

Multiple Linear Regression Model

Sachin Sharma

12/27/2021

Installing the libraries

```
library(tinytex)
library(ggplot2)
library(tidyverse)

## -- Attaching packages ----- tidyverse 1.3.1 --
## v tibble 3.1.6      v dplyr 1.0.7
## v tidyr 1.1.4      v stringr 1.4.0
## v readr 2.0.1      v forcats 0.5.1
## v purrr 0.3.4
## Warning: package 'tibble' was built under R version 4.1.2
## -- Conflicts ----- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()     masks stats::lag()
library(rvest)

## Warning: package 'rvest' was built under R version 4.1.2
##
## Attaching package: 'rvest'
## The following object is masked from 'package:readr':
##
##     guess_encoding
library(naniar)
library(corrplot)

## corrplot 0.90 loaded
```

In multiple linear regression the equation is like :

$$y = a_0 + a_1x_1 + a_2x_2 + \dots$$

where a_0 is the y intercept, and a_1 is slope , which can be compared with linear regression : $y = mx + c$
where c is y intercept and m is slope

In multiple linear regression model, we will plot scatter plot first to understand the relation between variables, whether or not the variables are linearly correlated:

Loading data:

You can download the dataset from kaggle: <https://www.kaggle.com/nehalbirla/vehicle-dataset-from-cardekho>

```
vehicle <- read.csv("car.csv")

head(vehicle)

##           name year selling_price km_driven  fuel seller_type
## 1   Maruti 800 AC 2007         60000    70000 Petrol  Individual
## 2 Maruti Wagon R LXI Minor 2007      135000    50000 Petrol  Individual
## 3   Hyundai Verna 1.6 SX 2012     600000   100000 Diesel  Individual
## 4   Datsun RediGO T Option 2017     250000    46000 Petrol  Individual
## 5   Honda Amaze VX i-DTEC 2014     450000   141000 Diesel  Individual
## 6   Maruti Alto LX BSIII 2007     140000   125000 Petrol  Individual
## transmission      owner
## 1      Manual  First Owner
## 2      Manual  First Owner
## 3      Manual  First Owner
## 4      Manual  First Owner
## 5      Manual Second Owner
## 6      Manual  First Owner

colnames(vehicle)

## [1] "name"          "year"          "selling_price" "km_driven"
## [5] "fuel"          "seller_type"   "transmission"  "owner"
```

Lets make scatter plot to see how strongly variables are correlated, we are interested in Mileage, lh and lc

```
str(vehicle)

## 'data.frame':   4340 obs. of  8 variables:
## $ name      : chr  "Maruti 800 AC" "Maruti Wagon R LXI Minor" "Hyundai Verna 1.6 SX" "Datsun RediGO T Option" ...
## $ year      : int   2007 2007 2012 2017 2014 2007 2016 2014 2015 2017 ...
## $ selling_price: int   60000 135000 600000 250000 450000 140000 550000 240000 850000 365000 ...
## $ km_driven  : int   70000 50000 100000 46000 141000 125000 25000 60000 25000 78000 ...
## $ fuel       : chr   "Petrol" "Petrol" "Diesel" "Petrol" ...
## $ seller_type: chr   "Individual" "Individual" "Individual" "Individual" ...
## $ transmission: chr   "Manual" "Manual" "Manual" "Manual" ...
## $ owner      : chr   "First Owner" "First Owner" "First Owner" "First Owner" ...

# we will first convert fuel,seller_type,transmission,owner in factor variable

vehicle$fuel <- as.factor(vehicle$fuel)
vehicle$seller_type <- as.factor(vehicle$seller_type)
```

```
vehicle$transmission <- as.factor(vehicle$transmission)
vehicle$owner <- as.factor(vehicle$owner)
```

Now converting the factor to numeric to check the correlation between variables

```
vehicle$fuel <- as.numeric(vehicle$fuel)
vehicle$seller_type <- as.numeric(vehicle$seller_type)
vehicle$transmission <- as.numeric(vehicle$transmission)
vehicle$owner <- as.numeric(vehicle$owner)

vehicle$name <- as.factor(vehicle$name)

vehicle$name <- as.numeric(vehicle$name)
```

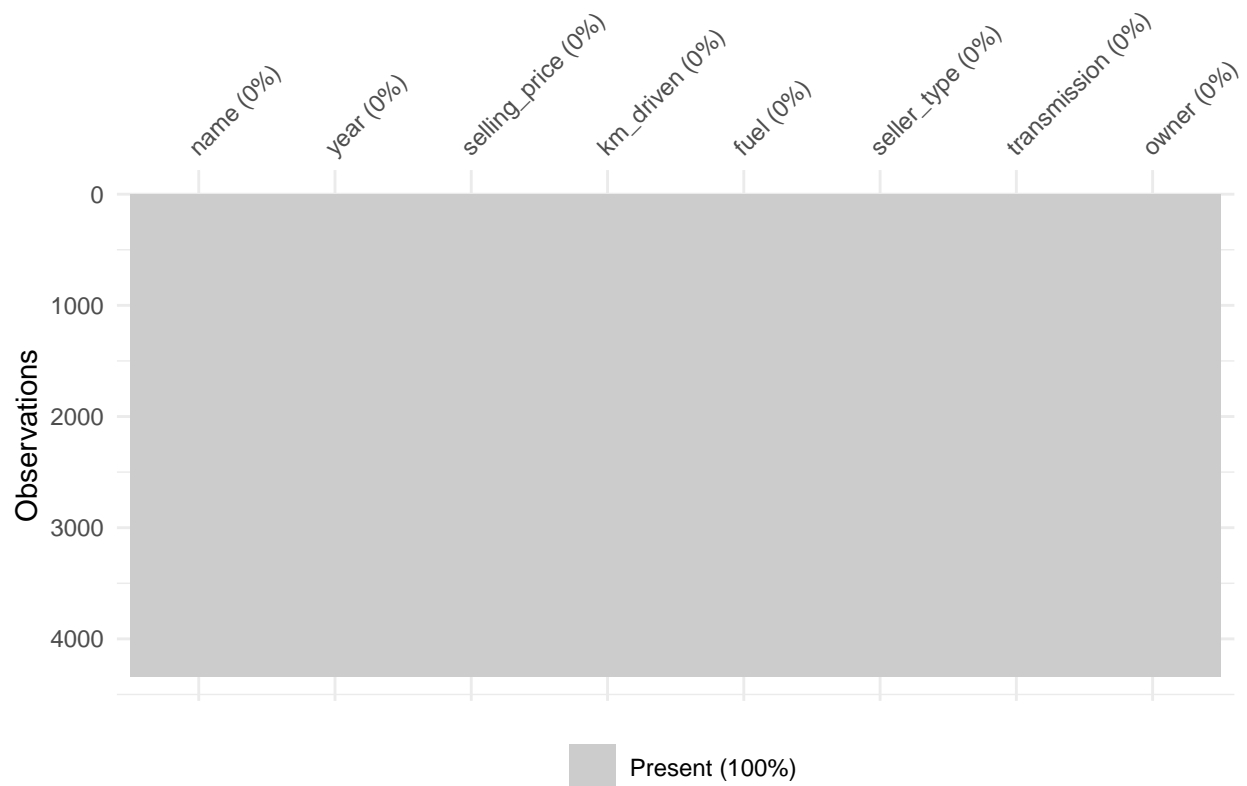
Now checking the missing values in the data if any :

```
sum(is.na(vehicle))
```

```
## [1] 0
```

Visualizing the missing values if any :

```
vis_miss(vehicle)
```

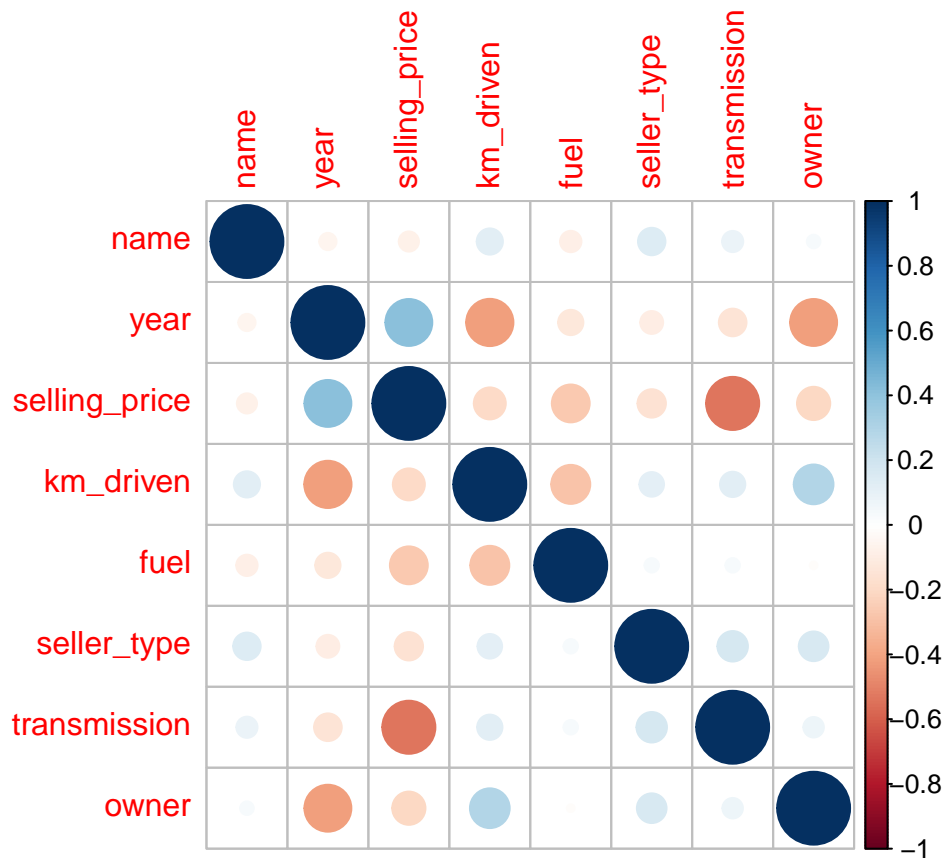


#Here we can see that there is no missing values in the data :

Lets draw a corrplot to see how varibales are related to each other :

```
cort <- cor(vehicle)

corrplot(cort)
```



```
colnames(vehicle)
```

```
## [1] "name"      "year"      "selling_price" "km_driven"
## [5] "fuel"      "seller_type" "transmission" "owner"
```

```
model_lm_1 <- lm(selling_price~.,data = vehicle)
```

```
summary(model_lm_1)
```

```
##
## Call:
## lm(formula = selling_price ~ ., data = vehicle)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1149547  -163275  -27635   115782   7527991
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -7.106e+07  3.803e+06 -18.687 < 2e-16 ***
## name        -4.301e+01  1.664e+01  -2.585  0.00976 **
## year         3.663e+04  1.881e+03  19.467 < 2e-16 ***
## km_driven    -9.648e-01  1.689e-01  -5.712 1.19e-08 ***
## fuel        -9.358e+04  4.717e+03 -19.837 < 2e-16 ***
## seller_type  -1.947e+04  1.477e+04  -1.318  0.18753
## transmission -8.838e+05  2.203e+04 -40.118 < 2e-16 ***
## owner       -1.710e+04  5.926e+03  -2.886  0.00392 **
```

```
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 429500 on 4332 degrees of freedom
## Multiple R-squared:  0.4498, Adjusted R-squared:  0.449
## F-statistic: 506 on 7 and 4332 DF,  p-value: < 2.2e-16

model_lm <- lm(selling_price~year+km_driven+fuel+seller_type+transmission+owner,data = vehicle)

summary(model_lm)

##
## Call:
## lm(formula = selling_price ~ year + km_driven + fuel + seller_type +
##     transmission + owner, data = vehicle)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1171869  -169573   -25668    119142   7549809
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -7.122e+07  3.805e+06 -18.718 < 2e-16 ***
## year           3.670e+04  1.883e+03  19.493 < 2e-16 ***
## km_driven     -9.974e-01  1.685e-01  -5.918 3.51e-09 ***
## fuel          -9.278e+04  4.710e+03 -19.698 < 2e-16 ***
## seller_type   -2.425e+04  1.467e+04  -1.653 0.09835 .
## transmission -8.868e+05  2.201e+04 -40.285 < 2e-16 ***
## owner         -1.676e+04  5.928e+03  -2.827 0.00472 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 429800 on 4333 degrees of freedom
## Multiple R-squared:  0.449, Adjusted R-squared:  0.4482
## F-statistic: 588.5 on 6 and 4333 DF,  p-value: < 2.2e-16

gc()

##           used (Mb) gc trigger (Mb) max used (Mb)
## Ncells 1213946 64.9   2027054 108.3  2027054 108.3
## Vcells 2382295 18.2    8388608  64.0   8388606  64.0
```

*** represents the significance of variable in the model.**

##*** represent highly significant##** represent significant##* represent less significant## no star :
represent no significance

Here we can see that seller_type is not significant for the model, we will remove this and update our model as follows :

```
model_lm1 <- lm(selling_price~year+km_driven+fuel+transmission+owner,data = vehicle)

summary(model_lm1)
```

```
##
## Call:
## lm(formula = selling_price ~ year + km_driven + fuel + transmission +
##      owner, data = vehicle)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1160675  -163426   -26086   115456   7562789
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -7.113e+07  3.805e+06 -18.692 < 2e-16 ***
## year         3.664e+04  1.883e+03  19.460 < 2e-16 ***
## km_driven    -1.016e+00  1.682e-01  -6.044 1.63e-09 ***
## fuel         -9.323e+04  4.704e+03 -19.820 < 2e-16 ***
## transmission -8.925e+05  2.175e+04 -41.036 < 2e-16 ***
## owner        -1.802e+04  5.880e+03  -3.064  0.0022 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 429800 on 4334 degrees of freedom
## Multiple R-squared:  0.4487, Adjusted R-squared:  0.448
## F-statistic: 705.3 on 5 and 4334 DF,  p-value: < 2.2e-16

model_lm2 <- lm(selling_price~year+km_driven+fuel+transmission,data = vehicle)

summary(model_lm2)

##
## Call:
## lm(formula = selling_price ~ year + km_driven + fuel + transmission,
##      data = vehicle)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1152345  -165847   -25623   114607   7567473
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -7.492e+07  3.602e+06 -20.799 < 2e-16 ***
## year         3.850e+04  1.783e+03  21.599 < 2e-16 ***
## km_driven    -1.085e+00  1.669e-01  -6.502 8.85e-11 ***
## fuel         -9.304e+04  4.708e+03 -19.762 < 2e-16 ***
## transmission -8.933e+05  2.177e+04 -41.037 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 430300 on 4335 degrees of freedom
## Multiple R-squared:  0.4475, Adjusted R-squared:  0.4469
## F-statistic: 877.6 on 4 and 4335 DF,  p-value: < 2.2e-16

model_lm3 <- lm(selling_price~year+km_driven+transmission,data = vehicle)

summary(model_lm3)
```

```
##
## Call:
## lm(formula = selling_price ~ year + km_driven + transmission,
##     data = vehicle)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1262774  -164383   -31830   103566   7443587
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -9.445e+07  3.616e+06 -26.119  <2e-16 ***
## year         4.803e+04  1.792e+03  26.803  <2e-16 ***
## km_driven    1.543e-01  1.614e-01   0.956    0.339
## transmission -9.153e+05  2.269e+04 -40.329  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 449200 on 4336 degrees of freedom
## Multiple R-squared:  0.3977, Adjusted R-squared:  0.3973
## F-statistic: 954.3 on 3 and 4336 DF,  p-value: < 2.2e-16
```