



Abstract: Linear Regression Model in R

ACD_ANR_PROJECT2.1

DONE BY : ASHWINI K
MENTOR NAME : ARVIND

Problem Statement

Imagine that the CEO of a DVD player sales company approaches you in order to predict the sale of DVDs. He also provides you the data such as the advertising budget (in thousands), sales (in thousands), number of times the song is played on the radio channel, Radio Mirchi per week and the attractiveness of the brand (rated on a scale of 1 to 10 by an independent agency).

Dataset -

:



Approach

: Splitting the data into 70 and 30

Inference

:

1 : it's a MLR, one response variable being sales and explanatory variables being advertise, attractiveness and plays.

2: $\text{sales} = \text{advertise}X1 + \text{attractiveness}X2 + \text{plays}X3$

3: Correlation between response and explanatory variables are having moderate

Source code with comments:

#Reading the data#

```
dvdsales<-read.csv("C:/All/R language/Final Project/fwdproject2/Sales_dataset.csv",header=TRUE)
dvdsales
View(dvdsales)
names(dvdsales)
class(dvdsales$attractiveness)
```

#Correlation between predictors and response variable#

```
cor(dvdsales$sales,dvdsales$advertise)
cor(dvdsales$sales,dvdsales$plays)
cor(dvdsales$sales,dvdsales$attractiveness)
```

#Conversion from numeric to factor

```
dvdsales$attractiveness<-as.factor(dvdsales$attractiveness)
class(dvdsales$attractiveness)
summary(dvdsales)
nrow(dvdsales)
```

#Split the data into 70 and 30

```
set.seed(1)
dvdds<-sample(nrow(dvdsales),nrow(dvdsales)*0.7)
```

#Conversion of changing qualifying into levels

```
dvdsales$attractiveness1<-ifelse(dvdsales$attractiveness< 3,1,ifelse(dvdsales$attractiveness<
5,2,ifelse(dvdsales$attractiveness< 7,3,ifelse(dvdsales$attractiveness< 9,4,5))))
dvdsales$attractiveness1<-as.factor(dvdsales$attractiveness1)
dvdsales$attractiveness<-as.factor(dvdsales$attractiveness)
summary(dvdsales)
```

#attractiveness isn't required since it's been qualifies into levels

```
dvdtrain<-dvdsales[dvdds,-4]
summary(dvdtrain)
dvdtest<-dvdsales[-dvdds,-4]
summary(dvdtest)
```

#target variable is not required in test and train input

```
dvdt<-dvdtest[, -2]
dvdt
```

#Model Creation

```
smod1<-lm(sales~.,data=dvdtrain)
summary(smod1)
```

#Verifying Assumptions in model

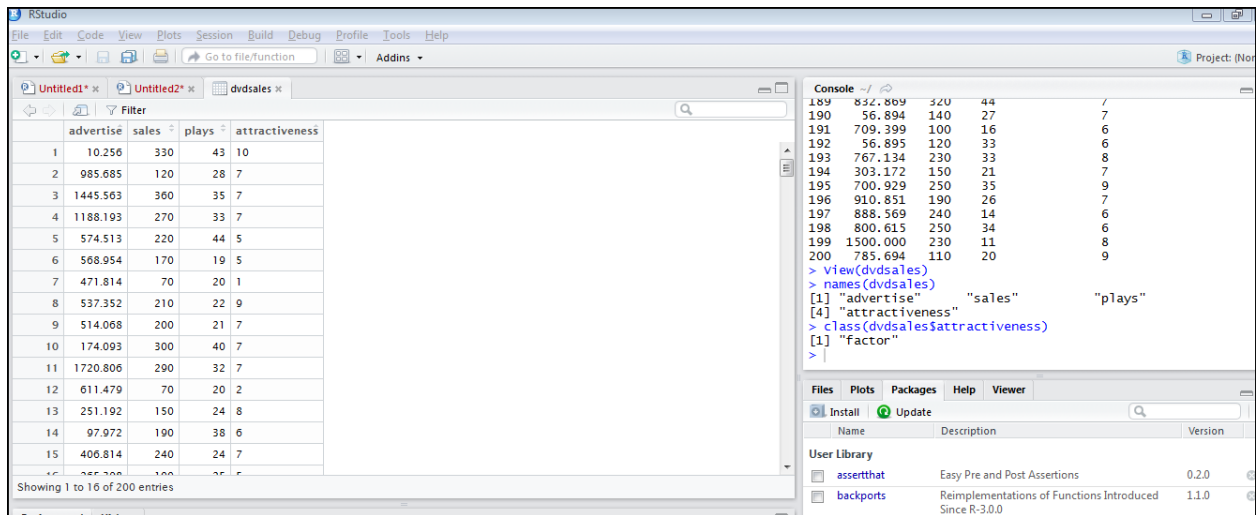
```
plot(smod1)
```

#Prediction in test data

```
dvdt$testres<-predict(smod1,newdata = dvdt)
summary(dvdt)
dvdt$testres
```

Screenshot's for the solution after running code :

dvdsales



The screenshot shows the RStudio interface. The main window displays a data frame with columns: advertise, sales, plays, and attractiveness. The console shows the output of the View function, which displays the first 20 rows of the data frame.

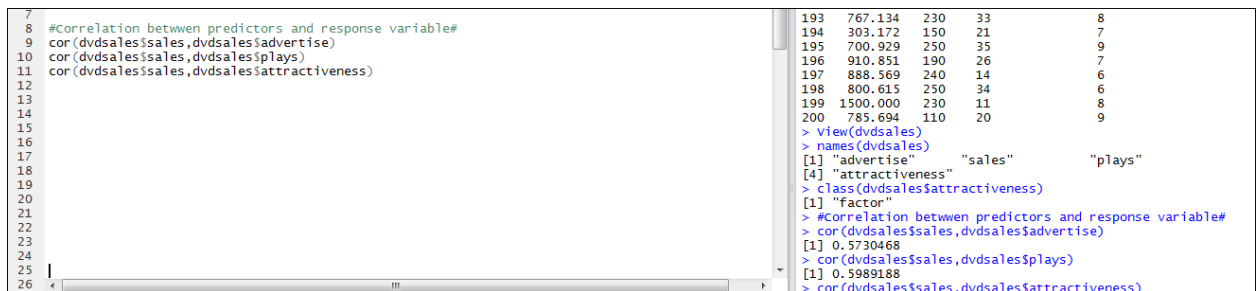
	advertise	sales	plays	attractiveness
1	10.256	330	43	10
2	985.685	120	28	7
3	1445.563	360	35	7
4	1188.193	270	33	7
5	574.513	220	44	5
6	568.954	170	19	5
7	471.814	70	20	1
8	537.352	210	22	9
9	514.068	200	21	7
10	174.093	300	40	7
11	1720.806	290	32	7
12	611.479	70	20	2
13	251.192	150	24	8
14	97.972	190	38	6
15	406.814	240	24	7

Showing 1 to 16 of 200 entries

Console output:

```
189 832.869 320 44 7
190 56.894 140 27 7
191 709.399 100 16 6
192 56.895 120 33 6
193 767.134 230 33 8
194 303.172 150 21 7
195 700.929 250 35 9
196 910.851 190 26 7
197 888.569 240 14 6
198 800.615 250 34 6
199 1500.000 230 11 8
200 785.694 110 20 9
> View(dvdsales)
> names(dvdsales)
[1] "advertise" "sales" "plays"
[4] "attractiveness"
> class(dvdsales$attractiveness)
[1] "factor"
> |
```

correlation between predictors



The screenshot shows the RStudio console with the following code and output:

```
7
8 #Correlation between predictors and response variable#
9 cor(dvdsales$sales,dvdsales$advertise)
10 cor(dvdsales$sales,dvdsales$plays)
11 cor(dvdsales$sales,dvdsales$attractiveness)
12
13
14
15
16
17
18
19
20
21
22
23
24
25 |
26
```

Console output:

```
193 767.134 230 33 8
194 303.172 150 21 7
195 700.929 250 35 9
196 910.851 190 26 7
197 888.569 240 14 6
198 800.615 250 34 6
199 1500.000 230 11 8
200 785.694 110 20 9
> View(dvdsales)
> names(dvdsales)
[1] "advertise" "sales" "plays"
[4] "attractiveness"
> class(dvdsales$attractiveness)
[1] "factor"
> #Correlation between predictors and response variable#
> cor(dvdsales$sales,dvdsales$advertise)
[1] 0.5730468
> cor(dvdsales$sales,dvdsales$plays)
[1] 0.5989188
> cor(dvdsales$sales,dvdsales$attractiveness)
```

conversion from numeric to factor

```

1 #creating the data#
2 dvdsales<-read.csv("C:/All/R language/Final Project/fwdproject2/Sales_dataset.csv",header=TRUE)
3 dvdsales
4 view(dvdsales)
5 names(dvdsales)
6 class(dvdsales$attractiveness)
7
8 #Correlation between predictors and response variable#
9 cor(dvdsales$sales,dvdsales$advertise)
10 cor(dvdsales$sales,dvdsales$plays)
11 cor(dvdsales$sales,dvdsales$attractiveness)
12
13
14
15
16
17
18 #Conversion from numeric to factor
19 dvdsales$attractiveness<-as.factor(dvdsales$attractiveness)
20 class(dvdsales$attractiveness)
21 view(dvdsales)
22 summary(dvdsales)
23 nrow(dvdsales)
24

```

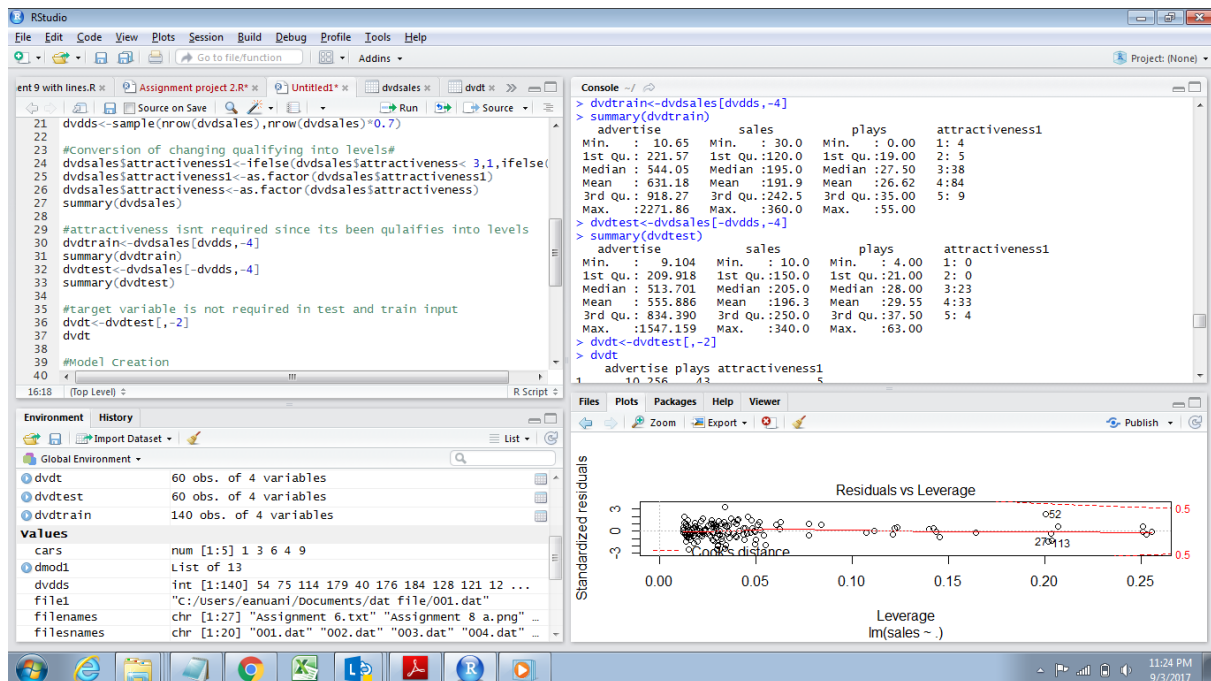
```

> class(dvdsales$attractiveness)
[1] "factor"
> view(dvdsales)
Error in view(dvdsales) : could not find function "view"
> summary(dvdsales)
  advertise      sales      plays
Min.   : 9.104   Min.   :10.0   Min.   :0.00
1st Qu.:215.609 1st Qu.:137.5   1st Qu.:19.75
Median :526.311 Median :200.0   Median :28.00
Mean   :608.592 Mean   :193.2   Mean   :27.50
3rd Qu.:903.380 3rd Qu.:250.0   3rd Qu.:36.00
Max.   :2271.860 Max.   :360.0   Max.   :63.00

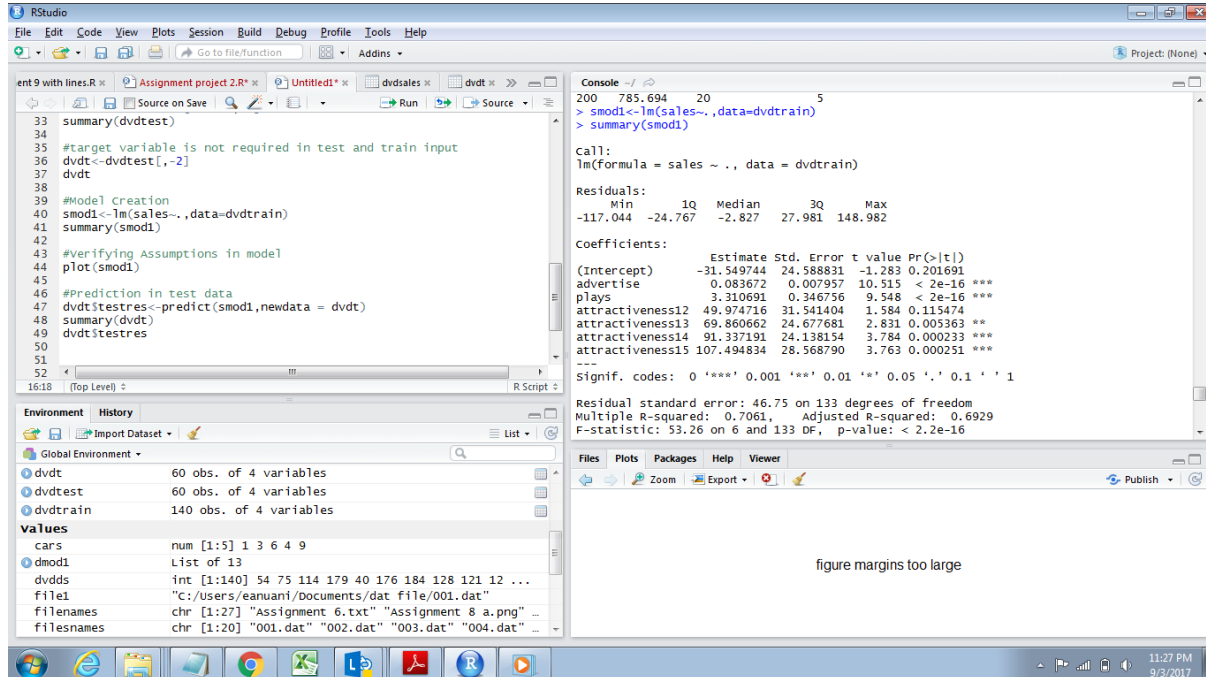
  attractiveness
 7      :73
 6      :44
 8      :44
 5      :17
 9      :12
 4      : 4
 (other): 6
> nrow(dvdsales)
[1] 200
>

```

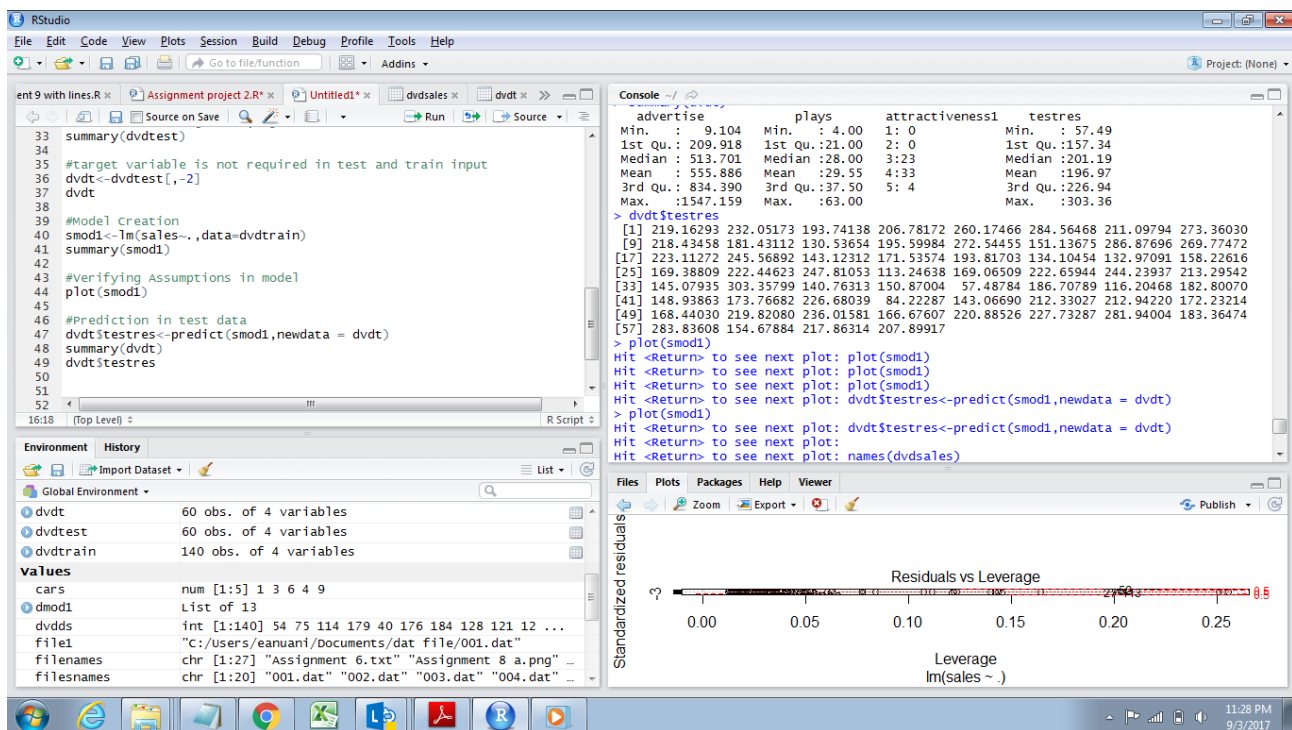
Attractiveness response removed



Model Creation



testres



CONCLUSION :

- Since P value is less than 0.5 the relationship is linear.
- Predictors, attractiveness 12 has no influence on target variable
- Advertise and plays have minimal error has their standard error is less
- Assumptions criteria are verified.

Predictors have direct relationship with target.

THANK YOU