# Electric Vehicle Pricing

Eric Wilson; Min Thiha Myo; Kevin Maldonado

12/7/2021

## Introduction

Electric Vehicles are rising in popularity amongst consumers in todays day and age. With climate change playing a big factor, several nations have made a consorted effort in investing in 'green technologies' which include electric vehicles. The development and improvement of battery technology has made the electric vehicle a viable and perhaps a more preferable choice for consumers. Some of the advantages of electric vehicles include saving money on gas, environmentally friendly, low maintenance cost. Local governments are also offering incentives to buy an electric vehicle. As there is with any new technology, there are challenges that electric vehicles face. Charging station scarcity and battery range are 2 big hurdles that owners will have to face. Despite the challenges that owners and manufactures might have to overcome, the rise of ownership of electric vehicles is undeniable and will only increase heading into the future.

Our project will take data from various electric vehicles that are available in the market and from that data we will perform a regression analysis to predict the price of electric vehicles based on different variables such as battery size, acceleration, range and efficiency.

## Data Preparation

### Data Source

Our data can be found through this link: https://www.kaggle.com/kkhandekar/cheapest-electric-cars

##Loading Dataset Our EV dataset contains 180 observations (rows) and 12 variables (columns).

### Data Cleanup

We now go into our data and start our clean up process. We start by converting several variables to numeric. Battery Size, PriceinUK, Acceleration, Top Speed, Range, Efficiency, Fast charge speed, and Price in Germany all get converted to numeric. Along with this conversion, we perform some regex (regular expression) to our data so that we can perform a successful conversion to numeric. Our 'Drive' variable will be converted to a factor. Note that since our data set comes from european sources the units will be in kilometers.

We then check for any missing values in our data set and remove them.

## Exploratory Data Analysis

Let's to a look at our data set now that we have removed missing values and cleaned it up.

```
## # A tibble: 6 x 11
##   Name            Acceleration TopSpeed Range Efficiency FastChargeSpeed Drive
##   <chr>                  <dbl>    <dbl> <dbl>      <dbl>           <dbl> <fct>
## 1 Opel Ampera-e            7.3      150   335        173             210 Front W~
## 2 Nissan Leaf              7.9      144   220        164             230 Front W~
## 3 Porsche Tayca~           2.8      260   390        215             860 All Whe~
## 4 Nissan e-NV20~          14        123   165        218             170 Front W~
## 5 Volkswagen ID~           8.9      160   275        164             260 Rear Wh~
## 6 BMW iX3                  6.8      180   385        192             520 Rear Wh~
## # ... with 4 more variables: NumberofSeats <dbl>, PriceinGermany <dbl>,
## #   PriceinUK <dbl>, BatterySize <dbl>
```

We have data available for EV prices in Germany (in euros) and prices in the UK (pound sterling ), but for the purpose of this report we will only be using prices in Germany as our response variable since it provides more data points.

## Summary of Variables

```
##    Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##   23.80   45.00   67.25   66.73   78.05  200.00


##    Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##   20490   38750   50890   59645   65096  215000


##    Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##   123.0   150.0   160.0   177.7   200.0   410.0


##    Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##   165.0   270.0   340.0   344.1   400.0   970.0


##   Acceleration       TopSpeed          Range         Efficiency
##   Min.   : 2.100  Min.   :123.0  Min.   :165.0  Min.   :104.0
##   1st Qu.: 5.400  1st Qu.:150.0  1st Qu.:270.0  1st Qu.:168.8
##   Median : 7.500  Median :160.0  Median :340.0  Median :189.0
##   Mean   : 7.670  Mean   :177.7  Mean   :344.1  Mean   :195.3
##   3rd Qu.: 9.125  3rd Qu.:200.0  3rd Qu.:400.0  3rd Qu.:217.2
##   Max.   :15.000  Max.   :410.0  Max.   :970.0  Max.   :281.0
##   FastChargeSpeed   BatterySize
##   Min.   : 120.0  Min.   : 23.80
##   1st Qu.: 290.0  1st Qu.: 45.00
##   Median : 440.0  Median : 67.25
##   Mean   : 486.6  Mean   : 66.73
##   3rd Qu.: 622.5  3rd Qu.: 78.05
##   Max.   :1410.0  Max.   :200.00
```
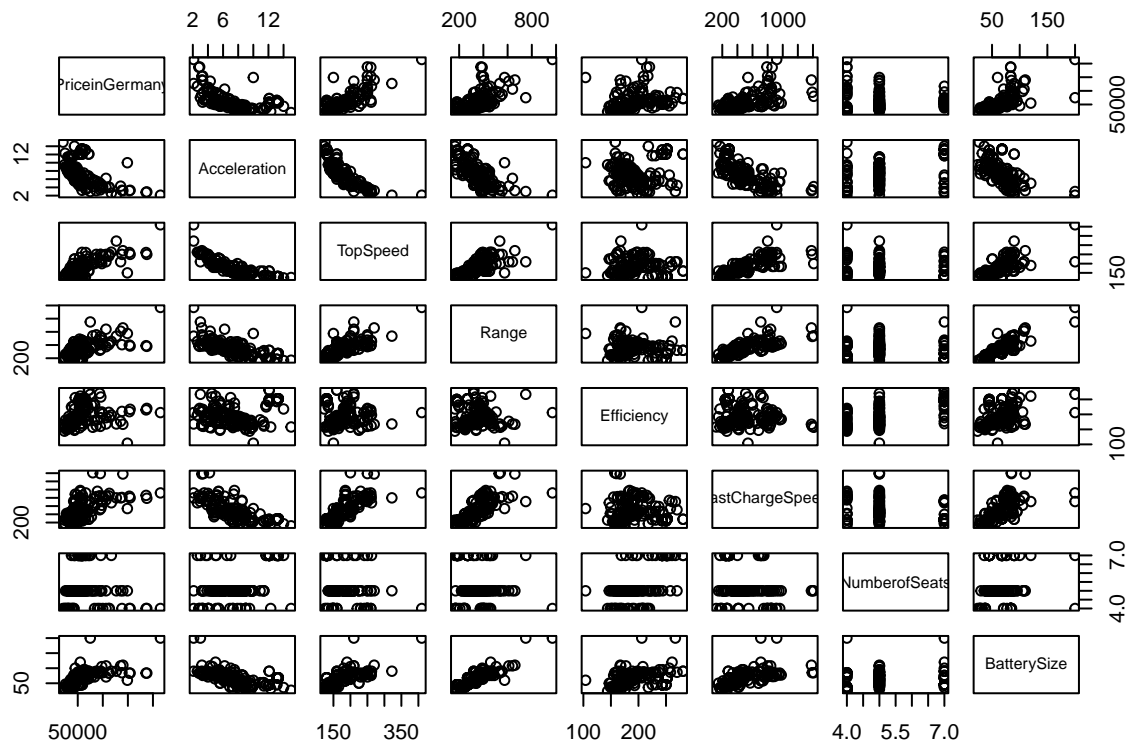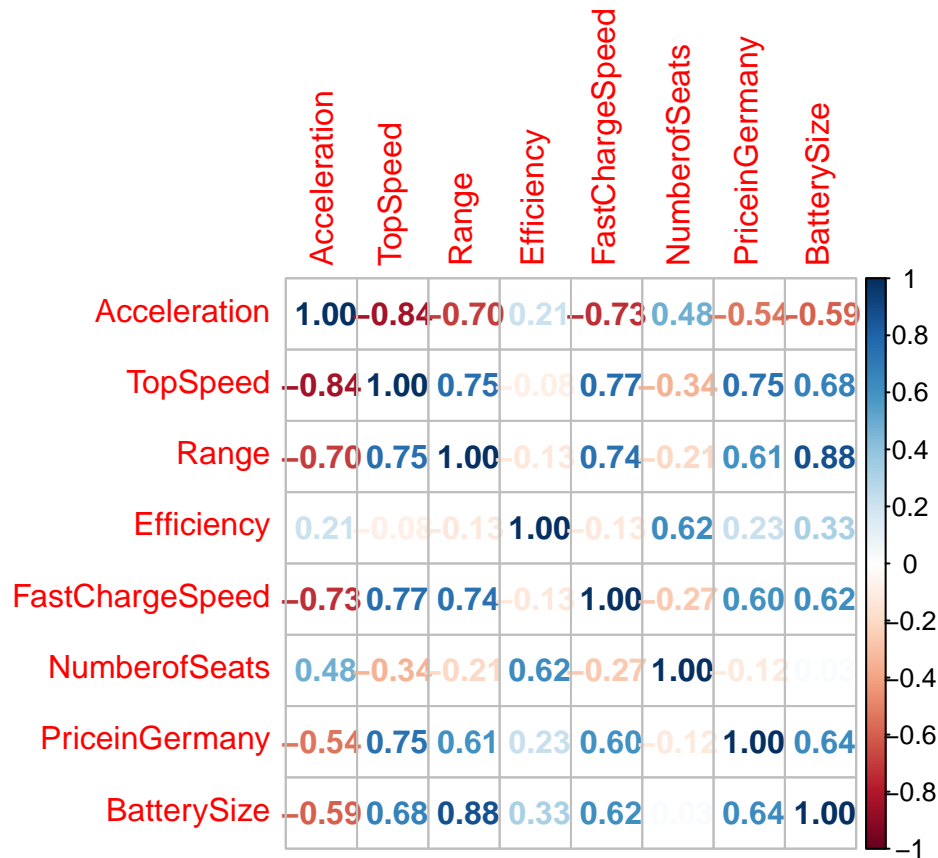
The electric vehicles in our data set range in price from €20,490 to €215,000 with an average price of €59,645. Battery size ranges from 23.8kWh to 200kWh with an average battery size of 66.73kWh The battery range ranges from 165km to 970km Top speed ranges from 123km/h to 410km/h with an average top speed of 177.7km/h

# Correlation



This scatter plot matrices shows the correlation that exists between variables. Starting from the top left, the variables move downward diagonally and is an axis label for the plots shown. With these scatterplots we can see that there is some correlation between PriceinGermany and top speed, range, fast charging speed and battery size. These plots might be a little difficult to read so lets do a different correlation plot.
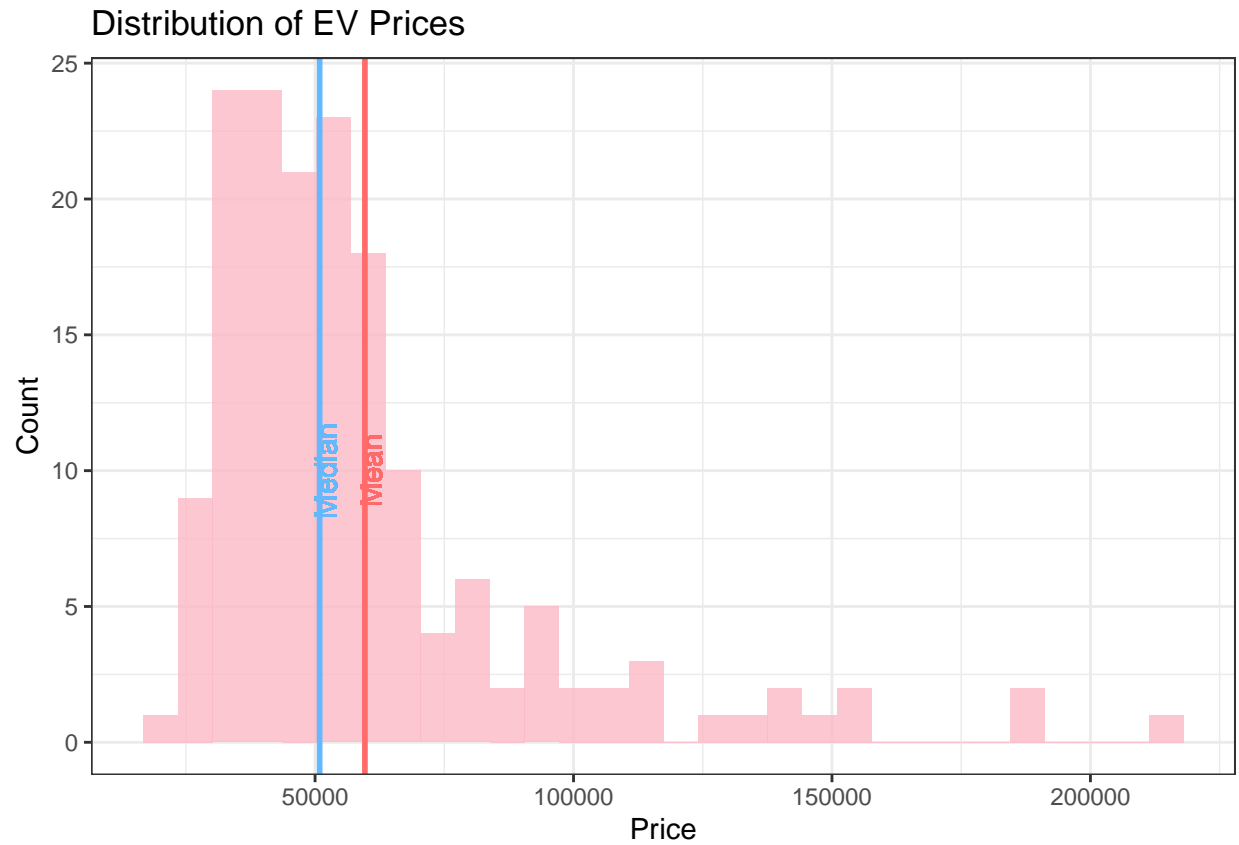
|  | Acceleration | TopSpeed | Range | Efficiency | FastChargeSpeed | NumberofSeats | PriceinGermany | BatterySize |
|---|---|---|---|---|---|---|---|---|
| Acceleration | 1.00 | -0.84 | -0.70 | 0.21 | -0.73 | 0.48 | -0.54 | -0.59 |
| TopSpeed | -0.84 | 1.00 | 0.75 | -0.08 | 0.77 | -0.34 | 0.75 | 0.68 |
| Range | -0.70 | 0.75 | 1.00 | -0.13 | 0.74 | -0.21 | 0.61 | 0.88 |
| Efficiency | 0.21 | -0.08 | -0.13 | 1.00 | -0.13 | 0.62 | 0.23 | 0.33 |
| FastChargeSpeed | -0.73 | 0.77 | 0.74 | -0.13 | 1.00 | -0.27 | 0.60 | 0.62 |
| NumberofSeats | 0.48 | -0.34 | -0.21 | 0.62 | -0.27 | 1.00 | -0.12 | 0.03 |
| PriceinGermany | -0.54 | 0.75 | 0.61 | 0.23 | 0.60 | -0.12 | 1.00 | 0.64 |
| BatterySize | -0.59 | 0.68 | 0.88 | 0.33 | 0.62 | 0.03 | 0.64 | 1.00 |

This correlation plot is also a correlation matrix. We us the argument "method = number" to show the coefficients in a number with different colors that describe their respective correlation. Confirming with what we saw in the scatterplot matrix, we can easily identify the correlation between PriceinGermany and top speed, range, fast charge speed and battery size.
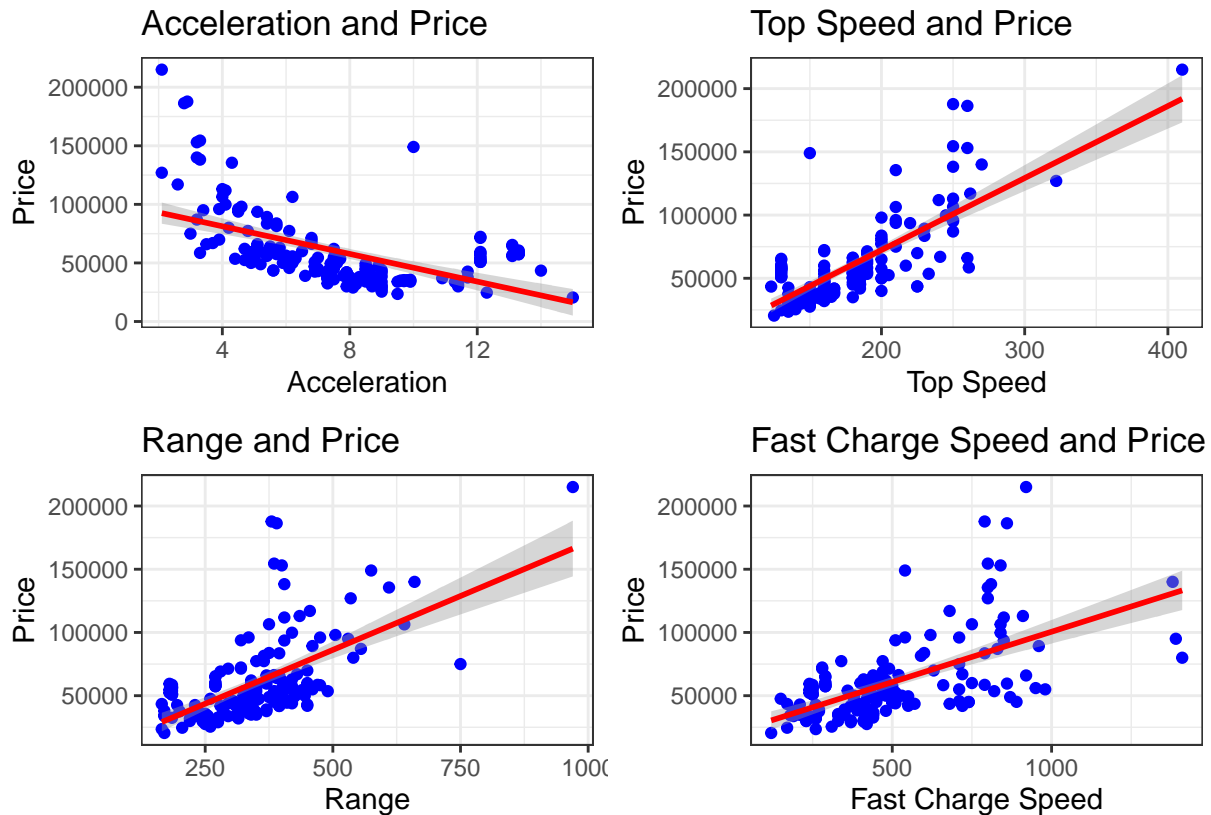
## Distribution of EV Prices

Taking a more in depth look in EV prices, we will graph the distribution of the prices as well as mark the median and mean

```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```

## Distribution of EV Prices



Now we will look at the relationship price has with acceleration, top speed, range and efficiency

```
## 'geom_smooth()' using formula 'y ~ x'
## 'geom_smooth()' using formula 'y ~ x'
## 'geom_smooth()' using formula 'y ~ x'
## 'geom_smooth()' using formula 'y ~ x'
```
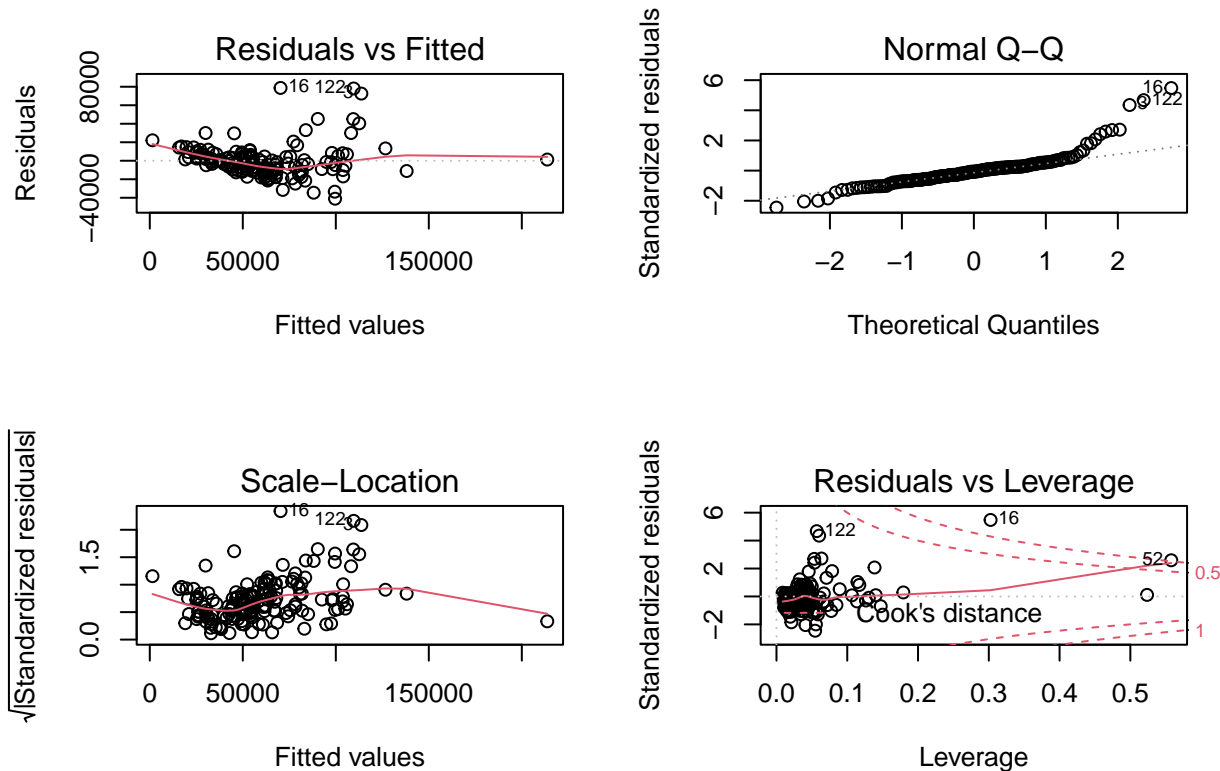
## Model Selection and Evaluation

Let's fit our data into the simplest model and perform residual analysis

```
##
## Call:
## lm(formula = PriceinGermany ~ ., data = ev_numeric)
##
## Residuals:
##    Min     1Q Median     3Q    Max
## -41048  -9454  -1264   4493  78697
##
## Coefficients:
##                  Estimate Std. Error t value Pr(>|t|)
## (Intercept)     -2.358e+05  2.510e+04  -9.394  < 2e-16 ***
## Acceleration     2.034e+03  9.816e+02   2.072   0.0399 *
## TopSpeed         5.398e+02  6.838e+01   7.893 4.86e-13 ***
## Range            4.547e+02  7.441e+01   6.110 7.64e-09 ***
## Efficiency       1.093e+03  1.410e+02   7.747 1.12e-12 ***
## FastChargeSpeed -2.057e+00  1.002e+01  -0.205   0.8376
## NumberofSeats   -9.750e+03  2.241e+03  -4.351 2.44e-05 ***
## BatterySize     -2.006e+03  3.582e+02  -5.602 9.34e-08 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```
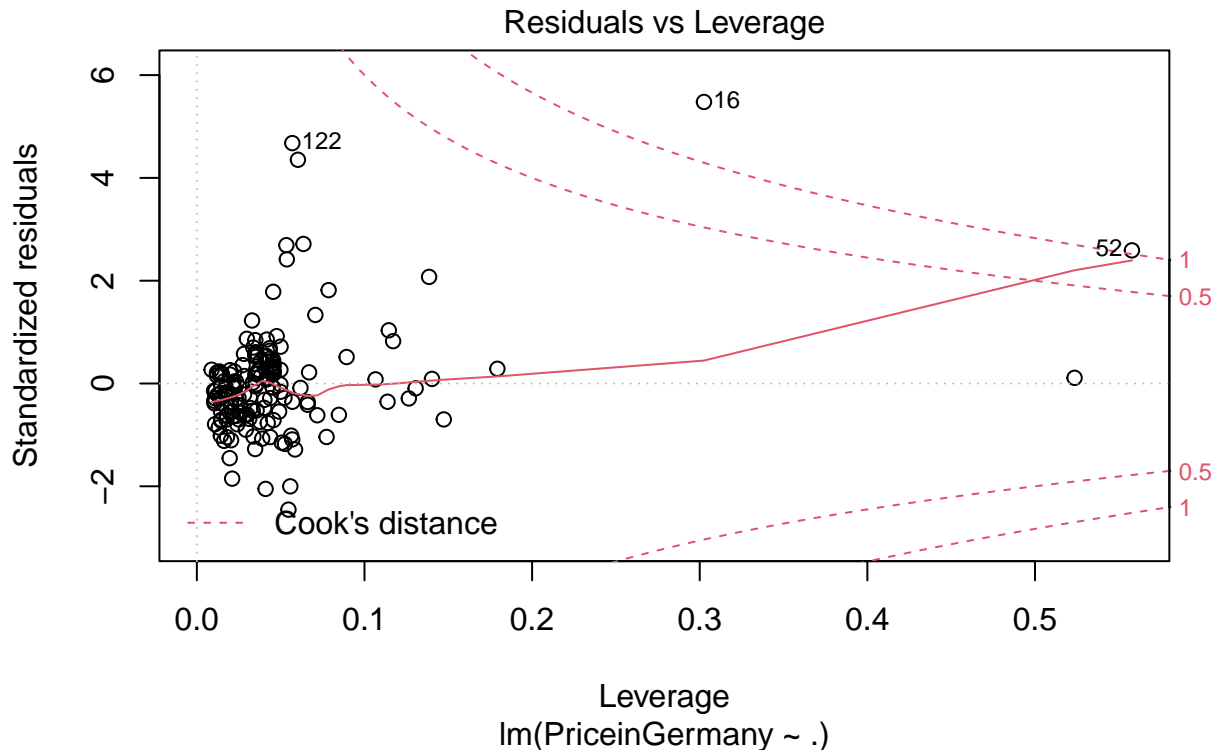
```
## 
## Residual standard error: 17200 on 156 degrees of freedom
## Multiple R-squared:  0.7331, Adjusted R-squared:  0.7211
## F-statistic:  61.2 on 7 and 156 DF,  p-value: < 2.2e-16
```



Our base model is highly significant with a good Adjusted R-squared of 0.7211. Every variable except Acceleration and FastChargeSpeed are significant. However, looking at the residual plots, we can see that there is a cosine pattern in the residuals. The normal qq plot also diverges away from the line.

## Leverage points

From our Price distribution, we can observe that there are already outliers in our dataset. However, we'd like to see how it impacts our model and if those outliers are leverage points. To do this, we utilize the Outliers vs Leverage Plot.

Residuals vs Leverage

lm(PriceinGermany ~ .)

Let's investigate row 16, 52, and 122

```
## # A tibble: 3 x 11
##   Name            Acceleration TopSpeed Range Efficiency FastChargeSpeed Drive
##   <chr>                  <dbl>    <dbl> <dbl>      <dbl>           <dbl> <fct>
## 1 Lightyear One             10      150   575        104             540 All Wh~
## 2 Tesla Cybertru~            3      210   750        267             710 All Wh~
## 3 Porsche Taycan~          2.9      250   380        220             790 All Wh~
## # ... with 4 more variables: NumberofSeats <dbl>, PriceinGermany <dbl>,
## #   PriceinUK <dbl>, BatterySize <dbl>
```

All Three Vehicles are premium vehicles with a higher price point and better performance than other vehicles. We do not want to remove data points unless necessary to avoid overfitting. However one vehicle stands out more than others which is 'Lightyear One'. It's price point is very high at 149000 EUD, however the stats do not follow the trend in our models. This is because Lightyear One is a concept Solar Vehicle. Thus, we will consider the vehicle to be a Leverage Point and remove the row.

## Checking Multicollinearity

One of the assumptions of a linear regression model is that there is little to no Multicollinearity among the variables. This means Observations are independent of each other. Here we will fit all the variables into a linear regression model and analyze the results.

From looking at the summary of the model, we can observe that 'FastChargeSpeed' coefficients have the wrong sign. From our plots, we see that Price increases as 'FastChargeSpeed' increases. From the pairs

plots, we can also observe that there is an obvious linear relationship between TopSpeed, Acceleration and Fast Charge Speed.

Thus, we will use a measure called VIF (Variance Inflation Factor) to further investigate

```
##     Acceleration         TopSpeed            Range      Efficiency FastChargeSpeed
##         4.675994         4.687021        41.221009       12.076341        3.369230
##    NumberofSeats      BatterySize
##         2.438362        42.062331
```
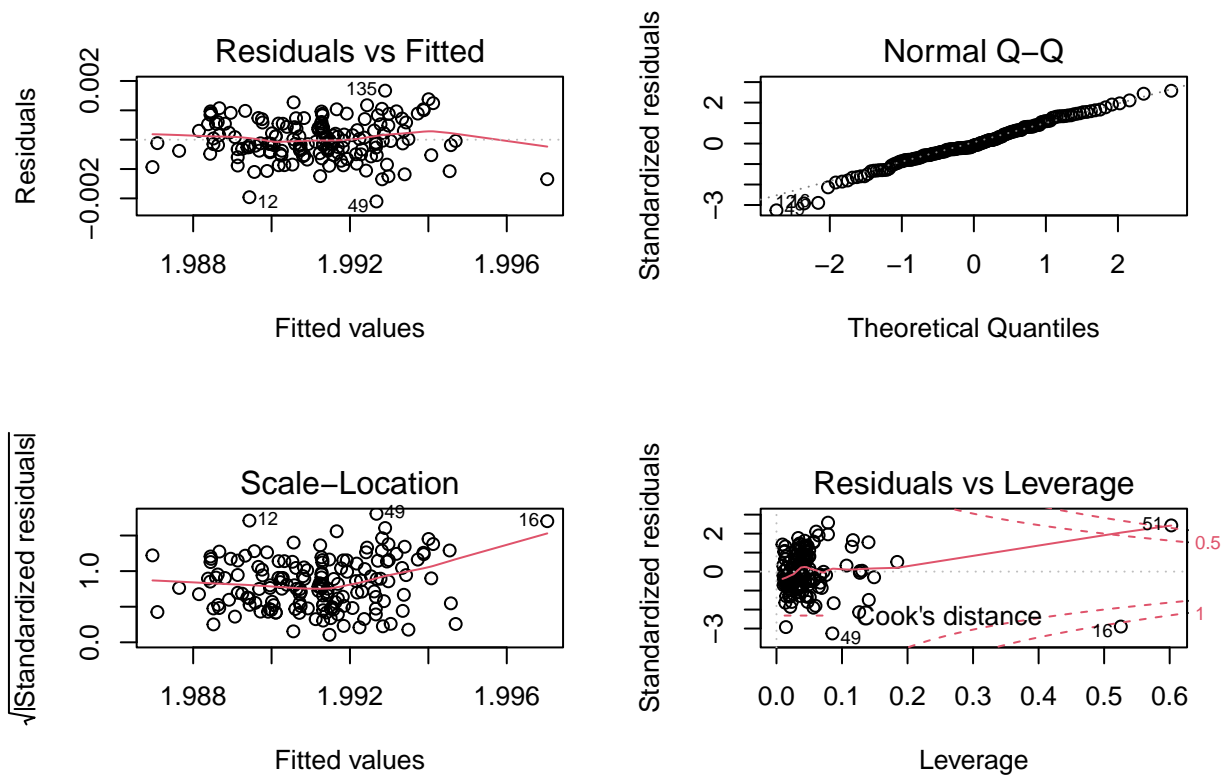
As a general rule of thumb, a VIF of 10 or greater is a cause for concern. We have three variables that have very concerning VIF scores.

## Transformation

It is imperative that we transform are variables so that we could produce a model that fits the assumptions of confidence intervals, normality and constant variance. We'll go ahead and use Power Transform to analyze the recommended transformations.

```
## bcPower Transformations to Multinormality
##                 Est Power Rounded Pwr Wald Lwr Bnd Wald Upr Bnd
## PriceinGermany    -0.7166        -0.5      -0.9467      -0.4864
## BatterySize        0.0024         0.0      -0.0035       0.0082
## Acceleration       0.3696         0.5       0.1239       0.6152
## TopSpeed          -0.9921        -1.0      -1.4089      -0.5752
## Range              0.0037         0.0      -0.0029       0.0103
## Efficiency         0.0087         0.0      -0.0121       0.0294
## FastChargeSpeed   -0.0310         0.0      -0.2356       0.1737
## NumberofSeats     -0.8567        -1.0      -1.6196      -0.0938
##
## Likelihood ratio test that transformation parameters are equal to 0
##   (all log transformations)
##                                   LRT df       pval
## LR test, lambda = (0 0 0 0 0 0 0 0) 65.46494  8 3.9051e-11
##
## Likelihood ratio test that no transformations are needed
##                                   LRT df       pval
## LR test, lambda = (1 1 1 1 1 1 1 1) 1675.886  8 < 2.22e-16
```

Strong evidence that $\lambda = 0$ for variable Battery Size, Range, Efficiency, and Fast Charged Speed so take the $log()$ of those variable to get closer to multivariate normality.

```
##      BatterySize     Acceleration          TopSpeed           Range      Efficiency
##       50.966857         4.660370          4.743997       51.455821       12.811970
## FastChargeSpeed   NumberofSeats
##        3.465666         2.460464
```
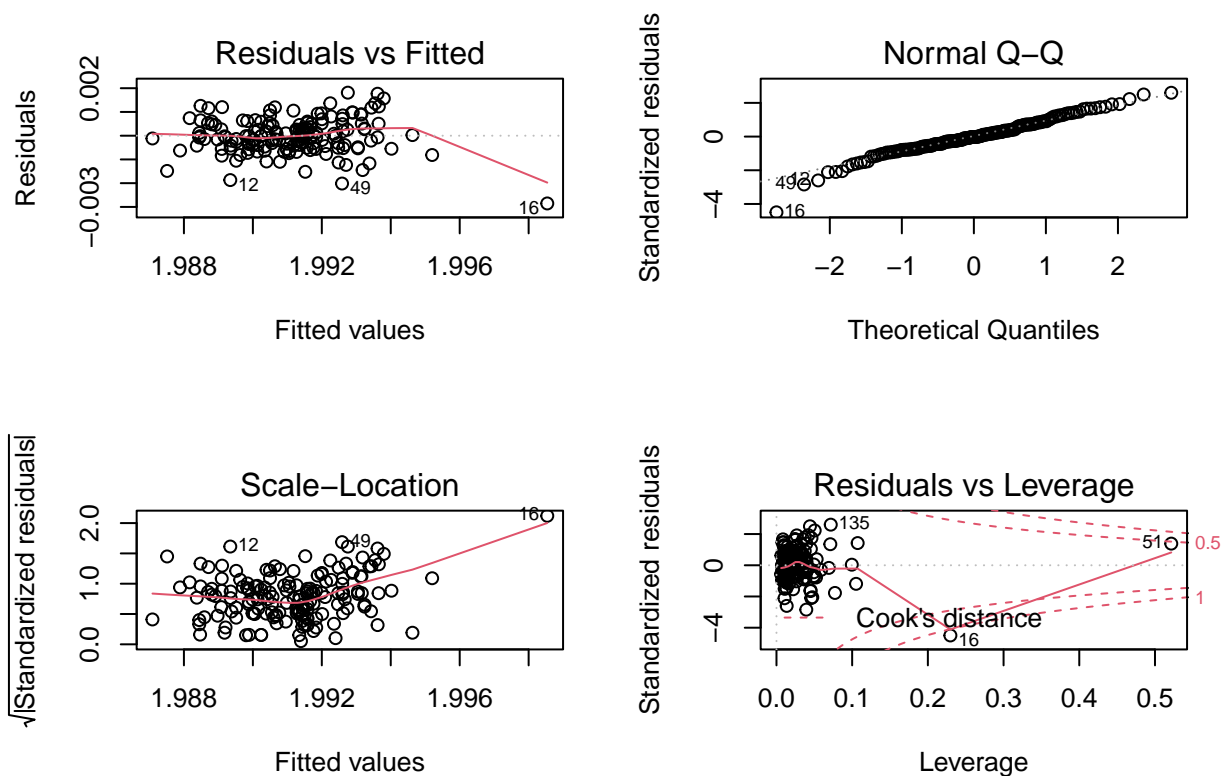
## Variable Selection

We've talked about our dataset having multi-colinearity before. To tackle that problem we will try a few
Variable Selection methods. We will be using sub-setting and analyzing the recommendations. While using
the regsubsets() function, we will particularly focus on Adjusted R-squared, CP values and BIC, as we cannot
use $R^2$ or RSS since higher variables will win everytime.

```
##             adjr2        bic         cp          rss BatterySize Acceleration
## 1  ( 1 ) 0.5012596 -104.21593 408.667619 2.588700e-04           *
## 1  ( 2 ) 0.4873063  -99.71829 424.549299 2.661125e-04
## 2  ( 1 ) 0.7727232 -228.24221 100.080330 1.172348e-04
## 2  ( 2 ) 0.6883048 -176.75755 195.568685 1.607798e-04
## 3  ( 1 ) 0.8041480 -248.42636  65.150076 1.003937e-04
## 3  ( 2 ) 0.7825658 -231.38683  89.409841 1.114567e-04
## 4  ( 1 ) 0.8377342 -275.02512  28.249925 8.265431e-05           *
## 4  ( 2 ) 0.8106304 -249.84737  58.524662 9.646032e-05           *            *
## 5  ( 1 ) 0.8537399 -287.89373  11.337637 7.402986e-05           *            *
## 5  ( 2 ) 0.8460940 -279.58803  19.823970 7.789983e-05           *
## 6  ( 1 ) 0.8571220 -287.65499   8.573665 7.185737e-05           *            *
## 6  ( 2 ) 0.8554328 -285.73924  10.436563 7.270690e-05           *            *
```

```
## 7  ( 1 ) 0.8585489 -285.24551   8.000000 7.068372e-05            *            *
##           TopSpeed Range Efficiency FastChargeSpeed NumberofSeats
## 1  ( 1 )
## 1  ( 2 )         *
## 2  ( 1 )         *            *
## 2  ( 2 )                      *               *
## 3  ( 1 )         *            *               *
## 3  ( 2 )         *     *      *
## 4  ( 1 )         *     *      *
## 4  ( 2 )               *      *
## 5  ( 1 )         *     *      *
## 5  ( 2 )         *     *      *                             *
## 6  ( 1 )         *     *      *                             *
## 6  ( 2 )         *     *      *               *
## 7  ( 1 )         *     *      *               *             *
```

Based on our results the 4(1) row looks the best. Note: The higher R^2 the better. The lower CP, the better. The lower BIC the better. We'll fit and test out the data accordingly

```
##
## Call:
## lm(formula = bcPower(PriceinGermany, p1$roundlam) ~ BatterySize +
##     TopSpeed + Range + Efficiency, data = ev_numeric)
##
## Residuals:
##        Min         1Q     Median         3Q        Max
## -2.854e-03 -4.171e-04 -1.704e-05  4.658e-04  1.812e-03
##
## Coefficients:
##               Estimate Std. Error  t value Pr(>|t|)
## (Intercept)  1.973e+00  9.949e-04 1983.223  < 2e-16 ***
## BatterySize -1.056e-04  1.423e-05   -7.420 6.77e-12 ***
## TopSpeed     2.594e-05  2.086e-06   12.434  < 2e-16 ***
## Range        2.362e-05  2.940e-06    8.033 2.06e-13 ***
## Efficiency   6.307e-05  4.811e-06   13.108  < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.0007233 on 158 degrees of freedom
## Multiple R-squared:  0.8417, Adjusted R-squared:  0.8377
## F-statistic: 210.1 on 4 and 158 DF,  p-value: < 2.2e-16
```
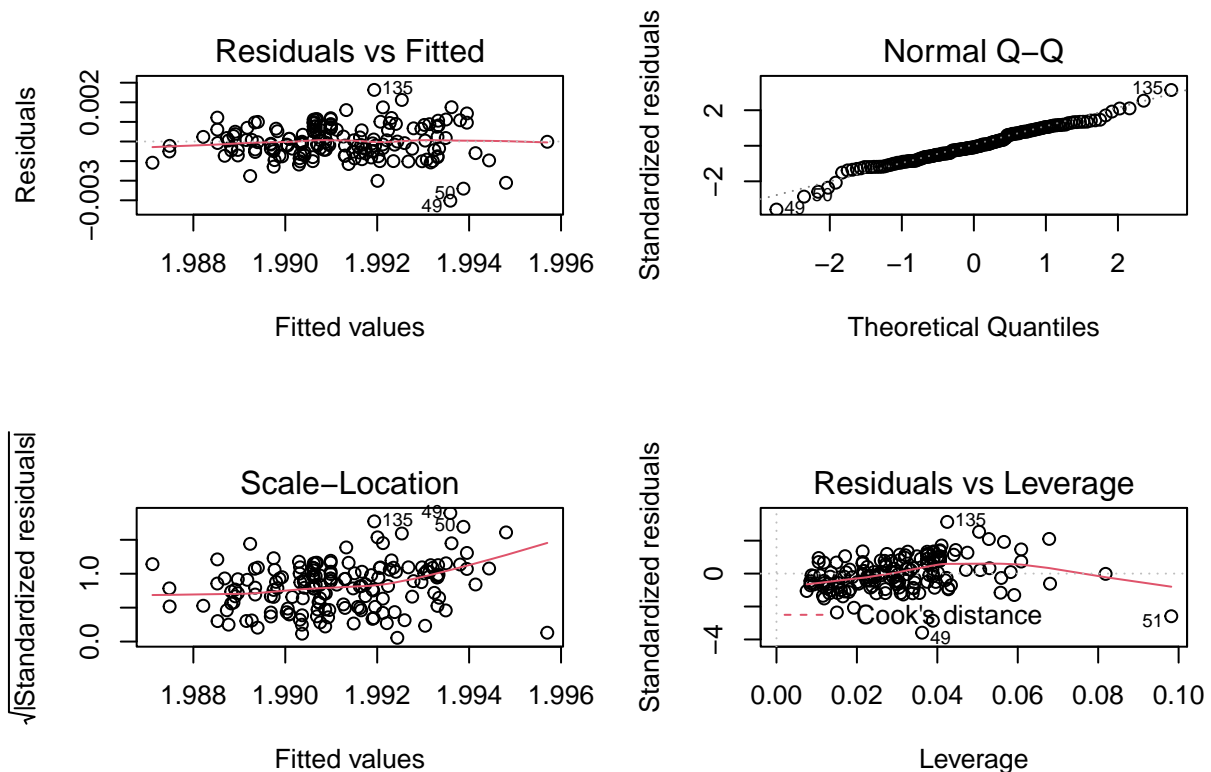
```
## [1] -1888.045
```

```
## BatterySize     TopSpeed       Range   Efficiency
##    37.53803      2.46139     35.51288      7.57899
```

Let's do transformations again.

```
## bcPower Transformations to Multinormality
##                Est Power Rounded Pwr Wald Lwr Bnd Wald Upr Bnd
## PriceinGermany   -0.6800        -0.5      -0.9046      -0.4554
## BatterySize       0.0030         0.0      -0.0016       0.0076
## TopSpeed         -0.8876        -1.0      -1.4046      -0.3706
## Range             0.0039         0.0      -0.0017       0.0095
## Efficiency        0.0101         0.0      -0.0085       0.0288
##
## Likelihood ratio test that transformation parameters are equal to 0
##  (all log transformations)
##                                  LRT df       pval
## LR test, lambda = (0 0 0 0 0) 44.35535   5 1.9617e-08
##
## Likelihood ratio test that no transformations are needed
##                                  LRT df       pval
## LR test, lambda = (1 1 1 1 1) 1502.21   5 < 2.22e-16
```

```
##
```

12

```
## Call:
## lm(formula = bcPower(PriceinGermany, p1$roundlam) ~ log(BatterySize) +
##     I(TopSpeed^-1) + log(Range) + log(Efficiency), data = ev_numeric)
##
## Residuals:
##       Min        1Q    Median        3Q       Max
## -3.020e-03 -5.339e-04 -7.643e-05  6.511e-04  2.632e-03
##
## Coefficients:
##                  Estimate Std. Error t value Pr(>|t|)
## (Intercept)       2.76124    0.32488   8.499 1.35e-14 ***
## log(BatterySize)  0.11634    0.04705   2.473   0.0145 *
## I(TopSpeed^-1)   -0.90567    0.09379  -9.656  < 2e-16 ***
## log(Range)       -0.11550    0.04703  -2.456   0.0151 *
## log(Efficiency)  -0.10988    0.04706  -2.335   0.0208 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.0008571 on 158 degrees of freedom
## Multiple R-squared:  0.7778, Adjusted R-squared:  0.7721
## F-statistic: 138.2 on 4 and 158 DF,  p-value: < 2.2e-16
```



```
##             adjr2       bic         cp          rss log.BatterySize.
## 1  ( 1 ) 0.5808465 -132.55323 141.358975 0.0002175607                *
## 1  ( 2 ) 0.4294332  -82.28517 249.859422 0.0002961514
```

```
## 2  ( 1 ) 0.7586573 -218.45419  14.868318 0.0001244903
## 2  ( 2 ) 0.7144000 -191.00921  46.385403 0.0001473193
## 3  ( 1 ) 0.7696946 -222.01270   7.983222 0.0001180545
## 3  ( 2 ) 0.7657600 -219.25150  10.767654 0.0001200714                    *
## 4  ( 1 ) 0.7721409 -219.68803   7.237815 0.0001160659                    *
## 4  ( 2 ) 0.7706551 -218.62856   8.282722 0.0001168228                    *
## 5  ( 1 ) 0.7755692 -218.10023   5.828076 0.0001135961                    *
## 5  ( 2 ) 0.7737800 -216.80596   7.078291 0.0001145017                    *
## 6  ( 1 ) 0.7767596 -214.91487   6.002646 0.0001122739                    *
## 6  ( 2 ) 0.7742251 -213.07478   7.762377 0.0001135485                    *
## 7  ( 1 ) 0.7753231 -209.82391   8.000000 0.0001122720                    *
##           I.Acceleration.0.5. I.TopSpeed..1. log.Range. log.Efficiency.
## 1  ( 1 )
## 1  ( 2 )                                         *
## 2  ( 1 )                                         *                       *
## 2  ( 2 )                        *                                        *
## 3  ( 1 )                                         *                       *
## 3  ( 2 )                                         *          *
## 4  ( 1 )                                         *          *            *
## 4  ( 2 )                                         *          *
## 5  ( 1 )                                         *          *            *
## 5  ( 2 )                        *                *          *            *
## 6  ( 1 )                        *                *          *            *
## 6  ( 2 )                                         *          *            *
## 7  ( 1 )                        *                *          *            *
##           log.FastChargeSpeed. I.NumberofSeats..1.
## 1  ( 1 )
## 1  ( 2 )
## 2  ( 1 )
## 2  ( 2 )
## 3  ( 1 )                   *
## 3  ( 2 )
## 4  ( 1 )
## 4  ( 2 )              *
## 5  ( 1 )              *
## 5  ( 2 )
## 6  ( 1 )              *
## 6  ( 2 )              *                    *
## 7  ( 1 )              *                    *


##
## Call:
## lm(formula = bcPower(PriceinGermany, p1$roundlam) ~ I(TopSpeed^-1) +
##     log(Efficiency) + log(FastChargeSpeed), data = ev_numeric)
##
## Residuals:
##       Min         1Q      Median         3Q        Max
## -3.035e-03 -5.345e-04 -5.264e-05  6.479e-04  2.591e-03
##
## Coefficients:
##                       Estimate Std. Error t value Pr(>|t|)
## (Intercept)          1.9580213  0.0029012 674.893  < 2e-16 ***
## I(TopSpeed^-1)      -0.8450610  0.1028330  -8.218 6.85e-14 ***
## log(Efficiency)      0.0063918  0.0004270  14.970  < 2e-16 ***
```

14

```
## log(FastChargeSpeed)  0.0007321  0.0002487   2.944  0.00372 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.0008617 on 159 degrees of freedom
## Multiple R-squared:  0.774,  Adjusted R-squared:  0.7697
## F-statistic: 181.5 on 3 and 159 DF,  p-value: < 2.2e-16
```