

Godavari College of Engineering, Jalgaon

Subject Name: Machine Learning

Practical No: 02

Date:

Title: Logistic Regression Analysis in R.

Aim: Study and implementation of logistic regression analysis in R.

Theory:

The Logistic Regression is a regression model in which the response variable (dependent variable) has categorical values such as True/False or 0/1. It actually measures the probability of a binary response as the value of response variable based on the mathematical equation relating it with the predictor variables.

The general mathematical equation for logistic regression is-

$$y = 1/(1+e^{-(a+b_1x_1+b_2x_2+b_3x_3+\dots)})$$

Following is the description of the parameters used –

- **y** is the response variable.
- **x** is the predictor variable.
- **a** and **b** are the coefficients which are numeric constants.

The function used to create the regression model is the **glm()** function.

Syntax:

The basic syntax for **glm()** function in logistic regression is-

`glm(formula, data, family)`

Following is the description of the parameters used –

- **formula** is the symbol presenting the relationship between the variables.
- **data** is the data set giving the values of these variables.
- **family** is R object to specify the details of the model. Its value is binomial for logistic regression.

Source Code:

The in-built data set "mtcars" describes different models of a car with their various engine specifications. In "mtcars" data set, the transmission mode (automatic or manual) is described by the column am which is a binary value (0 or 1). We can create a logistic regression model between the columns "am" and 3 other columns - hp, wt and cyl.

#Select some columns from mtcars.

```
Input ← mtcars[,c("am","cyl","hp","wt")]  
print(head(input))
```

#We use the **glm()** function to create the regression model and get its summary for analysis.

```
am.data = glm(formula = am ~ cyl + hp + wt, data = input, family = binomial)
print(summary(am.data))
```

Output:

```
> input <- mtcars[,c("am","cyl","hp","wt")]
> 
> print(head(input))
      am  cyl  hp   wt
Mazda RX4      1   6 110 2.620
Mazda RX4 Wag  1   6 110 2.875
Datsun 710     1   4  93 2.320
Hornet 4 Drive 0   6 110 3.215
Hornet Sportabout 0  8 175 3.440
Valiant        0   6 105 3.460
> input <- mtcars[,c("am","cyl","hp","wt")]
> 
> am.data = glm(formula = am ~ cyl + hp + wt, data = input, family = binomial)
> 
> print(summary(am.data))

Call:
glm(formula = am ~ cyl + hp + wt, family = binomial, data = input)

Deviance Residuals:
    Min       1Q   Median       3Q      Max
-2.17272  -0.14907  -0.01464   0.14116   1.27641

Coefficients:
            Estimate Std. Error z value Pr(>|z|)
(Intercept)  19.70288    8.11637   2.428  0.0152 *
cyl           0.48760    1.07162   0.455  0.6491
hp            0.03259    0.01886   1.728  0.0840 .
wt           -9.14947    4.15332  -2.203  0.0276 *
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

    Null deviance: 43.2297  on 31  degrees of freedom
Residual deviance:  9.8415  on 28  degrees of freedom
AIC: 17.841

Number of Fisher Scoring iterations: 8
> 
```

Conclusion:

In the summary as the p-value in the last column is more than 0.05 for the variables "cyl" and "hp", we consider them to be insignificant in contributing to the value of the variable "am". Only weight (wt) impacts the "am" value in this regression model.