1.087e+01  4.923e-01 -22.080  < 2e-16 ***
## balance           5.737e-03  2.319e-04  24.738  < 2e-16 ***
## income            3.033e-06  8.203e-06   0.370  0.71152
## factor(student)Yes -6.468e-01  2.363e-01  -2.738  0.00619 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## 
## (Dispersion parameter for binomial family taken to be 1)
## 
##     Null deviance: 2920.6  on 9999  degrees of freedom
## Residual deviance: 1571.5  on 9996  degrees of freedom
## AIC: 1579.5
## 
## Number of Fisher Scoring iterations: 8
```

```
mylogit1<- glm(default ~ balance+factor(student), data = Default, family = "binomial")
summary(mylogit1)
```

```
## 
## Call:
## glm(formula = default ~ balance + factor(student), family = "binomial",
##     data = Default)
## 
## Deviance Residuals:
##     Min       1Q   Median       3Q      Max
## -2.4578  -0.1422  -0.0559  -0.0203   3.7435
## 
## Coefficients:
##                    Estimate Std. Error z value Pr(>|z|)
## (Intercept)      -1.075e+01  3.692e-01 -29.116  < 2e-16 ***
## balance           5.738e-03  2.318e-04  24.750  < 2e-16 ***
## factor(student)Yes -7.149e-01  1.475e-01  -4.846 1.26e-06 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## 
## (Dispersion parameter for binomial family taken to be 1)
## 
##     Null deviance: 2920.6  on 9999  degrees of freedom
## Residual deviance: 1571.7  on 9997  degrees of freedom
## AIC: 1577.7
## 
## Number of Fisher Scoring iterations: 8
```

```
interlogit <- glm(default ~ balance+factor(student)+balance*factor(student), data = Default, family = "binomial")
summary(interlogit)
```

```
## 
## Call:
## glm(formula = default ~ balance + factor(student) + balance *
##     factor(student), family = "binomial", data = Default)
## 
## Deviance Residuals:
##     Min       1Q   Median       3Q      Max
## -2.4839  -0.1415  -0.0553  -0.0202   3.7628
## 
## Coefficients:
##                            Estimate Std. Error z value Pr(>|z|)
## (Intercept)              -1.087e+01  4.640e-01 -23.438   <2e-16 ***
## balance                   5.819e-03  2.937e-04  19.812   <2e-16 ***
## factor(student)Yes       -3.512e-01  8.037e-01  -0.437    0.662
## balance:factor(student)Yes -2.196e-04  4.781e-04  -0.459    0.646
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## 
## (Dispersion parameter for binomial family taken to be 1)
## 
##     Null deviance: 2920.6  on 9999  degrees of freedom
## Residual deviance: 1571.5  on 9996  degrees of freedom
## AIC: 1579.5
## 
## Number of Fisher Scoring iterations: 8
```

```
#Likelihood Ratio Test for testing the interation term
anova(mylogit1,interlogit,test='Chisq')
```

| | Resid. Df | Resid. Dev | Df | Deviance | Pr(>Chi) |
| | <dbl> | <dbl> | <dbl> | <dbl> | <dbl> |
| 1 | 9997 | 1571.682 | NA | NA | NA |
| 2 | 9996 | 1571.472 | 1 | 0.2094455 | 0.6472024 |

2 rows

```
lrtest(mylogit1,interlogit)
```

| | #Df <dbl> | LogLik <dbl> | Df <dbl> | Chisq <dbl> | Pr(>Chisq) <dbl> |
|---|---|---|---|---|---|
| 1 | 3 | -785.8408 | NA | NA | NA |
| 2 | 4 | -785.7361 | 1 | 0.2094455 | 0.6472024 |

2 rows

$$H_0 : \beta_3 = 0$$
reduced model is true (with no interation term)
$$H_1 : \beta_3 \neq 0$$
larger model is true(with interation term)

The likelihood ratio statistic is

$$\triangle G^2 = -2logL \text{ from the reduced model} - (-2logL \text{ from larger model})$$
$$= -2(-785.84) - (-2(-785.74)) = 0.2094$$
The p-value is $= 0.6472 > \alpha = 0.05$

Therefore, we reject the null hypothesis which means that the interaction term (student*balance) is insignificant to be in the model.

## Inclass Practice Problem

**Example:** The German Credit Data contains data on 6 variables and the classification whether an applicant is considered a Good or a Bad credit risk for 1000 loan applicants. A predictive model developed on this data is expected to provide a bank manager guidance for making a decision

whether to approve a loan to a prospective applicant based on his/her profiles.The independent variables are listed below

Creditability= (1 if good credit, 0 if bad credit)

Balance=Account Balance (Categorical variable with 4 levels)

$$Balance = \begin{cases} 1 \text{ if balance is more than } 5000 \\ 2 \text{ if balance is } 3001-5000 \\ 3 \text{ if balance is } 1001-3000 \\ 4 \text{ if balance is less than } 1000 \end{cases}$$

Duration= Duration of credit in months (months)

Employment=Length of current employment (years)

Amount=Credit amount (dollars)

Age=Age (year)

Build the logistic regression model for predicting the probability of hiring. Check whether interation terms should be added into the model or not.

```
library("readxl")
library(lmtest)# for lrtest() function
creditdata <- read_excel("c:/Users/thuntida.ngamkham/OneDrive - University of Calgary/dataset603/creditbi
lity.xlsx")

mylogit<-glm(Creditability~employment+Duration+Amount+Age+factor(Balance),data=creditdata,family="binomia
l")
summary(mylogit)
```

```
## 
## Call:
## glm(formula = Creditability ~ employment + Duration + Amount +
##     Age + factor(Balance), family = "binomial", data = creditdata)
## 
## Deviance Residuals:
##     Min       1Q   Median       3Q      Max
## -2.4298  -0.9897   0.4731   0.8329   1.7524
## 
## Coefficients:
##                   Estimate Std. Error z value Pr(>|z|)
## (Intercept)     -1.062e-01  3.320e-01  -0.320  0.74910
## employment       1.653e-01  6.553e-02   2.523  0.01165 *
## Duration        -3.459e-02  7.835e-03  -4.414 1.01e-05 ***
## Amount          -2.820e-05  3.267e-05  -0.863  0.38807
## Age              1.173e-02  7.094e-03   1.653  0.09829 .
## factor(Balance)2 5.441e-01  1.820e-01   2.990  0.00279 **
## factor(Balance)3 1.073e+00  3.329e-01   3.222  0.00127 **
## factor(Balance)4 1.992e+00  2.035e-01   9.787  < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## 
## (Dispersion parameter for binomial family taken to be 1)
## 
##     Null deviance: 1221.7  on 999  degrees of freedom
## Residual deviance: 1039.1  on 992  degrees of freedom
## AIC: 1055.1
## 
## Number of Fisher Scoring iterations: 4
```

```
mylogit1<-glm(Creditability~employment+Duration+factor(Balance),data=creditdata,family="binomial")
lrtest(mylogit1,mylogit)
```

| | #Df | LogLik | Df | Chisq | Pr(>Chisq) |
| | <dbl> | <dbl> | <dbl> | <dbl> | <dbl> |
|---|---|---|---|---|---|
| 1 | 6 | -521.2059 | NA | NA | NA |
| 2 | 8 | -519.5724 | 2 | 3.266944 | 0.1952504 |

2 rows

```
library("readxl")
library(lmtest)# for lrtest() function
creditdata <- read_excel("c:/Users/thuntida.ngamkham/OneDrive - University of Calgary/dataset603/creditbi
lity.xlsx")

mylogit<-glm(Creditability~employment+Duration+factor(Balance),data=creditdata,family="binomial")
summary(mylogit)
```

```
##
## Call:
## glm(formula = Creditability ~ employment + Duration + factor(Balance),
##     family = "binomial", data = creditdata)
##
## Deviance Residuals:
##     Min       1Q   Median       3Q      Max
## -2.3519  -1.0079   0.4854   0.8433   1.7237
##
## Coefficients:
##                   Estimate Std. Error z value Pr(>|z|)
## (Intercept)       0.216968   0.264923   0.819 0.412794
## employment        0.193602   0.063280   3.059 0.002217 **
## Duration         -0.039002   0.006197  -6.294 3.1e-10 ***
## factor(Balance)2  0.521296   0.180830   2.883 0.003942 **
## factor(Balance)3  1.098172   0.331798   3.310 0.000934 ***
## factor(Balance)4  1.983720   0.202793   9.782  < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##     Null deviance: 1221.7  on 999  degrees of freedom
## Residual deviance: 1042.4  on 994  degrees of freedom
## AIC: 1054.4
##
## Number of Fisher Scoring iterations: 4
```

```
interlogit1<-glm(Creditability~employment+Duration+factor(Balance)+employment*Duration+employment*factor
(Balance)+Duration*factor(Balance),data=creditdata,family="binomial")
summary(interlogit1)
```

```
##
## Call:
## glm(formula = Creditability ~ employment + Duration + factor(Balance) +
##     employment * Duration + employment * factor(Balance) + Duration *
##     factor(Balance), family = "binomial", data = creditdata)
##
## Deviance Residuals:
##     Min       1Q    Median       3Q      Max
## -2.5735  -0.9871   0.4659   0.8470   1.8135
##
## Coefficients:
##                             Estimate Std. Error z value Pr(>|z|)
## (Intercept)                0.7896253  0.5736386   1.377   0.1687
## employment                 0.1089534  0.1526150   0.714   0.4753
## Duration                  -0.0500824  0.0228850  -2.188   0.0286 *
## factor(Balance)2           0.1068732  0.5885954   0.182   0.8559
## factor(Balance)3           0.0487591  1.0796399   0.045   0.9640
## factor(Balance)4           0.0796115  0.6597457   0.121   0.9040
## employment:Duration       -0.0008868  0.0054809  -0.162   0.8715
## employment:factor(Balance)2  0.0129009  0.1474950   0.087   0.9303
## employment:factor(Balance)3  0.0966643  0.2935492   0.329   0.7419
## employment:factor(Balance)4  0.4558524  0.1771205   2.574   0.0101 *
## Duration:factor(Balance)2   0.0167505  0.0154671   1.083   0.2788
## Duration:factor(Balance)3   0.0379081  0.0344753   1.100   0.2715
## Duration:factor(Balance)4   0.0196034  0.0177147   1.107   0.2685
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##     Null deviance: 1221.7  on 999  degrees of freedom
## Residual deviance: 1031.2  on 987  degrees of freedom
## AIC: 1057.2
##
## Number of Fisher Scoring iterations: 5
```

```
#------------
interlogit2<-glm(Creditability~employment+Duration+factor(Balance)+employment*factor(Balance),data=credit
data,family="binomial")
summary(interlogit2)
```

```
##
## Call:
## glm(formula = Creditability ~ employment + Duration + factor(Balance) +
##     employment * factor(Balance), family = "binomial", data = creditdata)
##
## Deviance Residuals:
##     Min       1Q   Median       3Q      Max
## -2.6021  -1.0022   0.4509   0.8475   1.6650
##
## Coefficients:
##                              Estimate Std. Error z value Pr(>|z|)
## (Intercept)                  0.615270   0.371962   1.654  0.09810 .
## employment                   0.077850   0.100932   0.771  0.44052
## Duration                    -0.039772   0.006272  -6.342 2.27e-10 ***
## factor(Balance)2             0.426465   0.504263   0.846  0.39771
## factor(Balance)3             0.559255   1.000694   0.559  0.57625
## factor(Balance)4             0.423117   0.586048   0.722  0.47030
## employment:factor(Balance)2  0.024819   0.145862   0.170  0.86489
## employment:factor(Balance)3  0.163106   0.288961   0.564  0.57244
## employment:factor(Balance)4  0.480139   0.174436   2.753  0.00591 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##     Null deviance: 1221.7  on 999  degrees of freedom
## Residual deviance: 1033.5  on 991  degrees of freedom
## AIC: 1051.5
##
## Number of Fisher Scoring iterations: 5
```

```
#Likelihood Ratio Test
lrtest(interlogit2,interlogit1)
```

| | #Df<br><dbl> | LogLik<br><dbl> | Df<br><dbl> | Chisq<br><dbl> | Pr(>Chisq)<br><dbl> |
|---|---|---|---|---|---|
| 1 | 9 | -516.7544 | NA | NA | NA |
| 2 | 13 | -515.5883 | 4 | 2.332262 | 0.6749015 |

2 rows

```
anova(interlogit2,interlogit1,test='Chisq')
```

| Resid. Df | Resid. Dev | Df | Deviance | Pr(>Chi) |
| --- | --- | --- | --- | --- |
| <dbl> | <dbl> | <dbl> | <dbl> | <dbl> |
| 1 | 991 | 1033.509 | NA | NA | NA |
| 2 | 987 | 1031.177 | 4 | 2.332262 | 0.6749015 |

2 rows

```
#Likelihood Ratio Test
lrtest(mylogit,interlogit2)
```

| #Df | LogLik | Df | Chisq | Pr(>Chisq) |
| --- | --- | --- | --- | --- |
| <dbl> | <dbl> | <dbl> | <dbl> | <dbl> |
| 1 | 6 | -521.2059 | NA | NA | NA |
| 2 | 9 | -516.7544 | 3 | 8.902915 | 0.03060992 |

2 rows

```
anova(mylogit,interlogit2,test='Chisq')
```

| Resid. Df | Resid. Dev | Df | Deviance | Pr(>Chi) |
| --- | --- | --- | --- | --- |
| <dbl> | <dbl> | <dbl> | <dbl> | <dbl> |
| 1 | 994 | 1042.412 | NA | NA | NA |
| 2 | 991 | 1033.509 | 3 | 8.902915 | 0.03060992 |

2 rows

$$\hat{y} = \frac{e^{0.61527+0.07785X_1-0.039772X_2+0.426465X_{3i}+0.559255X_{4i}+0.423117X_{5i}+0.024819X_1*X_{3i}+0.163106X_1*X_{4i}+0.480139X_1}}{1+e^{0.61527+0.07785X_1-0.039772X_2+0.426465X_{3i}+0.559255X_{4i}+0.423117X_{5i}+0.024819X_1*X_{3i}+0.163106X_1*X_{4i}+0.48013}}$$

*where,*

$$Balance = \begin{cases} 1 \text{ if balance is more than } 5000 \\ 2 \text{ if balance is } 3001-5000 \\ 3 \text{ if balance is } 1001-3000 \\ 4 \text{ if balance is less than } 1000 \end{cases}$$

## Inclass Practice Problem

**Experience in hiring.** Suppose you are investigating the hiring practices of a particular firm. Build the logistic regression model for predicitng the probability of hiring. Check whether interation terms should be added into the model or not. The data are provided in **DISCRIM.csv file**

```
library(ROCR)# for ROC
```

```
## Loading required package: gplots
```

```
##
## Attaching package: 'gplots'
```

```
## The following object is masked from 'package:stats':
##
##      lowess
```

```
library(pROC)
```

```
## Type 'citation("pROC")' for a citation.
```

```
##
## Attaching package: 'pROC'
```
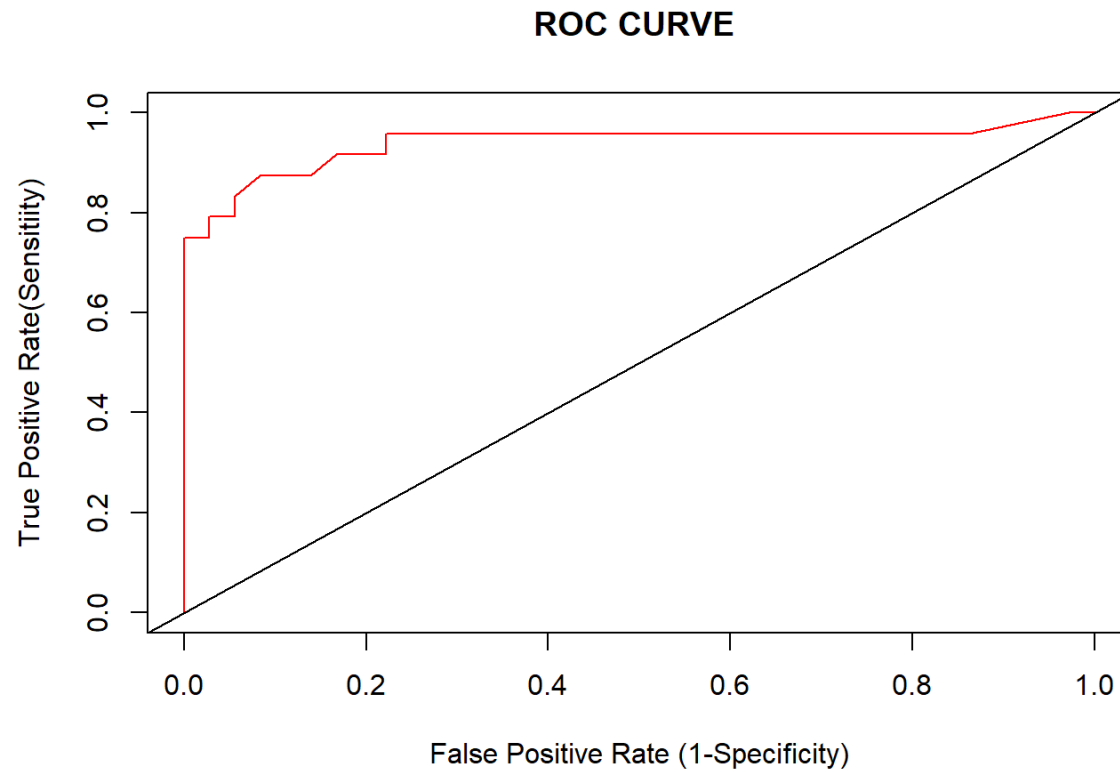
```
## The following objects are masked from 'package:stats':
##
##     cov, smooth, var
```

```
library(lmtest)# for lrtest() function
discrim=read.csv("c:/Users/thuntida.ngamkham/OneDrive - University of Calgary/dataset603/DISCRIM.csv", he
ader = TRUE)


mylogit1 <- glm(HIRE ~ factor(GENDER)+EXP+EDUC, data = discrim, family = "binomial")
summary(mylogit1)
```

```
##
## Call:
## glm(formula = HIRE ~ factor(GENDER) + EXP + EDUC, family = "binomial",
##     data = discrim)
##
## Deviance Residuals:
##     Min       1Q   Median       3Q      Max
## -1.5224  -0.4214  -0.1321   0.2853   3.2570
##
## Coefficients:
##                 Estimate Std. Error z value Pr(>|z|)
## (Intercept)      -8.2627     2.5427  -3.250 0.001156 **
## factor(GENDER)1   2.1482     0.9319   2.305 0.021160 *
## EXP               0.7602     0.2100   3.620 0.000295 ***
## EDUC              0.5509     0.2974   1.852 0.063984 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##     Null deviance: 80.761  on 59  degrees of freedom
## Residual deviance: 36.207  on 56  degrees of freedom
## AIC: 44.207
##
## Number of Fisher Scoring iterations: 6
```

```
# ROC&AUC for HIRE ~ factor(GENDER)+EXP
#--------ROC Curve-----------
prob=predict(mylogit1,type=c("response"))
pred<-prediction(prob,discrim$HIRE)
perf<-performance(pred,measure = "tpr",x.measure="fpr")
plot(perf,col=2,main="ROC CURVE ", xlab="False Positive Rate (1-Specificity)",ylab="True Positive Rate(Se
nsitiity)")
abline(0,1)
```
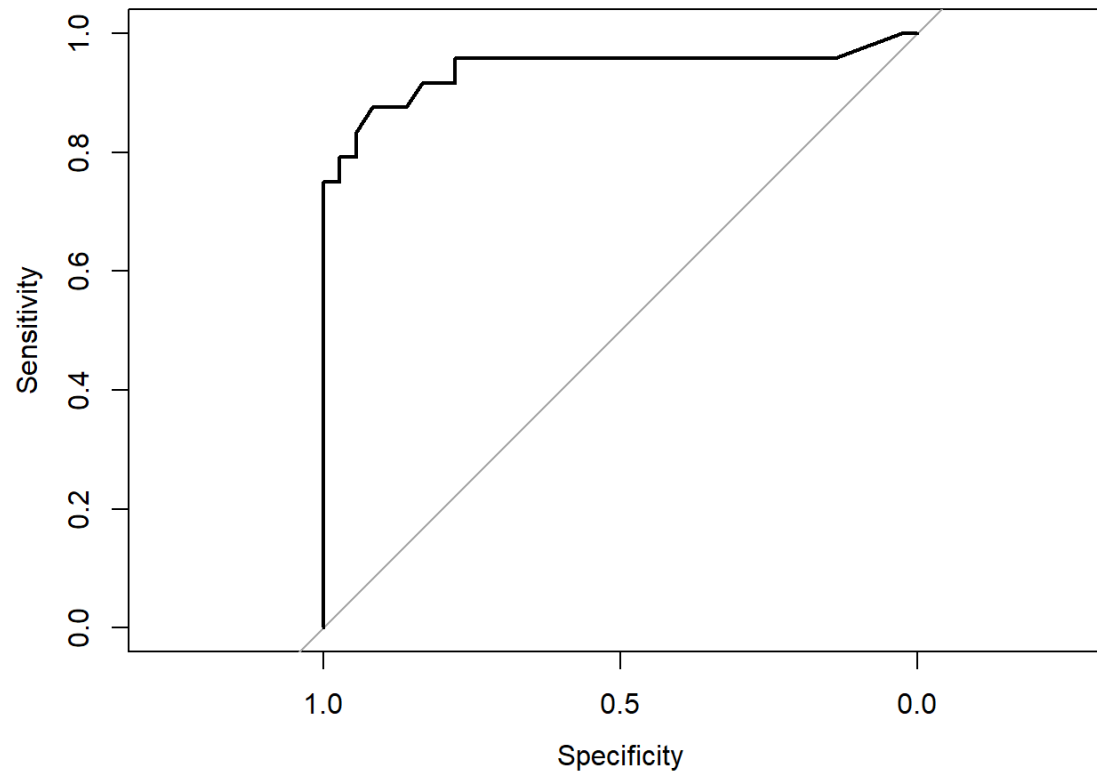
**ROC CURVE**



```
#---------AUC------------
roc<-roc(discrim$HIRE,prob)
```

```
## Setting levels: control = 0, case = 1
```

```
## Setting direction: controls < cases
```

```
plot(roc)
```
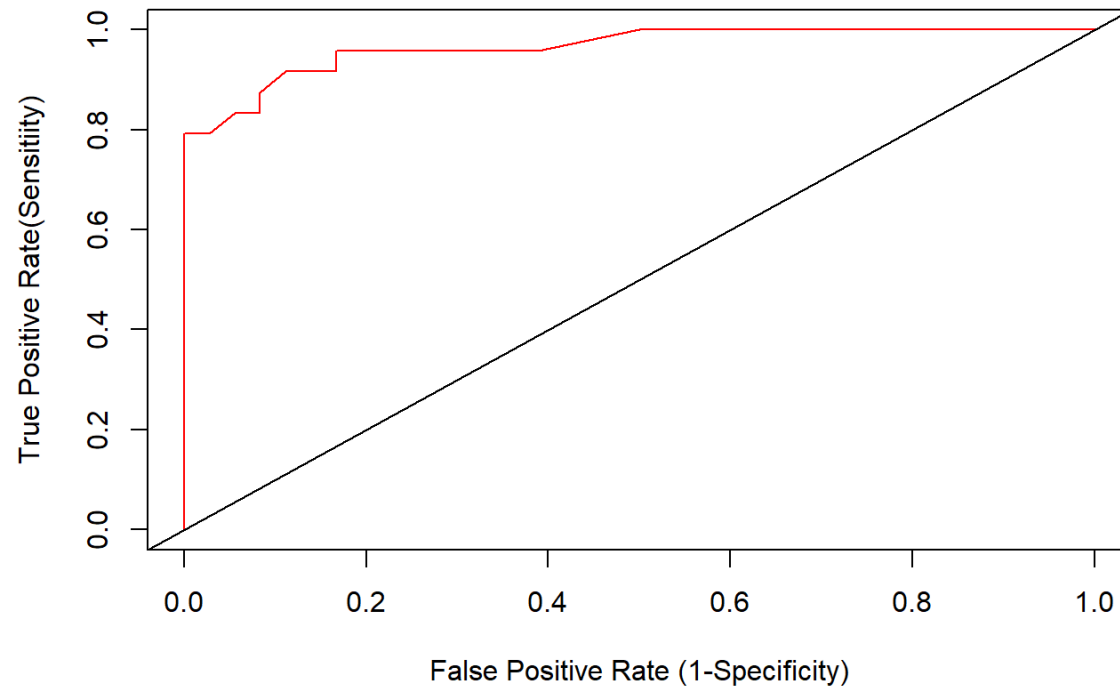


```
auc(roc)
```

```
## Area under the curve: 0.9398
```

```
fullinterlogit<-glm(HIRE ~ (factor(GENDER)+EXP+EDUC)^2, data = discrim, family = "binomial")
summary(fullinterlogit)
```

```
##
## Call:
## glm(formula = HIRE ~ (factor(GENDER) + EXP + EDUC)^2, family = "binomial",
##     data = discrim)
##
## Deviance Residuals:
##      Min        1Q    Median        3Q       Max
## -1.61085  -0.45627  -0.05627   0.00492   2.15132
##
## Coefficients:
##                      Estimate  Std. Error  z value  Pr(>|z|)
## (Intercept)            5.0141      6.5730    0.763     0.446
## factor(GENDER)1      -17.3153     12.3618   -1.401     0.161
## EXP                   -1.6162      1.3057   -1.238     0.216
## EDUC                  -1.8378      1.5193   -1.210     0.226
## factor(GENDER)1:EXP    1.6720      1.5265    1.095     0.273
## factor(GENDER)1:EDUC   2.3813      1.4568    1.635     0.102
## EXP:EDUC               0.4358      0.3000    1.452     0.146
##
## (Dispersion parameter for binomial family taken to be 1)
##
##     Null deviance: 80.761  on 59  degrees of freedom
## Residual deviance: 24.797  on 53  degrees of freedom
## AIC: 38.797
##
## Number of Fisher Scoring iterations: 9
```

```
# ROC&AUC for HIRE ~ factor(GENDER)+EXP
#--------ROC Curve-----------
prob=predict(fullinterlogit,type=c("response"))
pred<-prediction(prob,discrim$HIRE)
perf<-performance(pred,measure = "tpr",x.measure="fpr")
plot(perf,col=2,main="ROC CURVE ", xlab="False Positive Rate (1-Specificity)",ylab="True Positive Rate(Se
nsitiity)")
abline(0,1)
```

## ROC CURVE



```
#---------AUC------------
roc<-roc(discrim$HIRE,prob)
```

```
## Setting levels: control = 0, case = 1
## Setting direction: controls < cases
```
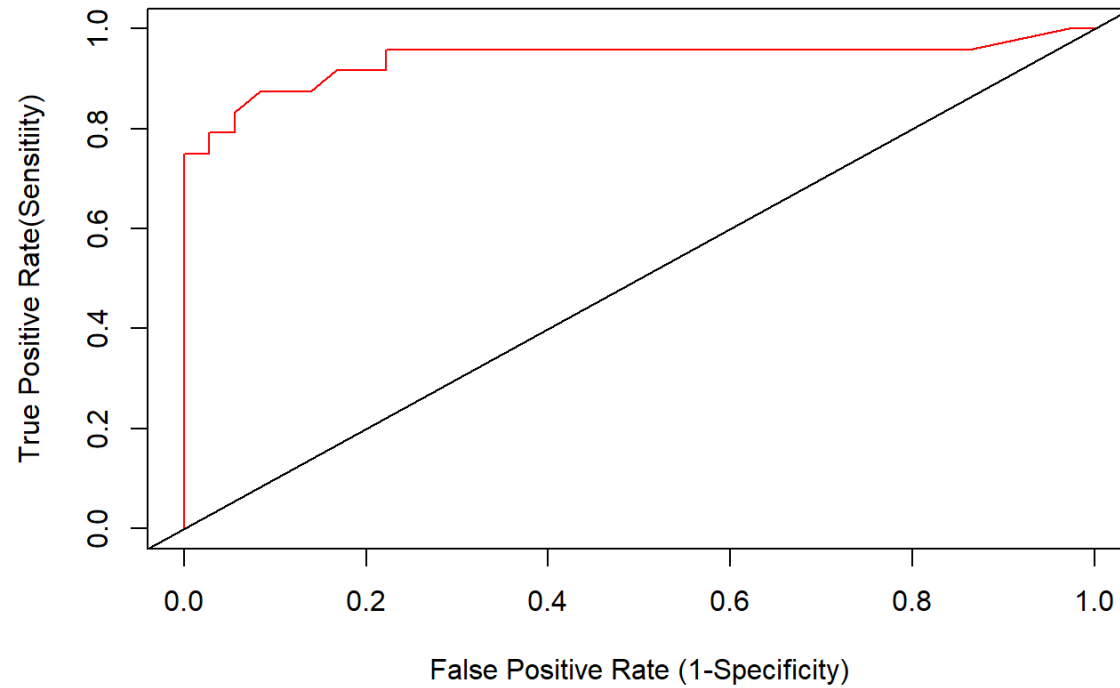
```
auc(roc)
```

```
## Area under the curve: 0.9653
```

```
interlogit1<-glm(HIRE ~ factor(GENDER)+EXP+EDUC, data = discrim, family = "binomial")
summary(interlogit1)
```

```
##
## Call:
## glm(formula = HIRE ~ factor(GENDER) + EXP + EDUC, family = "binomial",
##     data = discrim)
##
## Deviance Residuals:
##     Min       1Q   Median       3Q      Max
## -1.5224  -0.4214  -0.1321   0.2853   3.2570
##
## Coefficients:
##                 Estimate Std. Error z value Pr(>|z|)
## (Intercept)      -8.2627     2.5427  -3.250 0.001156 **
## factor(GENDER)1   2.1482     0.9319   2.305 0.021160 *
## EXP               0.7602     0.2100   3.620 0.000295 ***
## EDUC              0.5509     0.2974   1.852 0.063984 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##     Null deviance: 80.761  on 59  degrees of freedom
## Residual deviance: 36.207  on 56  degrees of freedom
## AIC: 44.207
##
## Number of Fisher Scoring iterations: 6
```

```
# ROC&AUC for HIRE ~ factor(GENDER)+EXP
#--------ROC Curve-----------
prob=predict(interlogit1,type=c("response"))
pred<-prediction(prob,discrim$HIRE)
perf<-performance(pred,measure = "tpr",x.measure="fpr")
plot(perf,col=2,main="ROC CURVE ", xlab="False Positive Rate (1-Specificity)",ylab="True Positive Rate(Se
nsitiity)")
abline(0,1)
```

## ROC CURVE



```
#---------AUC------------
roc<-roc(discrim$HIRE,prob)
```

```
## Setting levels: control = 0, case = 1
## Setting direction: controls < cases
```

```
auc(roc)
```

```
## Area under the curve: 0.9398
```

```
interlogit1<-glm(HIRE ~ factor(GENDER)+EXP+EDUC+factor(GENDER)*EXP+EXP*EDUC, data = discrim, family = "bi
nomial")
summary(interlogit1)
```

```
## 
## Call:
## glm(formula = HIRE ~ factor(GENDER) + EXP + EDUC + factor(GENDER) *
##     EXP + EXP * EDUC, family = "binomial", data = discrim)
## 
## Deviance Residuals:
##      Min        1Q    Median        3Q       Max
## -1.50750  -0.38876  -0.22780   0.06325   2.55853
## 
## Coefficients:
##                    Estimate Std. Error z value Pr(>|z|)
## (Intercept)         -2.3072     3.8982  -0.592    0.554
## factor(GENDER)1     -1.3277     2.4353  -0.545    0.586
## EXP                 -0.5710     0.8481  -0.673    0.501
## EDUC                -0.2905     0.7125  -0.408    0.683
## factor(GENDER)1:EXP  0.8743     0.6054   1.444    0.149
## EXP:EDUC             0.2015     0.1624   1.240    0.215
## 
## (Dispersion parameter for binomial family taken to be 1)
## 
##     Null deviance: 80.761  on 59  degrees of freedom
## Residual deviance: 30.248  on 54  degrees of freedom
## AIC: 42.248
## 
## Number of Fisher Scoring iterations: 7
```

```
anova(interlogit1,test='Chisq')
```

| | Df <int> | Deviance <dbl> | Resid. Df <int> | Resid. Dev <dbl> | Pr(>Chi) <dbl> |
|---|---|---|---|---|---|
| NULL | NA | NA | 59 | 80.76140 | NA |
| factor(GENDER) | 1 | 10.356555 | 58 | 70.40485 | 1.290158e-03 |
| EXP | 1 | 30.276510 | 57 | 40.12834 | 3.746356e-08 |
| EDUC | 1 | 3.921053 | 56 | 36.20728 | 4.768499e-02 |
| factor(GENDER):EXP | 1 | 4.064664 | 55 | 32.14262 | 4.378940e-02 |
| EXP:EDUC | 1 | 1.894606 | 54 | 30.24801 | 1.686833e-01 |

6 rows

```
interlogit2<-glm(HIRE ~ factor(GENDER)+EXP+EDUC+factor(GENDER)*EXP, data = discrim, family = "binomial")
summary(interlogit2)
```

```
##
## Call:
## glm(formula = HIRE ~ factor(GENDER) + EXP + EDUC + factor(GENDER) *
##      EXP, family = "binomial", data = discrim)
##
## Deviance Residuals:
##      Min        1Q     Median        3Q        Max
## -1.58945  -0.40000  -0.18130   0.05828    2.86935
##
## Coefficients:
##                    Estimate Std. Error z value Pr(>|z|)
## (Intercept)         -6.7924     2.5034  -2.713  0.00666 **
## factor(GENDER)1     -1.8914     2.4121  -0.784  0.43298
## EXP                  0.4880     0.2191   2.228  0.02591 *
## EDUC                 0.5511     0.3083   1.788  0.07384 .
## factor(GENDER)1:EXP  0.9941     0.6017   1.652  0.09852 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 80.761  on 59  degrees of freedom
## Residual deviance: 32.143  on 55  degrees of freedom
## AIC: 42.143
##
## Number of Fisher Scoring iterations: 7
```

```
#Likelihood Ratio Test
anova(mylogit1,interlogit1,test='Chisq')
```

| | Resid. Df | Resid. Dev | Df | Deviance | Pr(>Chi) |
| | <dbl> | <dbl> | <dbl> | <dbl> | <dbl> |
| 1 | 56 | 36.20728 | NA | NA | NA |
| 2 | 54 | 30.24801 | 2 | 5.95927 | 0.05081138 |

2 rows

```
lrtest(mylogit1,interlogit2)
```

| #Df | LogLik | Df | Chisq | Pr(>Chisq) |
|---|---|---|---|---|
| <dbl> | <dbl> | <dbl> | <dbl> | <dbl> |
| 1 | 4 | -18.10364 | NA | NA | NA |
| 2 | 5 | -16.07131 | 1 | 4.064664 | 0.0437894 |

2 rows

## Logistic Regression Assumptions

Logistic regression is widely used because it is a less restrictive than other techniques such as simple and multiple linear regression. Because of it, many researchers do think that LR has no an assumption at all.

**First**, logistic regression does not require _a linear relationship between the dependent and independent variables__.

**Second**, the error terms (residuals) do not need to be normally distributed_.

**Third**, homoscedasticity (constanct varaince) is not required.

**Finally**, the dependent variable in logistic regression is not measured on an interval or ratio scale.

However, there are some assumptions still apply.

First, binary logistic regression requires the dependent variable to be binary and ordinal logistic regression requires the dependent variable to be ordinal.

## Independence Assumption (Independent observations and errors)

Identical to linear regression, the assumption of independent errors states that errors should not be correlated for two observations. In logistic regression requires the observations to be independent of each other. In other words, the observations should not come from repeated measurements or matched data. This typically occurs when the data for both dependent and independent variables are observed sequentially over a period of time-called **time-series data**. Therefore, if cases are selected at random, the independent observations condition is met. If no time series data have been used, the independent errors condition is met.

## Linearity Assumption (Linear relationship between between response and predictors)

For linear regression the assumption is that the outcome variable has a linear relationship with the explanatory variables, but for logistic regression this is not possible because the outcome is binary.

## Multicolinearity

Logistic regression requires there to be little or no multicollinearity among the independent variables. This means that the independent variables should not be too highly correlated with each other. We can apply ggpairs and compute VIF from multiple linear regression to check for multicollinearity.
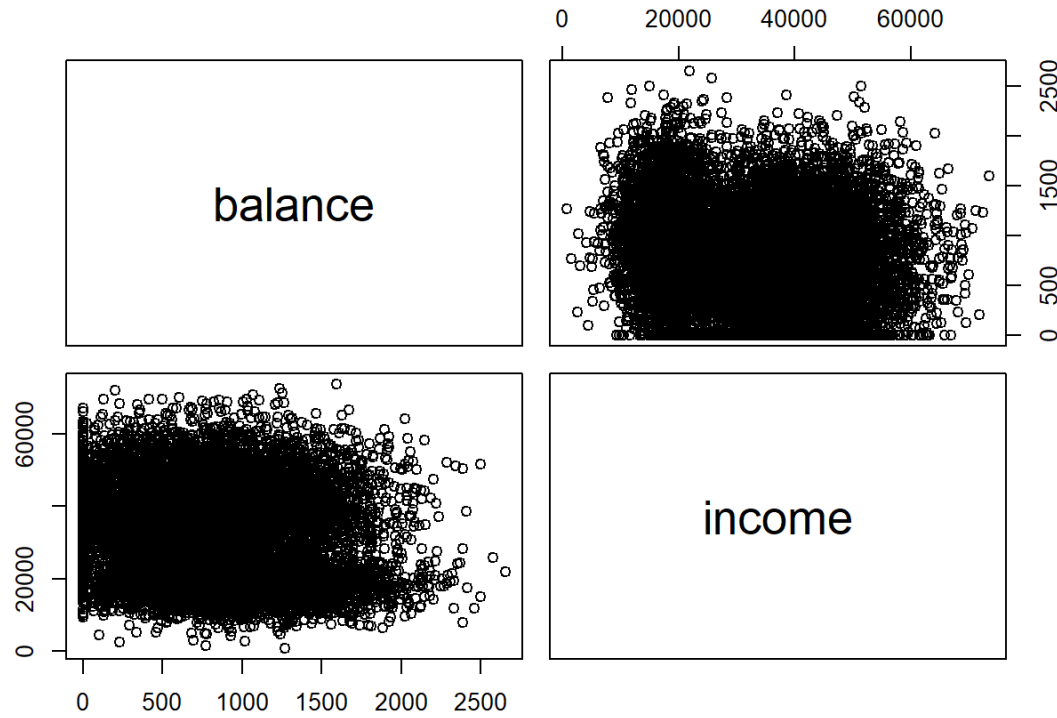
For example, using Default data to predict the probability of default, check Multicolinearity Assumption for the fitted model
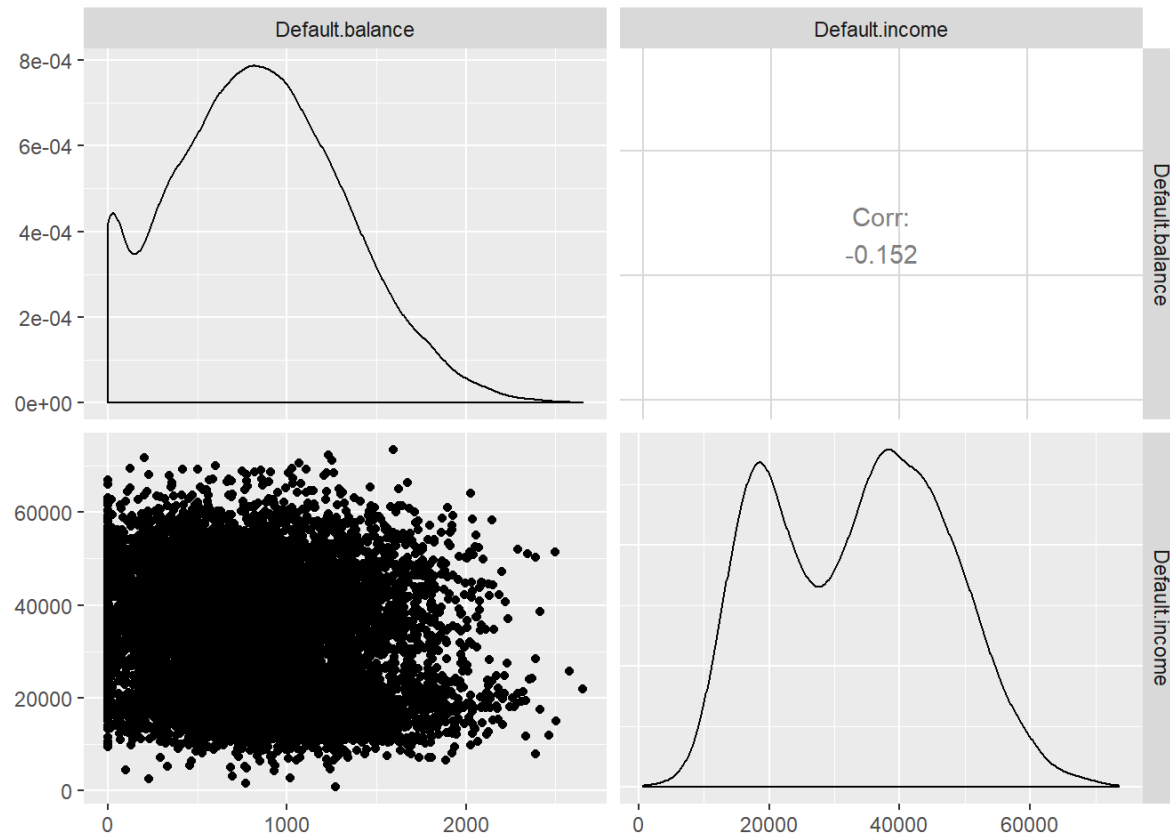
```
library(GGally)
```

```
## Loading required package: ggplot2
```

```
library(ISLR)
#Multicolinearity Assumption
pairs(~balance+income, data=Default)
```



```
defaultdata <-data.frame(Default$balance,Default$income)
ggpairs(defaultdata)
```

```r
library(mctest)
imcdiag(defaultdata,as.numeric(Default$default), method="VIF")
```

```
##
## Call:
## imcdiag(x = defaultdata, y = as.numeric(Default$default), method = "VIF")
##
##
##  VIF Multicollinearity Diagnostics
##
##                   VIF detection
## Default.balance 1.0237         0
## Default.income  1.0237         0
##
## NOTE:  VIF Method Failed to detect multicollinearity
##
##
## 0 --> COLLINEARITY is not detected by the test
##
## ==================================
```