

# Predictive Model Using Logistic Regression

Author: Data For Excellence (DFE), GPE, PSA International

Created on: 19 March 2018



## Use of this document

This file is created for the sole use of PSA Data Analytics Technical Workshop participants for demonstration and learning about the predictive modelling. All rights reserved.

## Scripting language used

This document is created using R Markdown, a scripting language available as open source from R Foundation.

## Dataset used in the model

The dataset, popularly known as "Adult" data, is publicly available in the UCI machine learning repository. The dataset is further modified for the purpose of making it useful for PSA training

## End of Introduction Section

---

---

---

---

## Start of Stage 1

### 1. The Business Understanding

- The income of a person is a function of many factors/attributes. Given enough data about these attributes, a supervised machine learning model could be developed.
- We want to predict who will earn more than 50k salary based on the 14 attributes of a person.
- The output is Yes/No or (1/0), where Yes or 1 indicate that the person will earn more than 50k. Since the output is a categorical variable, we will use Logistics Regression to predict if a person will earn 50k or not.

## End of Stage 1

---

---

---

---

## Start of Stage 2

### 2. Data Understanding

The dataset used in this project has 48,842 records and a binomial label indicating a salary of <50K or >50K USD. 76% of the records in the dataset have a class label of <50K.

### Data fields

AGE  
WORKCLASS  
FNLWGT  
EDUCATION  
EDUCATIONNUM  
MARITALSTATUS  
OCCUPATION

RELATIONSHIP  
RACE  
SEX  
CAPITALGAIN  
CAPITALLOSS  
HOURSPERWEEK  
NATIVECOUNTRY  
ABOVE50K

Loading all the required packages

```
library(dplyr)
library(InformationValue)
library(rmarkdown)
```

May need to load more libraries/packages depending on local computer/server

Loading the file into R data-frame

```
inputData <-
read.csv("https://github.com/laosze95/Training/raw/master/adult.csv")
# From Local PC use "D:/adult.csv"
# From Pramod Verma Github Page, use
"https://github.com/laosze95/Training/raw/master/adult.csv"
# From internet use "http://rstatistics.net/wp-
content/uploads/2015/09/adult.csv"

head(inputData)
```

##	AGE	WORKCLASS	FNLWGT	EDUCATION	EDUCATIONNUM	MARITALSTATUS
## 1	39	State-gov	77516	Bachelors	13	Never-married
## 2	50	Self-emp-not-inc	83311	Bachelors	13	Married-civ-spouse
## 3	38	Private	215646	HS-grad	9	Divorced
## 4	53	Private	234721	11th	7	Married-civ-spouse
## 5	28	Private	338409	Bachelors	13	Married-civ-spouse
## 6	37	Private	284582	Masters	14	Married-civ-spouse
##	OCCUPATION	RELATIONSHIP	RACE	SEX	CAPITALGAIN	CAPITALLOSS
## 1	Adm-clerical	Not-in-family	White	Male	2174	0
## 2	Exec-managerial	Husband	White	Male	0	0
## 3	Handlers-cleaners	Not-in-family	White	Male	0	0
## 4	Handlers-cleaners	Husband	Black	Male	0	0
## 5	Prof-specialty	Wife	Black	Female	0	0
## 6	Exec-managerial	Wife	White	Female	0	0
##	HOURSPERWEEK	NATIVECOUNTRY	ABOVE50K			
## 1	40	United-States	0			
## 2	13	United-States	0			
## 3	40	United-States	0			
## 4	40	United-States	0			
## 5	40	Cuba	0			
## 6	40	United-States	0			

```
inputData <- tbl_df(inputData)
```

## Looking at the structure of the data

```
dim(inputData)
```

```
## [1] 32561    15
```

```
class(inputData)
```

```
## [1] "tbl_df"      "tbl"        "data.frame"
```

```
str(inputData)
```

```
## Classes 'tbl_df', 'tbl' and 'data.frame':  32561 obs. of  15 variables:
## $ AGE          : int  39 50 38 53 28 37 49 52 31 42 ...
## $ WORKCLASS    : Factor w/ 9 levels " ?"," Federal-gov",...: 8 7 5 5 5 5 5
## $ FNLWGT       : int  77516 83311 215646 234721 338409 284582 160187
## $ EDUCATION    : Factor w/ 16 levels " 10th"," 11th",...: 10 10 12 2 10 13
## $ EDUCATIONNUM : int  13 13 9 7 13 14 5 9 14 13 ...
## $ MARITALSTATUS: Factor w/ 7 levels " Divorced"," Married-AF-spouse",...:
## $ OCCUPATION   : Factor w/ 15 levels " ?"," Adm-clerical",...: 2 5 7 7 11
## $ RELATIONSHIP : Factor w/ 6 levels " Husband"," Not-in-family",...: 2 1 2
## $ RACE         : Factor w/ 5 levels " Amer-Indian-Eskimo",...: 5 5 5 3 3 5
## $ SEX          : Factor w/ 2 levels " Female"," Male": 2 2 2 2 1 1 1 2 1
## $ CAPITALGAIN  : int  2174 0 0 0 0 0 0 0 14084 5178 ...
## $ CAPITALLOSS  : int  0 0 0 0 0 0 0 0 0 0 ...
## $ HOURSPERWEEK : int  40 13 40 40 40 40 16 45 50 40 ...
## $ NATIVECOUNTRY: Factor w/ 42 levels " ?"," Cambodia",...: 40 40 40 40 6
## $ ABOVE50K     : int  0 0 0 0 0 0 0 1 1 1 ...
```

```
summary(inputData)
```

##	AGE	WORKCLASS	FNLWGT
## Min.	:17.00	Private	:22696
## 1st Qu.:	:28.00	Self-emp-not-inc	:2541
## Median	:37.00	Local-gov	:2093
## Mean	:38.58	?	:1836
## 3rd Qu.:	:48.00	State-gov	:1298
## Max.	:90.00	Self-emp-inc	:1116
##		(Other)	:981
##	EDUCATION	EDUCATIONNUM	MARITALSTATUS

```

## HS-grad      :10501  Min.   : 1.00   Divorced      : 4443
## Some-college: 7291  1st Qu.: 9.00   Married-AF-spouse : 23
## Bachelors    : 5355  Median :10.00   Married-civ-spouse :14976
## Masters      : 1723  Mean    :10.08   Married-spouse-absent: 418
## Assoc-voc    : 1382  3rd Qu.:12.00   Never-married      :10683
## 11th         : 1175  Max.    :16.00   Separated          : 1025
## (Other)      : 5134                               Widowed            : 993
##
##              OCCUPATION                      RELATIONSHIP
## Prof-specialty :4140   Husband           :13193
## Craft-repair   :4099   Not-in-family     : 8305
## Exec-managerial:4066   Other-relative:   981
## Adm-clerical   :3770   Own-child         : 5068
## Sales          :3650   Unmarried         : 3446
## Other-service   :3295   Wife              : 1568
## (Other)        :9541
##
##              RACE                      SEX          CAPITALGAIN
## Amer-Indian-Eskimo: 311   Female:10771   Min.   : 0
## Asian-Pac-Islander:1039   Male  :21790   1st Qu.: 0
## Black           : 3124                               Median : 0
## Other           : 271                                Mean   :1078
## White           :27816                               3rd Qu.: 0
##                                                         Max.   :99999
##
## CAPITALLOSS  HOURS PER WEEK  NATIVE COUNTRY  ABOVE50K
## Min.   : 0.0   Min.   : 1.00   United-States:29170   Min.   :0.0000
## 1st Qu.: 0.0   1st Qu.:40.00   Mexico          : 643   1st Qu.:0.0000
## Median : 0.0   Median :40.00   ?              : 583   Median :0.0000
## Mean   : 87.3   Mean   :40.44   Philippines     : 198   Mean   :0.2408
## 3rd Qu.: 0.0   3rd Qu.:45.00   Germany         : 137   3rd Qu.:0.0000
## Max.   :4356.0   Max.   :99.00   Canada          : 121   Max.   :1.0000
##                                     (Other)       : 1709

```

- There are 14 attributes consisting of eight categorical and six continuous attributes. The work class describes the type of employer such as self-employed or federal and occupation describes the employment type such as farming, clerical or managerial.
- Education contains the highest level of education attained such as high school or doctorate.
- The relationship attribute has categories such as unmarried or husband and marital status has categories such as married or separated.
- The other nominal attributes are country of residence, gender and race.
- The continuous attributes are age, hours worked per week, education number (numeric representation of the education attribute), capital gain and loss, and a weight attribute which is a demographic score assigned to an individual based on information such as state of residence and type of employment.

- Some of the variables are not self-explanatory. The continuous variable `fnlwgt` represents final weight, which is the number of units in the target population that the responding unit represents.
- The variable `education_num` stands for the number of years of education in total, which is a continuous representation of the discrete variable `education`. The variable `relationship` represents the responding unit's role in the family.
- `Capital_gain` and `capital_loss` are income from investment sources other than wage/salary.
- For simplicity of this analysis, the weighting factor is discarded. Total number of years of education can represent by the highest education level completed. Role in the family can be assessed from gender and marital status. Thus, the following 3 variables are deleted `education`, `relationship`, and `fnlwgt`.

### Checking the class bias of the data

```
table(inputData$ABOVE50K)
```

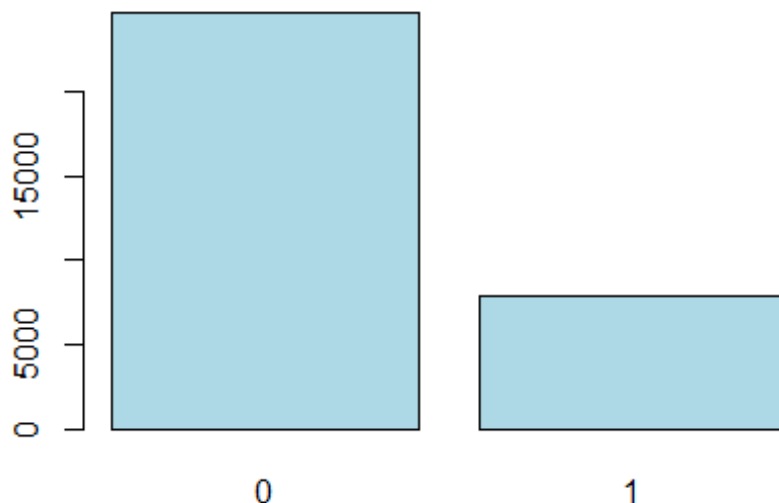
```
##
```

```
##      0      1
```

```
## 24720  7841
```

```
# histogram of age by income group
```

```
barplot(table(inputData$ABOVE50K), col = "lightblue")
```



Since there is a class bias, a condition observed when the proportion of events is much smaller than proportion of non-events. So we must sample the observations in approximately equal proportions to get better models.

## End of Stage 2

---

---

---

---

## Start of Stage 3

### 3. Data Preparation

First we want to clean up the data set to include only those variables which are important. From our data understanding, we know FNLWGT and RELATIONSHIP is not required.

```
inputData$FNLWGT <- NULL
inputData$RELATIONSHIP <- NULL
head(inputData$FNLWGT)

## Warning: Unknown or uninitialised column: 'FNLWGT'.
## NULL

head(inputData$RELATIONSHIP)

## Warning: Unknown or uninitialised column: 'RELATIONSHIP'.
## NULL
```

**Creating two sets of data from given data** \* Training set - For training the model \* Test set - For test and validation

#### Creating training data set

```
input_ones <- inputData[which(inputData$ABOVE50K == 1), ] # all 1's
input_zeros <- inputData[which(inputData$ABOVE50K == 0), ] # all 0's

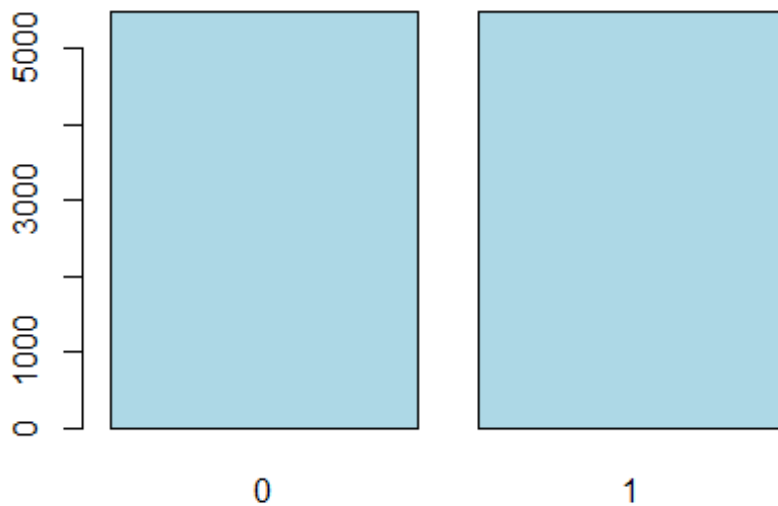
set.seed(100) # for repeatability of samples

input_ones_training_rows <- sample(1:nrow(input_ones), 0.7*nrow(input_ones))
# 1's for training
input_zeros_training_rows <- sample(1:nrow(input_zeros),
0.7*nrow(input_ones)) # 0's for training.

#Pick as many 0's as 1's
training_ones <- input_ones[input_ones_training_rows, ]
training_zeros <- input_zeros[input_zeros_training_rows, ]
```

```
# row bind the 1's and 0's
trainingData <- rbind(training_ones, training_zeros)

# Checking the bias on training data
barplot(table(trainingData$ABOVE50K), col = "lightblue")
```



```
head(trainingData)

## # A tibble: 6 x 13
##   AGE      WORKCLASS EDUCATION EDUCATIONNUM MARITALSTATUS
##   <int>      <fctr>    <fctr>      <int>      <fctr>
## 1   49 Self-emp-not-inc HS-grad         9 Divorced
## 2   44   State-gov     Masters        14 Married-civ-spouse
## 3   36   Federal-gov  Masters        14 Married-civ-spouse
## 4   51     Private   Assoc-voc       11 Married-civ-spouse
## 5   49   Local-gov   HS-grad         9 Married-civ-spouse
## 6   41     Private   HS-grad         9 Married-civ-spouse
## # ... with 8 more variables: OCCUPATION <fctr>, RACE <fctr>, SEX <fctr>,
## #   CAPITALGAIN <int>, CAPITALLOSS <int>, HOURSPERWEEK <int>,
## #   NATIVECOUNTRY <fctr>, ABOVE50K <int>
```

### Creating the test data set

```
test_ones <- input_ones[-input_ones_training_rows, ]
test_zeros <- input_zeros[-input_zeros_training_rows, ]

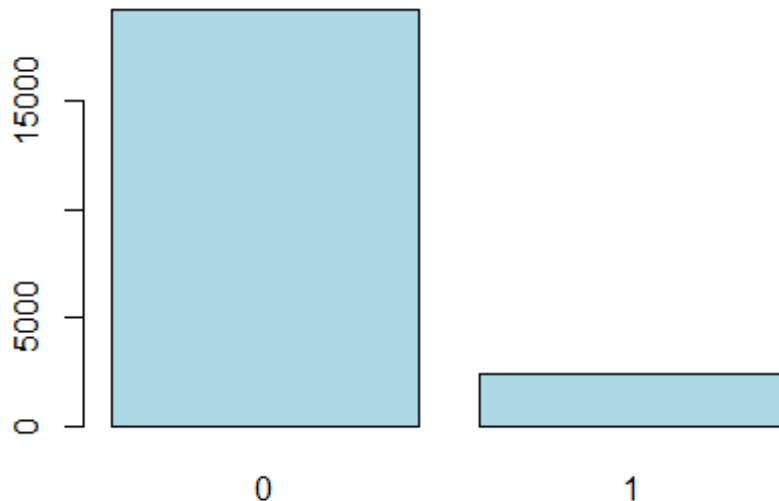
# row bind the 1's and 0's
```



```
testData <- rbind(test_ones, test_zeros)
```

*# We do not need to correct the bias on test data because model should take care of future uncertainty*

```
barplot(table(testData$ABOVE50K), col = "lightblue")
```



```
head(testData)
```

```
## # A tibble: 6 x 13
##   AGE  WORKCLASS      EDUCATION EDUCATIONNUM  MARITALSTATUS
##   <int> <fctr>      <fctr>      <int>      <fctr>
## 1   31   Private    Masters        14   Never-married
## 2   30 State-gov  Bachelors      13   Married-civ-spouse
## 3   56 Local-gov  Bachelors      13   Married-civ-spouse
## 4   31   Private  Some-college    10   Married-civ-spouse
## 5   43   Private  Some-college    10   Married-civ-spouse
## 6   42   Private   Doctorate     16   Married-civ-spouse
## # ... with 8 more variables: OCCUPATION <fctr>, RACE <fctr>, SEX <fctr>,
## #   CAPITALGAIN <int>, CAPITALLOSS <int>, HOURSPERWEEK <int>,
## #   NATIVECOUNTRY <fctr>, ABOVE50K <int>
```

## Feature Selection

- Now we want to know that out of 14 attributes, which are the most important one. There are many methods to find out the best attributes. We will use WOE (Weight of Evidence) method. The choice of feature selection is based on data types and model types.

**\*Weight of evidence (WOE)** is a measure of how much the evidence supports or undermines a hypothesis. WOE measures the relative risk of an attribute of binning level. The value depends on whether the value of the target variable is a non-event or an event.

## Compute Information Values

We will compute information values for both categorical and continuous variable. The continuous variable needs to be converted to categorical variable before we compute information value.

```
# segregate continuous and factor variables
factor_vars <- c("WORKCLASS", "EDUCATION", "MARITALSTATUS", "OCCUPATION",
"RELATIONSHIP", "RACE", "SEX", "NATIVECOUNTRY")
continuous_vars <- c("AGE", "FNLWGT", "EDUCATIONNUM", "HOURSPERWEEK",
"CAPITALGAIN", "CAPITALLOSS")

# initialization for the for IV results
iv_df <- data.frame(VARS=c(factor_vars, continuous_vars), IV=numeric(14))

# compute IV for categorical Variables

iv_df[iv_df$VARS == "WORKCLASS", "IV"] <- IV(X=inputData$WORKCLASS,
Y=inputData$ABOVE50K)[1]
iv_df[iv_df$VARS == "EDUCATION", "IV"] <- IV(X=inputData$EDUCATION,
Y=inputData$ABOVE50K)[1]
iv_df[iv_df$VARS == "MARITALSTATUS", "IV"] <- IV(X=inputData$MARITALSTATUS,
Y=inputData$ABOVE50K)[1]
iv_df[iv_df$VARS == "OCCUPATION", "IV"] <- IV(X=inputData$OCCUPATION,
Y=inputData$ABOVE50K)[1]
iv_df[iv_df$VARS == "RACE", "IV"] <- IV(X=inputData$RACE,
Y=inputData$ABOVE50K)[1]
iv_df[iv_df$VARS == "SEX", "IV"] <- IV(X=inputData$SEX,
Y=inputData$ABOVE50K)[1]
iv_df[iv_df$VARS == "NATIVECOUNTRY", "IV"] <- IV(X=inputData$NATIVECOUNTRY,
Y=inputData$ABOVE50K)[1]

# compute IV for Continuous Variables

iv_df[iv_df$VARS == "AGE", "IV"] <- IV(X=as.factor(inputData$AGE),
Y=inputData$ABOVE50K)[1]
iv_df[iv_df$VARS == "EDUCATIONNUM", "IV"] <-
IV(X=as.factor(inputData$EDUCATIONNUM), Y=inputData$ABOVE50K)[1]
iv_df[iv_df$VARS == "HOURSPERWEEK", "IV"] <-
IV(X=as.factor(inputData$HOURSPERWEEK), Y=inputData$ABOVE50K)[1]
iv_df[iv_df$VARS == "CAPITALGAIN", "IV"] <-
IV(X=as.factor(inputData$CAPITALGAIN), Y=inputData$ABOVE50K)[1]
```

```
iv_df[iv_df$VARS == "CAPITALLOSS", "IV"] <-
IV(X=as.factor(inputData$CAPITALLOSS), Y=inputData$ABOVE50K)[1]
```

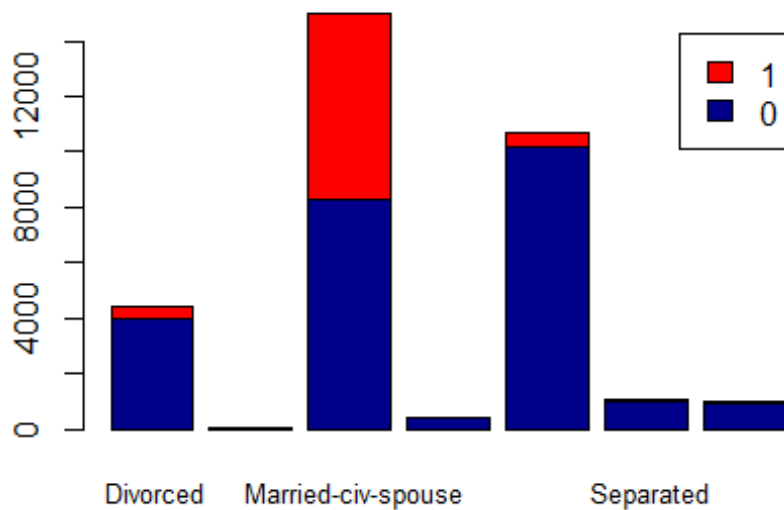
```
iv_df <- iv_df[order(-iv_df$IV), ] # sort
iv_df
```

```
##          VARS          IV
## 3  MARITALSTATUS 1.33882907
## 9          AGE 0.88214658
## 4   OCCUPATION 0.77622839
## 2   EDUCATION 0.74105372
## 11 EDUCATIONNUM 0.74105372
## 12 HOURSPERWEEK 0.49628770
## 13  CAPITALGAIN 0.31266990
## 7          SEX 0.30328938
## 14  CAPITALLOSS 0.20749663
## 1   WORKCLASS 0.16338802
## 8  NATIVECOUNTRY 0.07939344
## 6          RACE 0.06929987
## 5  RELATIONSHIP 0.00000000
## 10         FNLWGT 0.00000000
```

```
table(inputData$ABOVE50K, inputData$MARITALSTATUS )
```

```
##
##      Divorced  Married-AF-spouse  Married-civ-spouse
## 0          3980                13                8284
## 1           463                10                6692
##
##      Married-spouse-absent  Never-married  Separated  Widowed
## 0                   384            10192         959        908
## 1                   34             491         66         85
```

```
barplot(table(inputData$ABOVE50K, inputData$MARITALSTATUS
), col=c("darkblue", "red"), legend = TRUE, cex.names=0.8)
```



### End of Stage 3

---



---



---



---

### Start of Stage 4

## 4. Modelling

Building the Logistic Model using the most significant attributes which are

MARITALSTATUS  
 AGE  
 OCCUPATION  
 EDUCATION  
 EDUCATIONNUM  
 HOURSPERWEEK  
 CAPITALGAIN  
 SEX

However, we see that EDUCATION AND EDUCATIONNUM ARE HIGHLY CORELATED SO WE CAN PICK ONLY ONE

```
logitMod <- glm(ABOVE50K ~ MARITALSTATUS + AGE + OCCUPATION + EDUCATION +  
HOURSPERWEEK + CAPITALGAIN + SEX, data=trainingData,  
family=binomial(link="logit"))
```

```
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
```

```
# predicted scores
```

```
predicted <- predict(logitMod, testData, type="response")  
summary(logitMod)
```

```
##
```

```
## Call:
```

```
## glm(formula = ABOVE50K ~ MARITALSTATUS + AGE + OCCUPATION + EDUCATION +  
##   HOURSPERWEEK + CAPITALGAIN + SEX, family = binomial(link = "logit"),  
##   data = trainingData)
```

```
##
```

```
## Deviance Residuals:
```

```
##      Min       1Q   Median       3Q      Max  
## -3.6709  -0.5232  -0.0001   0.6198   3.3011
```

```
##
```

```
## Coefficients:
```

```
##              Estimate Std. Error z value Pr(>|z|)  
## (Intercept)    -6.401e+00  3.296e-01 -19.421 < 2e-16  
## MARITALSTATUS Married-AF-spouse    3.782e+00  9.181e-01  4.120 3.79e-05  
## MARITALSTATUS Married-civ-spouse    2.272e+00  9.649e-02 23.550 < 2e-16  
## MARITALSTATUS Married-spouse-absent  6.645e-02  3.133e-01  0.212 0.832043  
## MARITALSTATUS Never-married    -4.266e-01  1.160e-01 -3.676 0.000236  
## MARITALSTATUS Separated    -2.519e-01  2.177e-01 -1.157 0.247222  
## MARITALSTATUS Widowed     3.010e-02  2.084e-01  0.144 0.885141  
## AGE     3.010e-02  2.582e-03 11.656 < 2e-16  
## OCCUPATION Adm-clerical    7.243e-01  1.682e-01  4.305 1.67e-05  
## OCCUPATION Armed-Forces   -5.469e-02  1.489e+00 -0.037 0.970709  
## OCCUPATION Craft-repair    8.360e-01  1.611e-01  5.190 2.10e-07  
## OCCUPATION Exec-managerial  1.577e+00  1.615e-01  9.766 < 2e-16  
## OCCUPATION Farming-fishing  -6.030e-01  2.223e-01 -2.713 0.006668  
## OCCUPATION Handlers-cleaners  1.449e-01  2.258e-01  0.642 0.521137  
## OCCUPATION Machine-op-inspct  3.845e-01  1.837e-01  2.093 0.036393  
## OCCUPATION Other-service   -7.268e-02  1.947e-01 -0.373 0.708893  
## OCCUPATION Priv-house-serv  -3.540e+00  2.815e+00 -1.258 0.208530  
## OCCUPATION Prof-specialty    1.386e+00  1.661e-01  8.342 < 2e-16  
## OCCUPATION Protective-serv    1.130e+00  2.218e-01  5.096 3.48e-07  
## OCCUPATION Sales     9.440e-01  1.643e-01  5.745 9.22e-09  
## OCCUPATION Tech-support    1.457e+00  2.033e-01  7.167 7.69e-13  
## OCCUPATION Transport-moving  5.464e-01  1.832e-01  2.983 0.002857  
## EDUCATION 11th    -1.288e-01  2.957e-01 -0.435 0.663212  
## EDUCATION 12th    -2.787e-01  4.052e-01 -0.688 0.491579  
## EDUCATION 1st-4th  -1.037e+00  5.499e-01 -1.885 0.059371  
## EDUCATION 5th-6th  -5.688e-01  4.215e-01 -1.349 0.177210
```

## EDUCATION 7th-8th	-5.481e-01	3.288e-01	-1.667	0.095570
## EDUCATION 9th	-2.157e-01	3.715e-01	-0.580	0.561585
## EDUCATION Assoc-acdm	1.212e+00	2.606e-01	4.649	3.33e-06
## EDUCATION Assoc-voc	1.230e+00	2.508e-01	4.902	9.51e-07
## EDUCATION Bachelors	1.791e+00	2.305e-01	7.771	7.79e-15
## EDUCATION Doctorate	2.796e+00	3.396e-01	8.233	< 2e-16
## EDUCATION HS-grad	6.236e-01	2.224e-01	2.804	0.005048
## EDUCATION Masters	1.978e+00	2.487e-01	7.952	1.83e-15
## EDUCATION Preschool	-1.229e+01	1.250e+02	-0.098	0.921705
## EDUCATION Prof-school	3.214e+00	3.481e-01	9.232	< 2e-16
## EDUCATION Some-college	1.018e+00	2.263e-01	4.498	6.86e-06
## HOURSPERWEEK	3.448e-02	2.625e-03	13.133	< 2e-16
## CAPITALGAIN	3.227e-04	1.773e-05	18.202	< 2e-16
## SEX Male	1.616e-01	7.719e-02	2.094	0.036242
##				
## (Intercept)	***			
## MARITALSTATUS Married-AF-spouse	***			
## MARITALSTATUS Married-civ-spouse	***			
## MARITALSTATUS Married-spouse-absent				
## MARITALSTATUS Never-married	***			
## MARITALSTATUS Separated				
## MARITALSTATUS Widowed				
## AGE	***			
## OCCUPATION Adm-clerical	***			
## OCCUPATION Armed-Forces				
## OCCUPATION Craft-repair	***			
## OCCUPATION Exec-managerial	***			
## OCCUPATION Farming-fishing	**			
## OCCUPATION Handlers-cleaners				
## OCCUPATION Machine-op-inspct	*			
## OCCUPATION Other-service				
## OCCUPATION Priv-house-serv				
## OCCUPATION Prof-specialty	***			
## OCCUPATION Protective-serv	***			
## OCCUPATION Sales	***			
## OCCUPATION Tech-support	***			
## OCCUPATION Transport-moving	**			
## EDUCATION 11th				
## EDUCATION 12th				
## EDUCATION 1st-4th	.			
## EDUCATION 5th-6th				
## EDUCATION 7th-8th	.			
## EDUCATION 9th				
## EDUCATION Assoc-acdm	***			
## EDUCATION Assoc-voc	***			
## EDUCATION Bachelors	***			
## EDUCATION Doctorate	***			
## EDUCATION HS-grad	**			
## EDUCATION Masters	***			
## EDUCATION Preschool				

```
## EDUCATION Prof-school          ***
## EDUCATION Some-college         ***
## HOURSPERWEEK                   ***
## CAPITALGAIN                    ***
## SEX Male                        *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 15216.0  on 10975  degrees of freedom
## Residual deviance:  8541.7  on 10936  degrees of freedom
## AIC: 8621.7
##
## Number of Fisher Scoring iterations: 12
```

## End of Stage 4

---



---



---



---

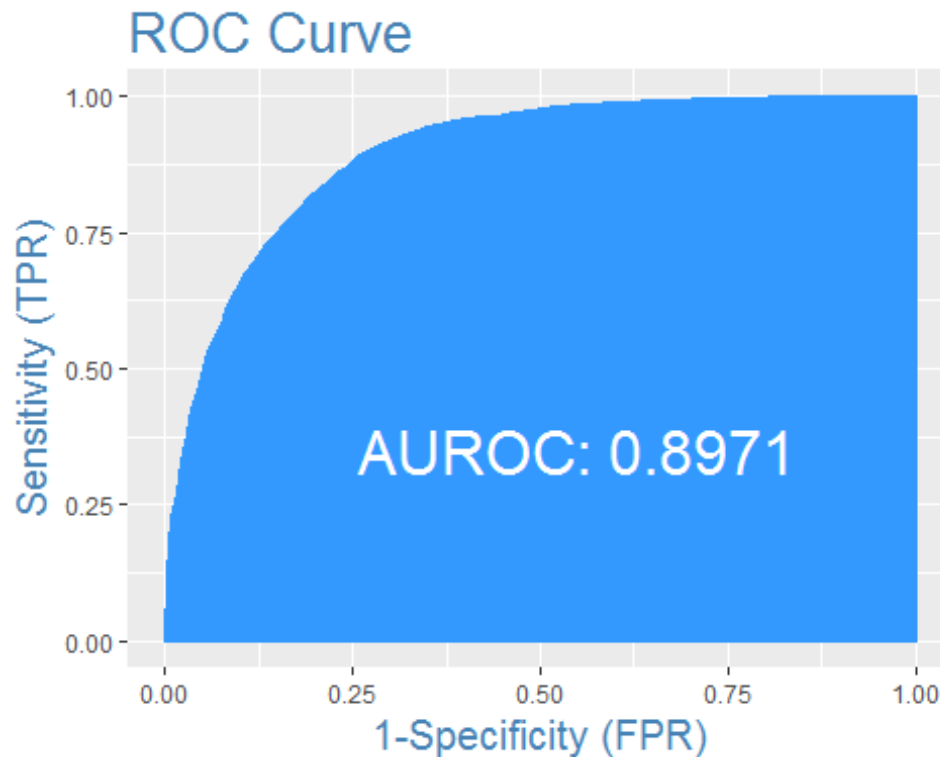
## Start of Stage 5

### 5. Evaluation

We need to evaluate the model using the test data. Evaluation checks a number of parameters for accuracy. In classification problems, we should be checking the following parameters

- ROC: Receiver Operating Characteristics Curve traces the percentage of true positives accurately predicted by a given logit model as the prediction probability cutoff is lowered from 1 to 0. For a good model, as the cutoff is lowered, it should mark more of actual 1's as positives and lesser of actual 0's as 1's.
- So for a good model, the curve should rise steeply, indicating that the TPR (Y-Axis) increases faster than the FPR (X-Axis) as the cutoff score decreases. Greater the area under the ROC curve, better the predictive ability of the model.

```
# The model has area under ROC curve 89.7%, which is pretty good
plotROC(testData$ABOVE50K, predicted)
```



#### Specificity and Sensitivity

- Sensitivity (or True Positive Rate) is the percentage of 1's (actuals) correctly predicted by the model
- Specificity is the percentage of 0's (actuals) correctly predicted.
- Specificity can also be calculated as  $1 - \text{False Positive Rate}$ .

```
sensitivity(testData$ABOVE50K, predicted)
```

```
## [1] 0.8283043
```

```
specificity(testData$ABOVE50K, predicted)
```

```
## [1] 0.796849
```

The above numbers are calculated on the validation sample that was not used for training the model. So, a truth detection rate of 82% on test data is good.

#### Confusion Matrix

```
cm <- as.data.frame(confusionMatrix(testData$ABOVE50K, predicted))
```

```
colnames(cm) <- c("Actual 0", "Actual 1")
```

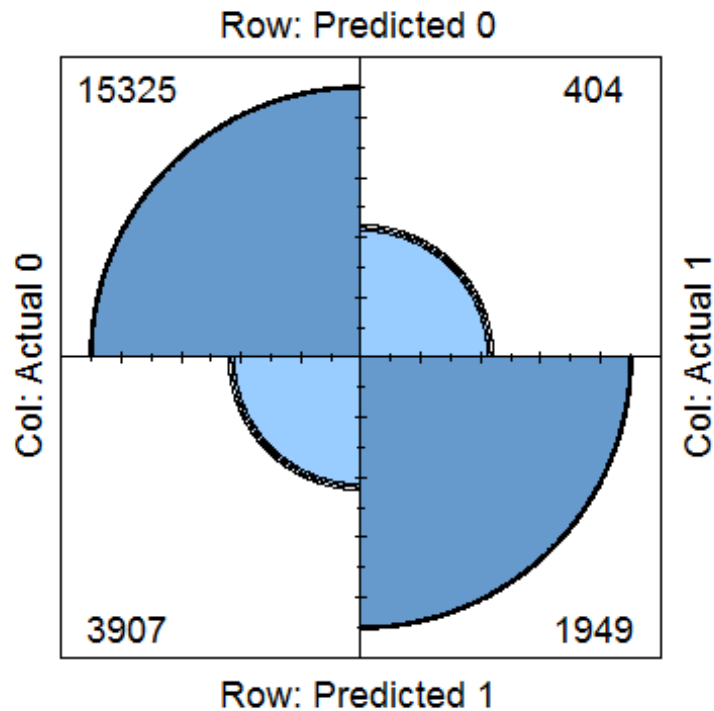
```
rownames(cm) <- c("Predicted 0", "Predicted 1")
```

```
cm
```

```
##           Actual 0 Actual 1
## Predicted 0    15325     404
## Predicted 1     3907    1949
```



```
fourfoldplot(as.matrix(cm))
```



## End of Stage 5

---

---

---

---

## Start of Stage 6

### 6. Deployment

- Creation of the model is generally not the end of the project. Even if the purpose of the model is to increase knowledge of the data, the knowledge gained will need to be organized and presented in a way that the customer can use it.
- It often involves applying "live" models within an organization's decision making processes. For example, real-time personalization of Web pages or repeated scoring of marketing databases.
- Depending on the requirements, the deployment phase can be as simple as generating a report or as complex as implementing a repeatable Data Analytics process across the enterprise.

- In many cases, it is the customer, not the data analyst, who carries out the deployment steps. However, even if the analyst will carry out the deployment effort, it is important for the customer to understand up front what actions need to be carried out in order to actually make use of the created models.
- The deployment of the model will depend on the IT/product architecture, with which it needs to be integrated. The model could run outside the IT/product architecture. The output could be integrated with the system using API or similar interface.
- If the model needs to be integrated with a product (like GTOS), then the product should be able to support ML algorithms. Deployment is driven by IT and engineering team with the support from the data scientist.

## End of Stage 6

---



---



---



---

## Start of Stage 7

### 7. Maintenance and Support

- A Data Analytic product could be created and deployed in less than a year. However, the maintenance and support of the product could run into years.
- This phase is very important because of changing nature of data and processes within an organisation. The data product may require fine tuning to accommodate the new realities.

#### Plan Maintenance and Support Roadmap

- Important if the Data Analytics results become part of the day-to-day business and IT environment
- Helps to avoid unnecessarily long periods of incorrect usage of Data Analytics results
- Needs a detailed plan on monitoring process
- Takes into account the specific type of deployment

## End of the Script